



Proceedings of the Ninth Annual
U.S. Army Conference
On Applied Statistics,
29-31 October 2003

Yasmin H. Said, Barry A. Bodd
EDITORS

Hosted by:
UNIVERSITY OF CALIFORNIA, DAVIS

Cosponsored by:
U.S. ARMY RESEARCH LABORATORY
U.S. ARMY RESEARCH OFFICE
TRADOC ANALYSIS CENTER – WHITE SANDS MISSILE RANGE
WALTER REED ARMY INSTITUTE OF RESEARCH
UNIFORMED SERVICES UNIVERSITY OF HEALTH SCIENCES

Cooperating Institutions:
LOS ALAMOS NATIONAL LABORATORY
GEORGE MASON UNIVERSITY
INSTITUTE FOR DEFENSE ANALYSIS
OFFICE OF NAVAL RESEARCH
INTERFACE FOUNDATION OF NORTH AMERICA, INC.

Army Research Laboratory

Aberdeen Proving Ground, MD 21005-5067

Proceedings of the Ninth Annual U.S. Army Conference On Applied Statistics, 29-31 October 2003

Yasmin H. Said, EDITOR

Department of Computational and Data Sciences, George Mason University

Barry A. Bodt, EDITOR

Computational and Information Sciences Directorate, ARL

Hosted by:

Francisco J. Samaniego

University of California, Davis

NINTH U.S. ARMY CONFERENCE ON APPLIED STATISTICS

Agenda and Table of Contents

Monday, October 27

TUTORIAL

An Introduction to Generalized Linear Mixed Models.....10
Charles McCulloch, University of California, San Francisco

Tuesday, October 28

TUTORIAL CONTINUED

Wednesday, October 29

GENERAL SESSION I

Chair: Francisco J. Samaniego, University of California, Davis

Hierarchical Models for Spatio-Temporally Correlated Public Health Data
Brad Carlin, University of Minnesota (Keynote Address)

Directional Distributions in Shape Analysis
S. Rao Jammalamadaka, University of California, Santa Barbara

SPECIAL SESSION I: MICROARRAYS

Chair: Edward Wegman, George Mason University

DNA Microarrays in Medicine: Expressing Yourself One Pixel At A Time.....174
J. Patrick Vandersluis, HealthRx Corporation

Transformations and Background Estimation in the Statistical Analysis of Microarrays
Karan Kafadar, University of Colorado, Denver

Statistical Methods of Detecting Differential Gene Expression
David R. Bickel, Medical College of Georgia

CONTRIBUTED SESSION I

Chair: Linda Moss, U.S. Army Research Laboratory

System Test Time Based on Lindstrom/Madden Approach for Continuous Data with Weighted Subsystems.....198

Thomas R. Walker, U.S. Army Evaluation Center

Information Integration for Stockpile Surveillance.....213

Alyson G. Wilson, Los Alamos National Laboratory

Application of Fisher's Combined Probability Test to the Validation of the AIM-9X Missile Model.....234

Art Fries, Institute for Defense Analyses

SPECIAL SESSION II: HOMELAND SECURITY

Chair: Art Fries, Institute for Defense Analyses

Research Within the Department of Homeland Security: Key Issues and Priorities.....275

Parney Albright, Department of Homeland Security

Data Confidentiality, Data Integration, Data Mining, Data Quality: Statistical Challenges for Counterterrorism

Alan Karr, National Institute of Statistical Science

Probabilistic Risk Assessment Methods in Evaluating the Efficacy of Security Procedures

Bernard Harris, University of Nebraska, Lincoln

CONTRIBUTED SESSION II

Chair: David Cruess, Uniformed Services University of the Health Sciences

A Taxonomy of Terrorisms.....289

Cheryl A. Loeb, Center for Technology and National Security Policy, National Defense University and James R. Thompson, Rice University

DoD's Role in Homeland Security: Experimental Opportunity and Experimental Results
Paul Deason, TRADOC Analysis Center-White Sands Missile Range.....329

Water Supply Infrastructure Vulnerability Assessment Methodologies.....358
Major John B. Willis, TRADOC Analysis Center (TRAC)
LTC Thomas M. Cioppa, TRADOC Analysis Center (TRAC)

CONTRIBUTED SESSION III

Chair: Joseph Collins, U.S. Army Research Laboratory

Signature-Related Results on System Reliability
Francisco J. Samaniego, University of California, Davis

Stochastic Measures of Fatigue Crack Damage
Bruce J. West, U.S. Army Research Office

Are Super Efficient Estimators Super Powerful?
Jayaram Sethuraman, Florida State University

GENERAL SESSION II

Chair: Carl Russell, CTR Analytics

Predictive Data Mining with Multiple Additive Regression Trees
Jerome Friedman, Stanford University

BANQUET ADDRESS

Quantitative Sensory Evaluation of Food and Wine
Hildegard Heymann, Professor of Sensory Science, Department of Viticulture and Enology, University of California-Davis

Thursday, October 30

GENERAL SESSION III

Chair: Douglas Tang, Uniformed Services University of the Health Sciences

Calibrated Probabilistic Weather Forecasting via Bayesian Model Averaging
Adrian Raftery, University of Washington

SPECIAL SESSION III: BEYOND SOFTWARE RELIABILITY

Organizers: Tom Zeberlein, Army Test and Evaluation Command and Alyson G. Wilson, Los Alamos National Laboratory

Current and Future Challenges for Software Reliability Assessment.....447
William Farr, Naval Surface Warfare Center

Useful Software Reliability Modeling Practices in Industry Environments.....472
Daniel R. Jeske, University of California at Riverside

CONTRIBUTED SESSION IV

Chair: Donald Gaver, Naval Postgraduate School

Test and Evaluation Challenges Posed by the New DoD 5000.1 Defense Acquisition Directive

Ernest Seglie, Office of the Secretary of Defense, Operational Test & Evaluation
Donald Gaver, Naval Postgraduate School
Patricia Jacobs, Naval Postgraduate School

CONTRIBUTED SESSION V

Chair: Robyn Lee, U.S. Army Center for Health Promotion and Preventive Medicine and U.S. Army Medical Research Institute of Chemical Defense

Significance of Clusters Based on Correlation Distance

Harry L. Hurd, Harry L. Hurd Associates and University of North Carolina, Chapel Hill

Building a Robust Explanatory Model from a Large Data Set

David Kim, United States Military Academy, West Point

Clustering by Local Skewing.....501
David W. Scott, Rice University

CONTRIBUTED SESSION VI

Chair: Jackie Telford, Johns Hopkins Applied Physics Laboratory

Hereditary Portfolio Optimization With Transaction Costs and Taxes: An Infinite-Dimensional HJB Variational Inequality.....513
Mou-Hsiung Chang, U.S. Army Research Office

Cost Estimating Relationship Regression Variance Study
Donald MacKenzie, Wyle Laboratories, Inc.

CONTRIBUTED SESSION VII

Chair: Thomas R. Walker, U.S. Army Evaluation Center

Estimation of Direct-Fire Munition Accuracy Parameters Using the EM Algorithm
David W. Webb, U.S. Army Research Laboratory

Objective Force Urban Operations Agent Based Simulation Experiment.....558
MAJ Lloyd P. Brown, USMC, and LTC Tom Cioppa, USA, U.S. Army TRADOC
Analysis Center-Monterey

*Evaluation of the Technology Readiness Level for the Autonomous Mobility of the
EXperimental Unmanned Vehicle (XUV)*
Barry A. Bodt, U.S. Army Research Laboratory

CONTRIBUTED SESSION VIII

Chair: Robert Burge, Walter Reed Army Institute of Research

Risk Factors for Dental Caries in the Army Population
Robyn B. Lee, M.S. and MAJ Georgia dela Cruz, DMD, MPH
U.S. Army Center for Health Promotion and Preventive Medicine

*Parsimonious Survival Analysis Models: Studying Early Attrition in the Armed Services
by Frailty and Time-Dependent Survival Analysis*.....585
Yuanzhang Li, Timothy Powers, and Margot Krauss, Walter Reed Army Institute of
Research

The Effect of Dosage Errors on the Performance of the Up and Down Method.....591
Douglas R. Sommerville, U.S. Army Chemical and Biological Center

Estimating Proportions with Small Samples
Douglas H. Frank, Indiana University of Pennsylvania

SPECIAL SESSION IV

ROBUST AND RESILIENT CRITICAL INFRASTRUCTURE

Organizers: Jagdish Chandra, George Washington University and Robert Launer, U.S. Army Research Office

Overview, Problem Description, Challenges.....622
Jagdish Chandra, George Washington University

Risk Assessment: A Game Theoretic Approach.....635
Vicki Bier, University of Wisconsin

Distributed Detection with Adversarial and Cooperative Sensors
Nozer Singpurwalla, George Washington University

Optimal Filtering Techniques for Intrusion Detection
Thomas Kurtz, University of Wisconsin

Optimizing Performance in Networked Systems.....658
Steve Robinson, University of Wisconsin

Friday, October 31

GENERAL SESSION IV

Chair: David Kim, United States Military Academy, West Point

Generalized Inference in Mixed Linear Models.....683
Hari Iyer, Colorado State University

CONTRIBUTED SESSION IX

Chair: Alyson G. Wilson, Los Alamos National Laboratory

Test Extremes Using the Random Effect with Variance Adjusting Model: Detecting Non-normal Changes in Military Attrition.....724

Yuanzhang Li, Timothy Powers, and Margot Krauss, Walter Reed Army Institute of Research

Knowledge Discovery in the Military (clinical paper)

Deborah Leishman, Los Alamos National Laboratory

Panel:

Carl Russell, CTR Analytics

David Scott, Rice University

Robert Launer, US Army Research Office

CONTRIBUTED SESSION X

Chair: Paul J. Deason, U.S. Army TRADOC Analysis Center WSMR

Camouflaging Data: An Information Theoretic Perspective

Sallie Keller-McNulty, Los Alamos National Laboratory

Nozer D. Singpurwalla, George Washington University

Multivariate Equivalent Tests with Lognormal Distributions

Karan Kafadar, University of Colorado, Denver

Adaptive Estimation for Inverse Problems with Noisy Operators

Nick Hengartner, Los Alamos National Laboratory

Laurent Cavalier, U. Aix-Marseille II

GENERAL SESSION V

Chair: Barry Bodt, U.S. Army Research Laboratory

Bayesian Methods for Assessing System Reliability

Mike Hamada, Los Alamos National Laboratory

An Introduction to Generalized Linear Mixed Models

by

Charles E. McCulloch

Division of Biostatistics

Department of Epidemiology and Biostatistics

University of California, San Francisco

Army Conference on Applied Statistics, Napa 2003

© 2003 Charles Elliott McCulloch

Outline

- 1) Introduction
 - a) Example: Back pain
 - b) Overview of workshop
 - c) Hierarchical modeling
- 2) Review: Linear Mixed Models (LMMs)
 - a) Example: Fecal fat
 - b) Example: Propranolol
 - i) Analysis
 - ii) Correlations
 - iii) Shrinkage estimators
 - c) Fixed versus random factors
 - d) Best Linear Unbiased Prediction
 - e) Estimation, tests and software in LMMs
- 3) Review: Generalized Linear Models (GLMs)
 - a) Example: Potato flour dilutions
 - b) GLM Basics
 - c) Example: snap beans
 - d) Transforming versus linking
 - e) Estimation, tests and software for GLMs
- 4) Introduction to GLMMs
 - a) Example: skin cancer
- 5) Modeling in GLMMs

- a) Progabide and seizures
 - b) Cartoons and learning disabilities
 - c) Photosynthesis in corn
 - d) Chestnut leaf blight
 - e) Combat vehicle design
 - f) Troponin and cardiac damage
- 6) Features of GLMMs
- a) Consequences of model assumptions
 - b) Marginal versus conditional models
- 7) Inference for GLMMs
- a) Maximum likelihood
 - b) Conditional inference
 - c) Generalized estimating equations
 - d) Penalized quasi-likelihood
 - e) Best Prediction for GLMMs
 - f) Software
- 8) Case Studies
- a) Breeding Bird Survey
 - b) Progabide and seizures
 - c) Potomac River Fever in horses
 - d) Chestnut Leaf Blight
- 9) Summary/Discussion

Approximate Schedule

Morning: 8:30-12:00, Lunch 12:00-13:15, Afternoon: 13:15-14:00

Monday

Morning: Introduction, review of linear mixed models (LMMs) and generalized linear models (GLMs)

Afternoon: Introduction to generalized linear mixed models (GLMMs); GLMM modeling.

Tuesday

Morning: GLMM modeling (cont), features of GLMMs, inference for GLMMs.

Afternoon: Case studies and Summary/Discussion

References

Abramowitz, M. and Stegun, I.A. (1964). Handbook of Mathematical Functions. National Bureau of Standards, Washington, D.C.

Abu-Libdeh, H., Turnbull, B. and Clark, L.C. (1990) Analysis of multi-type recurrent events in longitudinal studies: Application to a skin cancer prevention trial. *Biometrics* **46**: 1017-1034.

Aitken, M. (1999). A general maximum likelihood analysis of variance components. *Biometrics* **55**: 117-128

Atwill, E.R., H.O. Mohammed, J.W. Lopez, C.E. McCulloch, and E.J. Dubovi. Cross-sectional evaluation of environmental, host, and management factors associated with the risk of seropositivity to *Ehrlichia risticii* in horses of New York State. *American Journal of Veterinary Research*, *57*: 278-285, 1996.

Borsch-Supan, A. and Hajivassiliou, V. (1993). Smooth unbiased multivariate probability simulators for maximum likelihood estimation of limited dependent variables. *Journal of Econometrics* **58**: 347-368.

Booth, J., and Hobert, J. (1998). Standard errors of prediction in generalized linear mixed models. *Journal of the American Statistical Association* **93**: 262-272.

Booth, J.G. and Hobert, J.H. (1999), "Maximizing generalized linear mixed model likelihoods with an automated Monte Carlo EM algorithm", *Journal of the Royal Statistical Society B* *62*: 265-285.

Breslow, N.E. and Clayton, D.G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association* **88**: 9-25.

Breslow, N.E. and Lin, X. (1994). Bias correction in generalized linear mixed models with a single component of dispersion. *Biometrika* **82**:81-91.

Carey, V., Zeger, S.L. and Diggle, P. (1993). Modelling multivariate binary data with alternating logistic regressions. *Biometrika* **80**: 517-526.

Casella, G.C. and Berger, R.L. (1994). Estimation with selected binomial information, or Do you really believe Dave Winfield is batting .471? *Journal of the American Statistical Association* **89**: 1080-1090.

Conaway, M.R. (1989). Analysis of repeated categorical measurements with conditional likelihood methods. *Journal of the American Statistical Association* **89**: 53-62.

Conaway, M.R. (1990). A random effects model for binary data. *Biometrics* **46**: 317-328.

- Cortesi, P., and Milgroom, M. (1998). Genetics of vegetative incompatibility in *Cryphonectria parasitica* *Appl. Environ Microbiol.* **64**: 2988-2994.
- Cortesi, P., Milgroom, M., and Bisiach, M. (1996). Distribution and diversity of vegetative incompatibility types in subpopulations of *Cryphonectria parasitica* in Italy. *Mycological Research*, **100**: 1087-1093.
- Cortesi, P., McCulloch, C, Song, H., Lin, H. and Milgroom, M. (2001). Genetic control of horizontal virus transmission in the chestnut blight fungus, *Cryphonectria parasitica*. *Genetics*, **159**: 107-118.
- Cox, D.R. and Snell E.J. (1989). *Analysis of Binary Data*, 2nd Edition. Chapman and Hall, London.
- Crowder, M.J. (1978). Beta-binomial ANOVA for proportions. *Applied Statistics* **27**:34-37.
- Devore, J. and Peck, R. (1993). *Statistics*. Duxbury, Belmont, CA.
- Diggle, P., Liang, K.-Y., Zeger, S.L. and Heagerty, P (2002). *Longitudinal Data Analysis*, 2nd Ed. Oxford University Press, Oxford.
- Drum, M. and McCullagh, P. (1993). REML estimation with exact covariance in the logistic mixed model. *Biometrics* **49**: 677-689.
- Engel, B. and Keen, A. (1994). A simple approach for the analysis of generalized linear mixed models. *Statistica Neerlandica* **48**:1-22.
- Fitzmaurice, G.M. (1995). A caveat concerning independence estimating equations with multivariate binary data. *Biometrics* **51**: 309-317.
- Geyer, C.J. (1994). Estimating normalizing constants and reweighting mixtures in Markov chain Monte Carlo. Technical Report No. 568, School of Statistics, University of Minnesota.
- Geyer, C.J. and Thompson, E.A. (1992). Constrained Monte Carlo maximum likelihood for dependent data. *Journal of the Royal Statistical Society, Series B*, **54**: 657-699.
- Gilks, W.R., Wang, C.C., Yvonnet, B., and Coursaget, P. (1993). Random-effects models for longitudinal data using Gibbs sampling. *Biometrics* **49**: 441-453.
- Gilmour, A.R., Anderson, R.D. and Rae, A.L. (1984). The analysis of binary data by a generalized linear mixed model. *Biometrika* **72**: 593-599.
- Hamet, P., Kuchel, O., Cuhe, J.L., Boucher, R., and Genest, J. (1973). Effect of propranolol on cyclic AMP excretion and plasma renin activity in labile essential hypertension. *Canadian Medical Association Journal* **1**: 1099-1103.

- Heagerty, P. and Lele, S. (1998). A composite likelihood approach to binary spatial data. *Journal of the American Statistical Association* **93**: 1099-1111.
- Heagerty, P. (1999). Marginally specified logistic-normal models for longitudinal binary data. *Biometrics* **55**: 688-698.
- Heagerty, P. and Kurland, B. (2001). Misspecified maximum likelihood estimates and generalized linear mixed models. *Biometrika* **88**: 973-986.
- Hedeker, D. and Gibbons, R.D. (1994). A random-effects ordinal regression model for multilevel analysis. *Biometrics* **50**: 933-944.
- James, F.C., McCulloch, C.E., and Wiedenfeld, D.A. (1996). New approaches to the analysis of population trends in land birds. *Ecology* **77**: 13-27.
- Karim, M.R., Zeger, S.L. (1992). Generalized linear models with random effects; Salamander mating revisited. *Biometrics* **48**: 631-644.
- Kauermann, G. and Carroll, R.J. (2001). A note on the efficiency of sandwich covariance matrix estimation. *J. Amer. Stat. Association* **79**: 1387-1396.
- Korff, M.V., Barlow, W., Cherkin, D., and Deyo, R.A. (1994). Effects of practice style in managing back pain. *Ann. Internal Med.* **121**: 187-95.
- Kuk, A.Y.C. (1995). Asymptotically unbiased estimation in generalized linear models with random effects. *Journal of the Royal Statistical Society, Series B* **57**:395-407.
- Lee, Y., and Nelder, J.A. (1996). Hierarchical generalized linear models (Disc: p656-678). *Journal of the Royal Statistical Society, Series B*, **58**: 619-656.
- Liang, K.-Y. and Zeger, S.L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika* **73**:13-22
- Liang, K.-Y., Zeger, S.L., and Qaqish, B.H. (1992). Multivariate regression analyses for categorical data. *Journal of the Royal Statistical Society, Series B*, **54**: 673-687.
- Lin, X. and Breslow, N.E. (1996). Bias correction in generalized linear mixed models with multiple components of dispersion. *Journal of the American Statistical Association* **91**: 1007-1016.
- Lindsey, J.K. and Lambert, P. (1998). On the appropriateness of marginal models for repeated measures in clinical trials. *Statistics in Medicine* **17**: 447-469.
- Lipsitz, S.R., Fitzmaurice, G.M., Orav, E.J., and Laird, N.M. (1994). Performance of generalized estimating equations in practical situations. *Biometrics* **50**: 270-278.

McCullagh, P. and Nelder, J. (1989). *Generalized Linear Models*, 2nd Ed. Chapman and Hall, London.

McCulloch, C.E. (1994). Maximum likelihood estimation of variance components for binary data. *Journal of the American Statistical Association* **89**: 330-335.

McCulloch, C.E. (1997). Maximum likelihood algorithms for generalized linear mixed models. *Journal of the American Statistical Association* **92**: 162-170.

McCulloch, C.E. and Searle, S.R. (2000). *Generalized, Linear, and Mixed Models*. Wiley, New York.

McGilchrist, C. A. (1993). REML estimation for survival models with frailty. *Biometrics* **49**: 221-225.

McGilchrist, C.A. (1994). Estimation in generalized mixed models. *Journal of the Royal Statistical Society, Series B* **56**: 61-69.

McGilchrist, C.A. and Yau, K. K. W. (1995). The derivation of BLUP, ML, REML estimation methods for generalised linear mixed models. *Communications in Statistics. Theory and Methods* **24**: 2963-2980.

McLean, R.A., Sanders, W.L., and Stroup, W.W. (1991). A unified approach to mixed linear models. *American Statistician* **45**: 54-64.

Natarajan, R. and McCulloch, C.E. (1995). A note on existence of the posterior distribution for a class of mixed models for binomial responses. *Biometrika* **82**:639-643.

Neuhaus, J. M., Lesperance, M. L. (1996). Estimation efficiency in a binary mixed-effects model setting. *Biometrika* **83**: 441-446.

Prentice, R.L. (1988). Correlated binary regression with covariates specific to each binary observation. *Biometrics* **44**: 1033-1048.

Robinson, G.K. (1991). That BLUP is a good thing - the estimation of random effects. *Statistical Sciences* **6**:15-51.

Ruppert, D., Cressie, N., and Carroll, R. (1989). A transformation/weighting model for estimating Michaelis-Menten parameters. *Biometrics* **45**: 637-656.

Ruppert, D., Reish, R.L., Deriso, R.B., and Carroll, R.J. (1984). Optimization using stochastic approximation and Monte Carlo simulation (with application to harvesting of Atlantic Menhaden). *Biometrics* **40**: 535-545.

Schall, R. (1991). Estimation in generalized linear models with random effects. *Biometrika* **78**:719-727.

Self and Liang, K.-Y. (1987). Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association* **82**: 605-610.

Searle, S.R., Casella, G., and McCulloch, C.E. (1992). *Variance Components*. Wiley, New York.

Snedecor, G.W. and Cochran, W.G. (1989). *Statistical Methods*, 8th Edition. Iowa State University Press, Ames, Iowa.

Stram, D.O. and Lee, J.W. (1994). Variance components testing in the longitudinal mixed effects model. *Biometrics* **50**:1171-1177.

Waclawiw, M.A., and Liang, K.-Y. (1993). Prediction of random effects in the generalized linear model. *Journal of the American Statistical Association* **88**: 171-178.

Wolfinger, R.W. (1994). Laplace's approximation for nonlinear mixed models. *Biometrika* **80**: 791-795.

Zeger, S.L. and Liang, K.-Y. (1986). Longitudinal data analysis for discrete and continuous outcomes. *Biometrics* **42**: 121-130.

Zeger, S.L. and Karim, M.R. (1991). Generalized linear models with random effects: a Gibbs sampling approach. *Journal of the American Statistical Association* **86**: 79-86.

1. Introduction

Example: Practice style and back pain (Korff, Barlow, Cherkin, and Deyo, 1994).

Forty-four primary care physicians in a large HMO were classified according to their practice style in treating back pain management (low, moderate or high frequency of prescription of pain medication and bed rest). An average of 24 patients per physician was followed for 2 years (1 month, 1 year and 2 year followups) after the indexed visit. Outcome measures included functional measures (pain intensity, activity limitation days, etc.), patient satisfaction (e.g., “After your visit with the doctor, you fully understood how to take care of your back problem”), and cost.

Q1. Does practice style influence function, satisfaction or cost?

Q2. How much of the variability in the responses is due to physician?

Q3. How well are individual physicians performing with regard to effectiveness and cost?

Model for log(cost) at year 1: Incorporating predictors of practice style of the physician, age of the patient, whether the back pain was cervical or thoracic or neither (yes=1, no=0).

Model for “understand how to care for your back” at year 1: Incorporating predictors of practice style of the physician, age of the patient, whether the back pain was cervical or thoracic or neither (yes=1, no=0).

1. b. Overview

- Hierarchical modelling
- Review of Linear Mixed Models (LMMs) and Generalized Linear Models (GLMs)
- Examples of Generalized Linear Mixed Models (GLMMs)
- Modeling using GLMMs
- Features of GLMMs
- Inference methods
- Case studies
- Discussion and summary

1. c. Hierarchical Modeling

Hierarchical data: Data (responses and/or predictors) collected from different levels within a study. Other terminology for the same or related ideas: repeated measures data, longitudinal data, clustered data, multilevel data.

Example 1: Practice style and back pain (Korff, Barlow, Cherkin, and Deyo, 1994).

Example 2: The Educational Testing Service in the past has offered guidance to both law school admissions officers and to potential applicants to law school via their Law School Validity Studies. One aspect of this has been to create a simple index that allows admission officers to screen applicants and for applicants to gauge the likelihood of acceptance to a law school before applying.

Two of the indicators used for predicting success in law school are the LSAT score (ranging from 200 to 800) and the undergraduate GPA (UGPA). A form of combining the LSAT and UGPA, which has successfully been used in the past to predict first year performance at law school, has been:

$$\text{Predicted performance} = \text{LSAT} + (\text{mult}) \times \text{UGPA}$$

Where mult is a multiplier chosen to reflect the relative importance of LSAT and UGPA and which might be dependent on the school doing the admissions. For example, a multiplier of 200 might make sense since it puts both GPA

(typically in a range of 1.0 to 4.0) and LSAT on the same scale.

In practice the multipliers have been estimated from data taken from admitted students and this is done separately for each law school. The estimation was often done by a multiple regression of first year performance on both LSAT and UGPA.

Q. What is the best way to estimate the multipliers for each school?

Example 3: Lack of digestive enzymes in the intestine can cause bowel absorption problems. This will be indicated by excess fat in the feces. Pancreatic enzyme supplements can be given to ameliorate the problem. Does the supplement form make a difference? (Graham DY, Enzyme replacement therapy of exocrine pancreatic insufficiency in man. *NEJM*, **296**: 1314-17, 1977 – But note: sex information made up for illustration.)

PatID / Sex	Fecal Fat (g/day)				Avg
	Pill type				
	None	Tablet	Capsule	Coated Capsule	
1 – M	44.5	7.3	3.4	12.4	16.900
2 – M	33.0	21.0	23.1	25.4	25.625
3 – M	19.1	5.0	11.8	22.0	14.475
4 – F	9.4	4.6	4.6	5.8	6.100
5 – F	71.3	23.3	25.6	68.2	47.100
6 – F	51.2	38.0	36.0	52.6	44.450
Avg	38.08	16.53	17.42	31.07	25.775

Example 4: Propranolol and hypertension
(Hamet, *et al*, 1975)

Below are data from an early, double-blind trial of the effect of propranolol on labile hypertension. Blood pressure was measured under the drug and a placebo both in the upright and recumbent positions.

Patient	Blood Pressure (mmHg)				Ave.
	Recumbent		Upright		
	Placebo	Propran.	Placebo	Propran	
1	96	71	73	87	81.75
2	96	85	104	76	90.25
3	92	89	83	90	88.50
4	97	110	101	85	98.25
5	104	85	112	94	98.75
6	100	73	101	93	91.75
7	93	81	88	85	86.75
Ave.	96.86	84.86	94.57	87.14	90.86

Q1: Does Propranolol have the same influence in recumbent and upright positions?

Q2: If the answer to Q1 is yes, is it effective?

Analysis Approaches

Basic tenet: Don't go beyond standard and accepted statistical practices unless necessary.

Applied in this context: Do we need hierarchical models?

The usual statistical methods (multiple regression, basic ANOVA, logistic regression, and many others) assume observations are independent.

Important idea: observations taken within the same subgroup in a hierarchy are often more similar to one another than to observations in different subgroups, other things being equal. [correlated]

Also getting the correlation assumptions wrong in a statistical analysis is often a very serious mistake.

Simple Analysis Strategies

What strategies might we employ in analyzing data from a hierarchical format?

1. Separate analyses for each subgroup.
2. Analyses at the lowest level in the hierarchy.
3. Analyses at the highest level in the hierarchy.
4. Derived variables.

Let's consider an example of each of these and advantages and disadvantages.

1. Separate analyses for each subgroup.

Fecal fat example?

The law school example follows this approach by calculating separate multipliers for each law school. Here are the multipliers estimated for selected law schools for three consecutive years and pooling the data across years.

Law School	Separate Years			Pooled Years	
	Year 1	Year 2	Year 3	Years 1-2	Years 2-3
1	2507	301	105	526	164
2	-24	49	153	5	116
3	179	118	98	149	107

Law schools 1 and 2 were selected as being somewhat extreme and 3 was “typical”.

2. Analyses at the lowest level in the hierarchy.

For the back pain example this corresponds to analyzing each observation on the patient and attributing to each one the higher level characteristics, e.g., an observation taken from a “low” doctor.

3. Analyses at the highest level in the hierarchy.

For the back pain example, this corresponds to calculating the average value of the response for each doctor across all patients and time periods.

4. Derived variable approach.

For the fecal fat example we would calculate several new responses: (1) the difference between the none and tablet observations for each patient, (2) the difference between the none and capsule observations for each patient, (3) the difference between the and tablet and coated tablet observations for each patient. These new responses are then subjected to one-sample t-tests.

When to Use Hierarchical Models

The use of hierarchical/mixed models is clearly indicated in three situations:

1. When the correlation structure is of primary interest.
2. When we wish to “borrow strength” across the levels of a hierarchy in order to improve estimates.

(81 law schools and one year of data versus 2 years of data)

3. When dealing with highly unbalanced correlated data.

2. Review: Linear Mixed Models (LMMs)

Analysis of the fecal fat example (Stata)

```
summ
Variable |      Obs      Mean   Std. Dev.   Min      Max
-----+-----
   fecfat |       24   25.775   20.00214    3.4    71.3
     patid |       24     3.5   1.744557     1     6
   pilltype |       24     2.5   1.14208     1     4

. sort pilltype

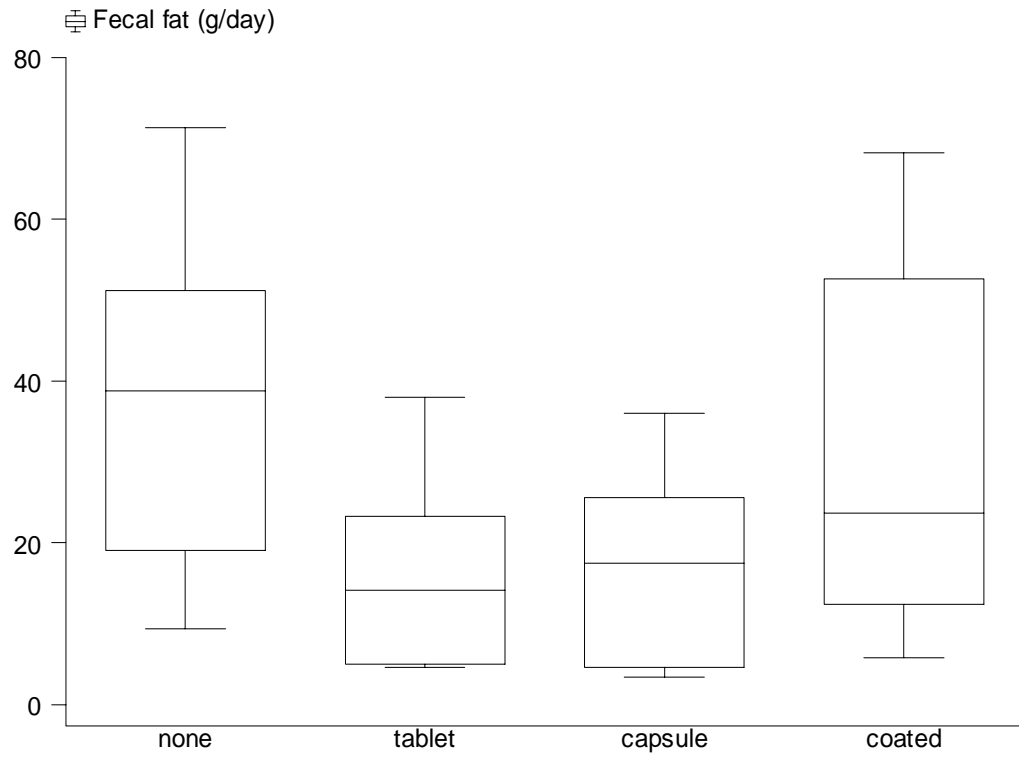
. by pilltype: summarize fecfat

-> pilltype= none
Variable |      Obs      Mean   Std. Dev.   Min      Max
-----+-----
   fecfat |         6  38.08333   22.47447    9.4    71.3

-> pilltype= tablet
Variable |      Obs      Mean   Std. Dev.   Min      Max
-----+-----
   fecfat |         6  16.53333   13.32091    4.6     38

-> pilltype= capsule
Variable |      Obs      Mean   Std. Dev.   Min      Max
-----+-----
   fecfat |         6  17.41667   12.93745    3.4     36

-> pilltype= coated
Variable |      Obs      Mean   Std. Dev.   Min      Max
-----+-----
   fecfat |         6  31.06667   24.2641    5.8    68.2
```



Analyses ignoring sex effects

ANOVA (*wrong analysis*)

```

. xi: regr fecfat i.pilltype
i.pilltype          Ipillt_1-4    (naturally coded; Ipillt_1 omitted)

```

Source	SS	df	MS	Number of obs =	24
Model	2008.6017	3	669.533901	F(3, 20) =	1.86
Residual	7193.36328	20	359.668164	Prob > F =	0.1687
Total	9201.96498	23	400.085434	R-squared =	0.2183
				Adj R-squared =	0.1010
				Root MSE =	18.965

fecfat	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
Ipillt_2	-21.55	10.9494	-1.968	0.063	-44.39005 1.29005
Ipillt_3	-20.66667	10.9494	-1.887	0.074	-43.50672 2.173384
Ipillt_4	-7.016668	10.9494	-0.641	0.529	-29.85672 15.82338
_cons	38.08333	7.742396	4.919	0.000	21.93298 54.23369


```

. testparm Ipill*

( 1)  Ipillt_2 = 0.0
( 2)  Ipillt_3 = 0.0
( 3)  Ipillt_4 = 0.0

      F( 3, 20) =      1.86
      Prob > F =      0.1687

```

Hierarchical analysis

```
. xi: xtgee fecfat i.pilltype, i(patid)
i.pilltype          Ipillt_1-4   (naturally coded; Ipillt_1 omitted)
```

```
Iteration 1: tolerance = 1.108e-15
```

```
GEE population-averaged model
Group variable:          patid
Link:                    identity
Family:                  Gaussian
Correlation:            exchangeable
Scale parameter:        299.7235
Number of obs           =      24
Number of groups        =       6
Obs per group: min     =       4
                      avg     =      4.0
                      max     =       4
Wald chi2(3)           =      22.53
Prob > chi2            =      0.0001
```

fecfat	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
Ipillt_2	-21.55	5.451781	-3.953	0.000	-32.23529	-10.86471
Ipillt_3	-20.66667	5.451781	-3.791	0.000	-31.35196	-9.981373
Ipillt_4	-7.016668	5.451781	-1.287	0.198	-17.70196	3.668626
_cons	38.08333	7.067808	5.388	0.000	24.23068	51.93598

```
. testparm Ipillt*
```

```
( 1)  Ipillt_2 = 0.0
( 2)  Ipillt_3 = 0.0
( 3)  Ipillt_4 = 0.0
```

```
      chi2( 3) =    22.53
      Prob > chi2 =    0.0001
```

Hierarchical analysis (variation)

```
. xi: xtgee fecfat i.pilltype, i(patid) robust
i.pilltype          Ipillt_1-4   (naturally coded; Ipillt_1 omitted)
```

```
Iteration 1: tolerance = 1.662e-15
```

```
GEE population-averaged model
Group variable:          patid
Link:                    identity
Family:                  Gaussian
Correlation:             exchangeable
Scale parameter:        299.7235
Number of obs           =      24
Number of groups        =       6
Obs per group: min     =       4
                      avg     =      4.0
                      max     =       4
Wald chi2(3)           =      11.71
Prob > chi2            =      0.0084
```

(standard errors adjusted for clustering on patid)

fecfat	Semi-robust		z	P> z	[95% Conf. Interval]	
	Coef.	Std. Err.				
Ipillt_2	-21.55	6.931847	-3.109	0.002	-35.13617	-7.96383
Ipillt_3	-20.66667	7.349407	-2.812	0.005	-35.07124	-6.262094
Ipillt_4	-7.016668	5.246295	-1.337	0.181	-17.29922	3.265881
_cons	38.08333	9.175163	4.151	0.000	20.10034	56.06632

```
. testparm Ipillt*
```

```
( 1) Ipillt_2 = 0.0
( 2) Ipillt_3 = 0.0
( 3) Ipillt_4 = 0.0
```

```
      chi2( 3) =    11.71
Prob > chi2 =    0.0084
```

Analyses incorporating sex effects

ANOVA (*wrong analysis*)

```
. xi: regr fecfat i.pilltype i.sex
i.pilltype      _Ipilltype_1-4      (naturally coded; _Ipilltype_1 omitted)
i.sex           _Isex_0-1           (naturally coded; _Isex_0 omitted)
```

Source	SS	df	MS	Number of obs =	24
Model	3110.21668	4	777.554169	F(4, 19) =	2.43
Residual	6091.7483	19	320.618332	Prob > F =	0.0837
Total	9201.96498	23	400.085434	R-squared =	0.3380
				Adj R-squared =	0.1986
				Root MSE =	17.906

fecfat	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
_Ipilltype_2	-21.55	10.33793	-2.08	0.051	-43.18753 .0875334
_Ipilltype_3	-20.66667	10.33793	-2.00	0.060	-42.3042 .970867
_Ipilltype_4	-7.016668	10.33793	-0.68	0.505	-28.6542 14.62087
_Isex_1	13.55	7.31002	1.85	0.079	-1.750047 28.85005
_cons	31.30833	8.172851	3.83	0.001	14.20236 48.41431

Hierarchical analysis

```
. xi: xtgee fecfat i.pilltype i.sex, i(patid)
i.pilltype      _Ipilltype_1-4      (naturally coded; _Ipilltype_1 omitted)
i.sex           _Isex_0-1           (naturally coded; _Isex_0 omitted)
```

Iteration 1: tolerance = 1.219e-15

```
GEE population-averaged model
Group variable:          patid
Link:                    identity
Family:                  Gaussian
Correlation:             exchangeable
Scale parameter:        253.8228
Wald chi2(4)            = 24.00
Prob > chi2              = 0.0001
```

fecfat	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
_Ipilltype_2	-21.55	5.451781	-3.95	0.000	-32.23529 -10.86471
_Ipilltype_3	-20.66667	5.451781	-3.79	0.000	-31.35196 -9.981373
_Ipilltype_4	-7.016668	5.451781	-1.29	0.198	-17.70196 3.668626
_Isex_1	13.55	11.16389	1.21	0.225	-8.330816 35.43082
_cons	31.30833	8.570992	3.65	0.000	14.5095 48.10717

Hierarchical analysis (variation)

```
. xi: xtgee fecfat i.pilltype i.sex, i(patid) robust
i.pilltype      _Ipilltype_1-4      (naturally coded; _Ipilltype_1 omitted)
i.sex           _Isex_0-1           (naturally coded; _Isex_0 omitted)
```

Iteration 1: tolerance = 1.219e-15

```
GEE population-averaged model
Group variable:      patid          Number of obs      =      24
Link:                identity       Number of groups   =      6
Family:              Gaussian       Obs per group: min =      4
Correlation:         exchangeable   avg                =      4.0
Scale parameter:    253.8228        max                =      4
Wald chi2(4)        =      12.80
Prob > chi2         =      0.0123
```

(standard errors adjusted for clustering on patid)

fecfat	Semi-robust					
	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_Ipilltype_2	-21.55	6.931847	-3.11	0.002	-35.13617	-7.96383
_Ipilltype_3	-20.66667	7.349407	-2.81	0.005	-35.07124	-6.262094
_Ipilltype_4	-7.016668	5.246295	-1.34	0.181	-17.29922	3.265881
_Isex_1	13.55	12.22942	1.11	0.268	-10.41923	37.51923
_cons	31.30833	4.918175	6.37	0.000	21.66889	40.94778

Notes

- Hierarchical data structures are common.
- They lead to correlated data.
- Ignoring the correlation can be a serious error.

Fixed versus Random Factors

Definition: If a distribution is assumed for the levels of a factor it is random. If the values are fixed, unknown constants it is a fixed factor.

Ramifications:

- Scope of inference
Inferences can be made on a statistical basis to the *population* from which the levels of the random factor have been selected.
- Incorporation of correlation in the model
Observations that share the same level of the random effect are being modeled as correlated.
- Accuracy of estimates
Using random factors involves making extra assumptions but gives more accurate estimates.
- Estimation method
Different estimation methods must be used.

How to decide in practice?

SAS Proc MIXED philosophy:
fixed factors → MODEL statement
random factors → RANDOM statement
additional temporal and spatial correlation
→ REPEATED statement

SAS program for the Propranolol Example

```
data propran;
input bp patient upright drug;
cards;
96      1      0      0
71      1      0      1
73      1      1      0
87      1      1      1
96      2      0      0
85      2      0      1
104     2      1      0
76      2      1      1
      .
      .
      .
92      3      0      0
93      6      1      1
93      7      0      0
81      7      0      1
88      7      1      0
85      7      1      1

proc mixed;
  class patient upright drug;
  model bp=upright drug upright*drug;
  estimate "blup pat 1" | patient 1 ;
  estimate "blup pat 2" | patient 0 1;
  estimate "blup pat 3" | patient 0 0 1 ;
  estimate "blup pat 4" | patient 0 0 0 1;
  random patient;
run;
```

SAS Output for the Propranolol Data

The SAS System The MIXED Procedure

Class Level Information

Class	Levels	Values
PATIENT	7	1 2 3 4 5 6 7
UPRIGHT	2	0 1
DRUG	2	0 1

REML Estimation Iteration History

Iteration	Evaluations	Objective	Criterion
0	1	142.68756055	
1	1	141.94268164	0.00000000

Convergence criteria met.

Covariance Parameter Estimates (REML)

Cov Parm	Estimate
PATIENT	15.79761905
Residual	85.79761905

Model Fitting Information for BP

Description	Value
Observations	28.0000
Res Log Likelihood	-93.0259
Akaike's Information Criterion	-95.0259
Schwarz's Bayesian Criterion	-96.2039
-2 Res Log Likelihood	186.0517

Tests of Fixed Effects

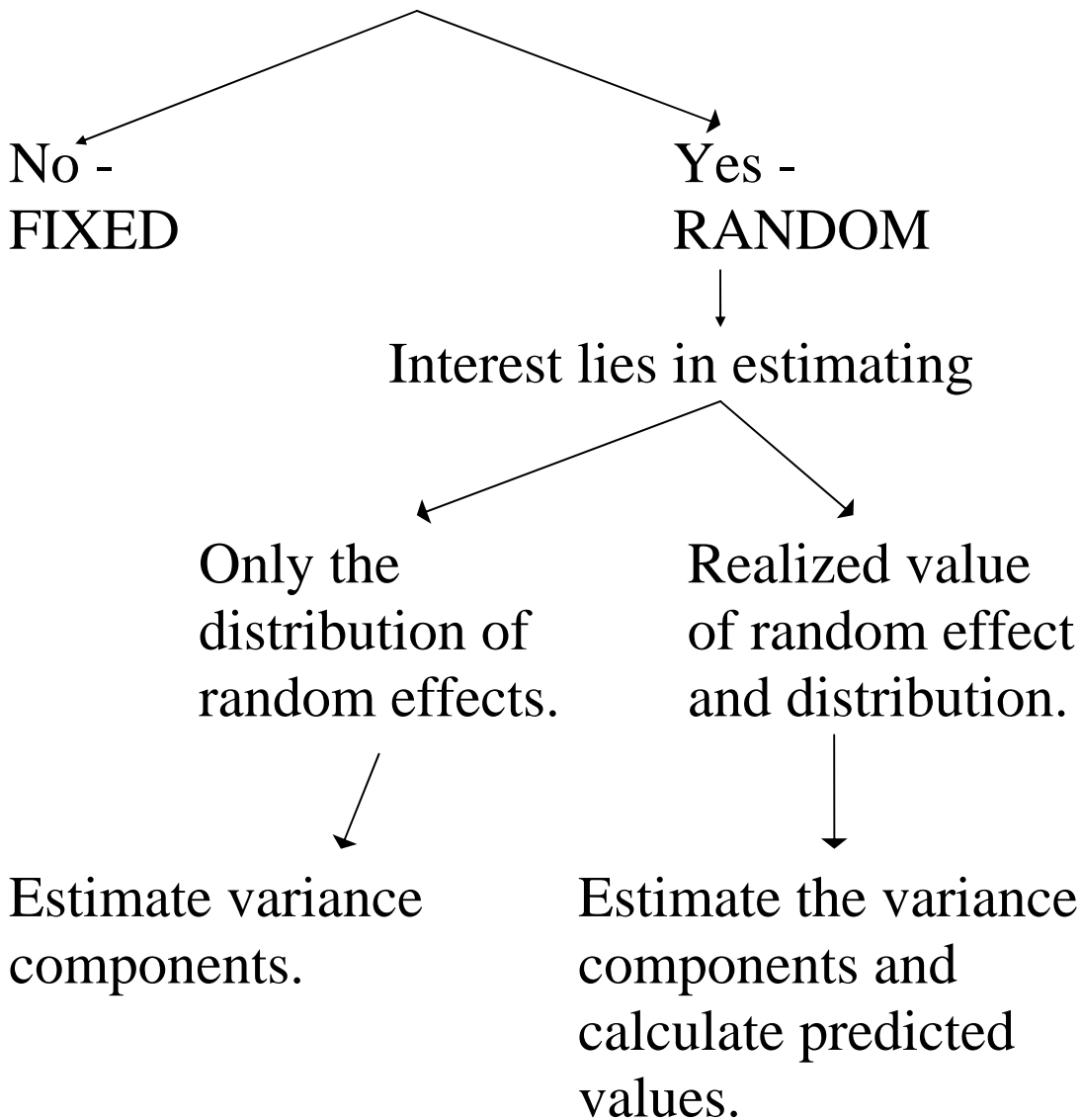
Source	NDF	DDF	Type III F	Pr > F
UPRIGHT	1	18	0.00	1.0000
DRUG	1	18	7.70	0.0125
UPRIGHT*DRUG	1	18	0.43	0.5221

ESTIMATE Statement Results

Parameter	Estimate	Std Error	DF	t	Pr > t
blup pat 1	-3.86262200	3.17088923	18	-1.22	0.2389
blup pat 2	-0.25750813	3.17088923	18	-0.08	0.9362
blup pat 3	-0.99973746	3.17088923	18	-0.32	0.7562
blup pat 4	3.13554021	3.17088923	18	0.99	0.3358

Flowchart

Willing to assume the effects come from a distribution?



Assuming a factor is random involves extra assumptions but allows broader inferences.

Correlation in Mixed Models

Model:

Y_{ijk} = blood pressure for person k in condition (i,j) .

$$= \mu + p_k + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$$

Covariance:

$$\begin{aligned} \text{Cov}(Y_{ijk}, Y_{i'j'k}) &= \text{Cov}(p_k, p_k) \\ &= \text{Var}(p_k) \end{aligned}$$

$$\text{correlation} = \text{Var}(p_k) / [\text{Var}(p_k) + \text{Var}(\varepsilon_{ijk})]$$

Predicting the random effect

What if we assume a factor is random, but are interested in the individual levels of the random effects?

For the balanced data situation of the Propranolol data, the form of the best linear unbiased predictor is relatively simple and informative:

$$E[p_k | \mathbf{Y}] = E[p_k | \bar{Y}_{..k}] = \text{best predictor}$$

$$\begin{bmatrix} p_k \\ \bar{Y}_{..k} \end{bmatrix} \sim \mathbf{N} \left(\begin{bmatrix} 0 \\ \bar{\mu} \end{bmatrix}, \begin{bmatrix} \sigma_p^2 & \\ \sigma_p^2 & \sigma_p^2 + \sigma^2/n \end{bmatrix} \right),$$

where $\bar{\mu}$ is $\mu + \bar{\alpha} + \bar{\beta} + (\bar{\alpha}\bar{\beta})$.

$$E[p_k | \bar{Y}_{..k}] = E[p_k] + \text{cov}(p_k, \bar{Y}_{..k}) \text{var}(\bar{Y}_{..k})^{-1} (\bar{Y}_{..k} - \bar{\mu})$$

$$= 0 + \sigma_p^2 (\sigma_p^2 + \sigma^2/n)^{-1} (\bar{Y}_{..k} - \bar{\mu})$$

$$= \frac{\sigma_p^2}{\sigma_p^2 + \sigma^2/n} (\bar{Y}_{..k} - \bar{\mu})$$

Best Linear Unbiased Prediction

$$\text{BLUP}(p_k) = \frac{\sigma_p^2}{\sigma_p^2 + \sigma^2/n} (\bar{Y}_{..k} - \bar{Y}_{...})$$

[Shrinkage]

Similarly,

$$\begin{aligned} \text{BLUP}(\bar{\mu} + p_k) &= \bar{Y}_{...} + \frac{\sigma_p^2}{\sigma_p^2 + \sigma^2/n} (\bar{Y}_{..k} - \bar{Y}_{...}) \\ &= \bar{Y}_{...} \frac{\sigma^2/n}{\sigma_p^2 + \sigma^2/n} + \bar{Y}_{..k} \frac{\sigma_p^2}{\sigma_p^2 + \sigma^2/n} \\ &= \bar{Y}_{...} \alpha + \bar{Y}_{..k} (1-\alpha) \\ \alpha &= \frac{\sigma^2/n}{\sigma_p^2 + \sigma^2/n} \end{aligned}$$

[Weighted average]

In practice: EBLUP (Estimated BLUP)

$$\text{EBLUP}(p_k) = \frac{\hat{\sigma}_p^2}{\hat{\sigma}_p^2 + \hat{\sigma}^2/n} (\bar{Y}_{..k} - \bar{Y}_{...})$$

Numerical illustration:

$$\begin{aligned}\text{EBLUP}(p_1) &= \frac{\hat{\sigma}_p^2}{\hat{\sigma}_p^2 + \hat{\sigma}^2/n} (\bar{Y}_{..k} - \bar{Y}_{...}) \\ &= \frac{15.7976}{15.7976 + 85.7976/4} (81.75 - 90.86) \\ &= .424(-9.11) \\ &= -3.863\end{aligned}$$

$$\begin{aligned}\text{EBLUP}(\bar{\mu} + p_1) &= \bar{Y}_{...} + \frac{\hat{\sigma}_p^2}{\hat{\sigma}_p^2 + \hat{\sigma}^2/n} (\bar{Y}_{..k} - \bar{Y}_{...}) \\ &= 90.86 - 3.863 \\ &= 86.99\end{aligned}$$

Contrast this with the mean for the first patient, which is 81.75.

BLUPs In Linear Mixed Models

The best predicted value of a random effect given the data is $\tilde{\mathbf{u}} = E[\text{random effect}|\text{data}]$.

A BLUP minimizes MSE of prediction among linear unbiased predictors:

$$\text{minimize } E[(\tilde{\mathbf{u}} - \mathbf{u})^2]$$

among $\tilde{\mathbf{u}}$ which are linear in \mathbf{Y} and for which $E[\tilde{\mathbf{u}} - \mathbf{u}] = 0$.

For linear mixed models the best predictor is

$$\tilde{\mathbf{u}}_{BP} = \mathbf{DZ}'\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}),$$

while the best linear unbiased predictor is

$$\tilde{\mathbf{u}}_{BLUP} = \mathbf{DZ}'\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

“Shrinkage” estimator.

Bottom line: can be interested in the specific levels of a random factor.

Estimation and Tests in LMMs

Estimation of parameters by maximum likelihood or restricted maximum likelihood. Maximize the log of the likelihood.

Tests of fixed effects via approximate F- tests (SAS PROC MIXED).

Basic idea: Consider $H_0: \mathbf{k}'\boldsymbol{\beta} = 0$.

Could do a likelihood ratio test or a Wald test.

$$\text{var}(\mathbf{k}'\hat{\boldsymbol{\beta}}) \cong \mathbf{k}'(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})\mathbf{k}$$

$$\frac{\mathbf{k}'\hat{\boldsymbol{\beta}}}{\sqrt{\mathbf{k}'(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})\mathbf{k}}} \sim N(0,1) \text{ under } H_0$$

But need $\hat{\mathbf{V}}$ in place of \mathbf{V} .

Distribution?

Tests of variances of random effects

When using a maximum likelihood analysis the typical tests are based on the improvement in the maximized value of the log likelihood. The difference in twice the log likelihood is compared to a chi-square distribution to test for statistical significance. For testing whether a variance component is equal to zero the usual method must be slightly modified. Ordinarily we would take the difference in log likelihoods of the models with and without the random effect and compare that directly to a χ_1^2 cutoff point. The modification is to either calculate a p-value and then cut it in half, or to compare to a cutoff point with twice the nominal α level.

Why? The intuition is that testing

$$H_0: \sigma_p^2 = 0 \text{ versus } H_0: \sigma_p^2 > 0$$

is a one-sided test. The usual test is inherently two-sided and must be adjusted to reflect this fact.

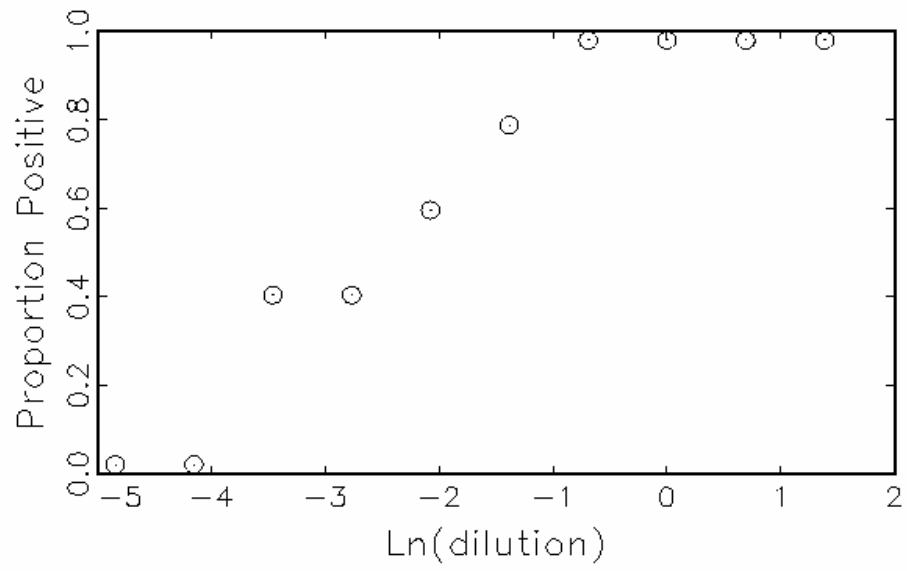
3. Review: Generalized Linear Models (GLMs)

Example: (from Finney, *Statistical Method in Bioassay*, 3rd Ed.). Study of growth of *Bacillus mesentericus* spores grown in dilutions of a potato flour suspension.

Dilution (g/100ml)	Spore Growth		Proportion
	Number of plates	Number positive	
1/128	5	0	0.0
1/64	5	0	0.0
1/32	5	2	0.4
1/16	5	2	0.4
1/8	5	3	0.6
1/4	5	4	0.8
1/2	5	5	1.0
1	5	5	1.0
2	5	5	1.0
4	5	5	1.0

Analysis?

Plot of Potato Flour Data



Analysis of potato flour data

Logistic Regression

Notation: Let x_i be the $\ln(\text{dilution})$ for the i th series and let Y_i be the number of positive plates.

Distribution: $Y_i \sim \text{indep. Binomial}(5, p(x_i))$,
which has mean $5p(x_i)$.

Model: $\ln(p(x_i)/(1-p(x_i))) = \alpha + \beta x_i$

S-shaped function of x

When $x = -\alpha/\beta$, $\alpha + \beta x = 0$,
and $p(x) = 1/2$.

Loglikelihood: $\sum_i y_i(\alpha + \beta x_i) - \log(1 + \exp(\alpha + \beta x_i))$

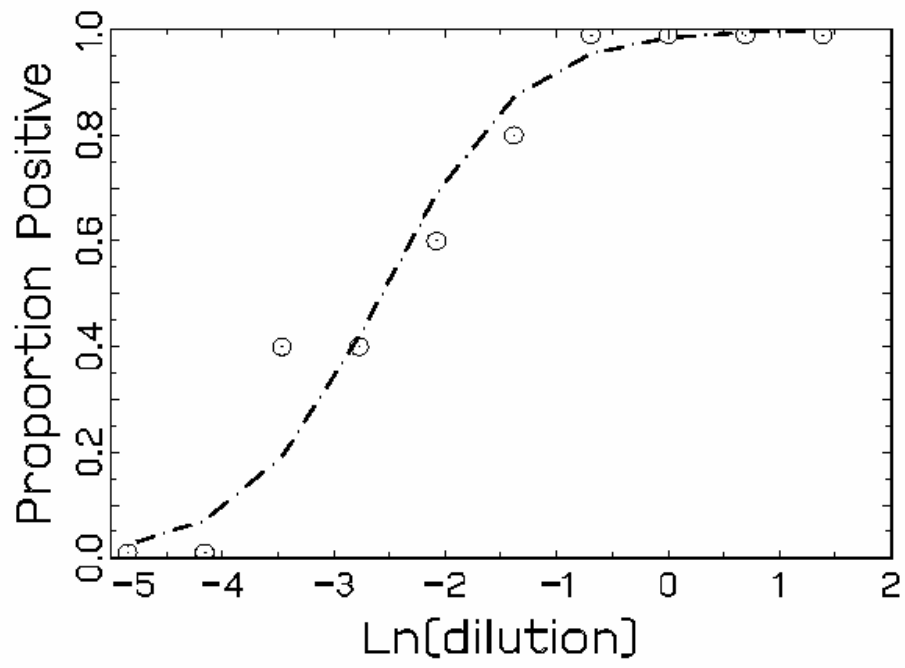
Maximum likelihood estimates:

$$\hat{\alpha} = 4.17$$

$$\hat{\beta} = 1.62$$

for $\ln(\text{dilution})$ which achieves 50%
positive results: $-4.17/1.62 = -2.57$

$$\exp(-2.57) = 0.076 \approx 1/13$$



GLMs

Dissect the modeling process into three distinct components:

1. What is the distribution of the data?
2. What aspect of the problem will be modelled?
3. What are the predictors?

In our example:

1. No. of successes in 5 trials => Binomial
2. log odds = $\ln(p/(1-p))$
3. $\ln(\text{dilution})$

GLMs

general case

$Y \sim$ distribution

$\mu =$ mean of Y

$g(\mu) = X\beta$

link function $g(\cdot)$

covariates $X\beta$

our example

$Y \sim$ Binomial

$np =$ mean of Y

$\ln(p/(1-p)) = \alpha + \beta x$

logit link

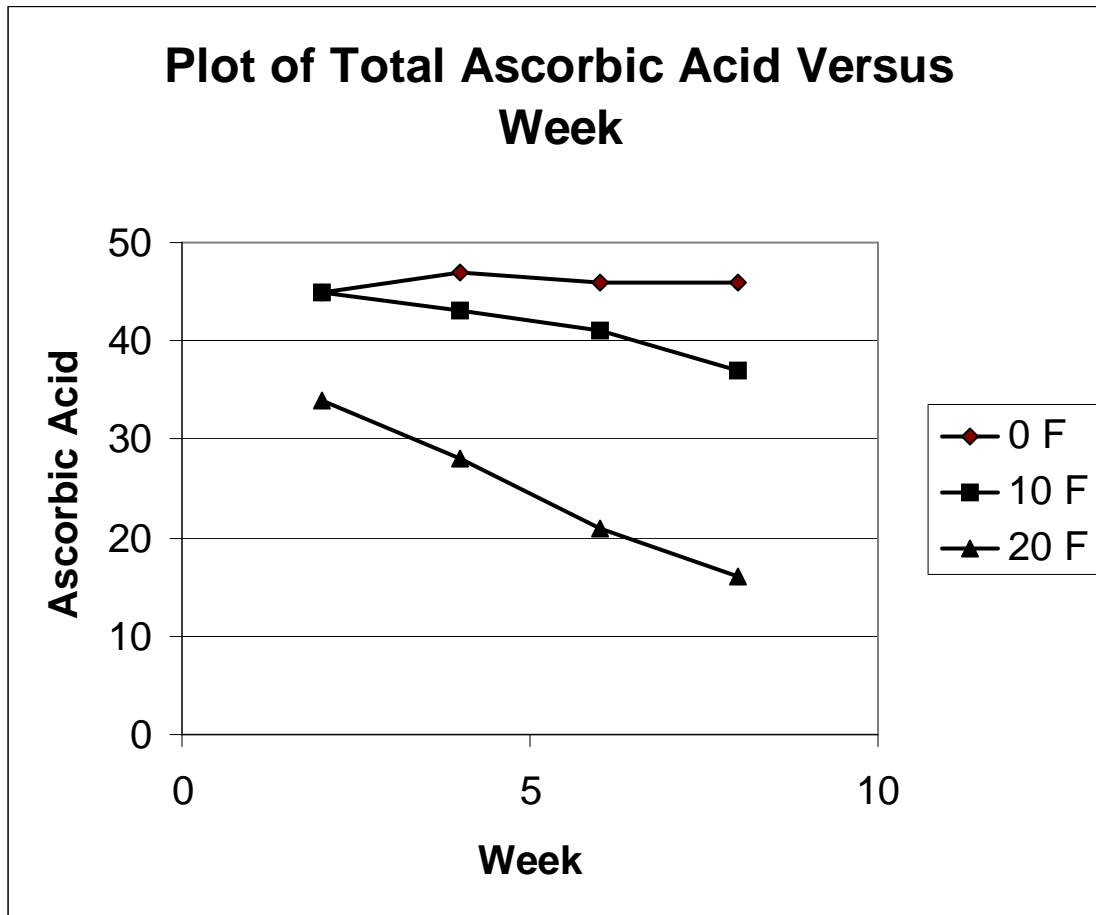
one predictor x

Example: (from Snedecor and Cochran, Sec 16.9, via McCullagh and Nelder, *Generalized Linear Models*). Amount of ascorbic acid remaining in snap beans after 2,4,6, and 8 weeks of storage at 0, 10 or 20 °F (a 3×4 factorial with 3 replicates per treatment combination).

Sum of three ascorbic acid determinations for each of 12 treatments on snap beans

Temp	Weeks of storage				Average
	2	4	6	8	
0	45	47	46	46	46.0
10	45	43	41	37	41.5
20	34	28	21	16	24.8
Ave	41.3	39.3	36.0	33.0	37.4

Here is a graph of the results.



Analysis: McCullagh and Nelder assume that the variance in ascorbic acid determination is constant on the original scale and wish to fit a model with exponential decline through time.

If Y_{ij} = average value at the i th temperature for week t_j , then a possible model is

$$Y_{ij} \sim \text{Normal}(\exp\{\alpha - \beta_i t_j\}, \sigma^2)$$

This is a generalized linear model for a Normal distribution with constant variation and with log link. The model has a common intercept and different slopes through time for each storage temperature.

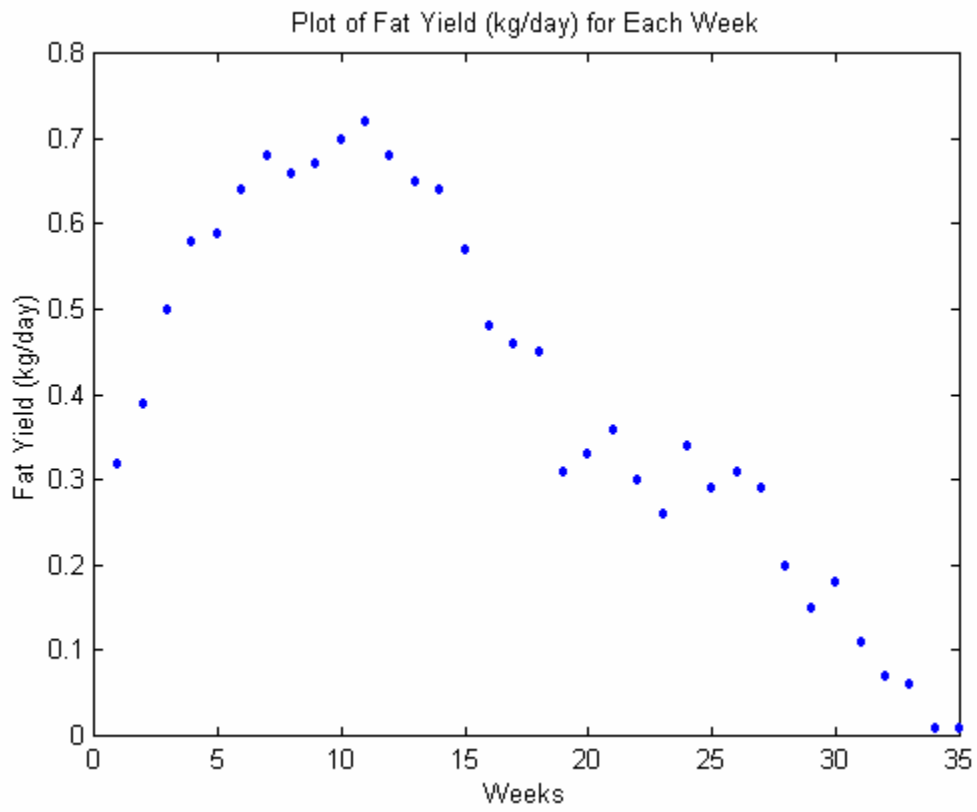
Another possibility: log transform.

$$\ln(Y_{ij}) \sim \text{Normal}(\alpha - \beta_i t_j, \tau^2)$$

Transform or Link?

Example: Average daily fat yield (kg/day) from milk from a single cow for each of 35 weeks.

0.31	0.39	0.50	0.58	0.59	0.64
0.68	0.66	0.67	0.70	0.72	0.68
0.65	0.64	0.57	0.48	0.46	0.45
0.31	0.33	0.36	0.30	0.26	0.34
0.29	0.31	0.29	0.20	0.15	0.18
0.11	0.07	0.06	0.01	0.01	



A typical model:

Fat yield “=” $\alpha t^\beta e^{\gamma t}$ where t=week

Transform:

$$\ln Y_i \sim N(\ln(\alpha) + \beta \ln(t_i) + \gamma t_i, \sigma^2)$$

$$\ln Y_i = \ln(\alpha) + \beta \ln(t_i) + \gamma t_i + \varepsilon_i$$

$$E[\ln Y_i] = \ln(\alpha) + \beta \ln(t_i) + \gamma t_i$$

$$Y_i = \alpha t_i^\beta e^{\gamma t_i} e^{\varepsilon_i}$$

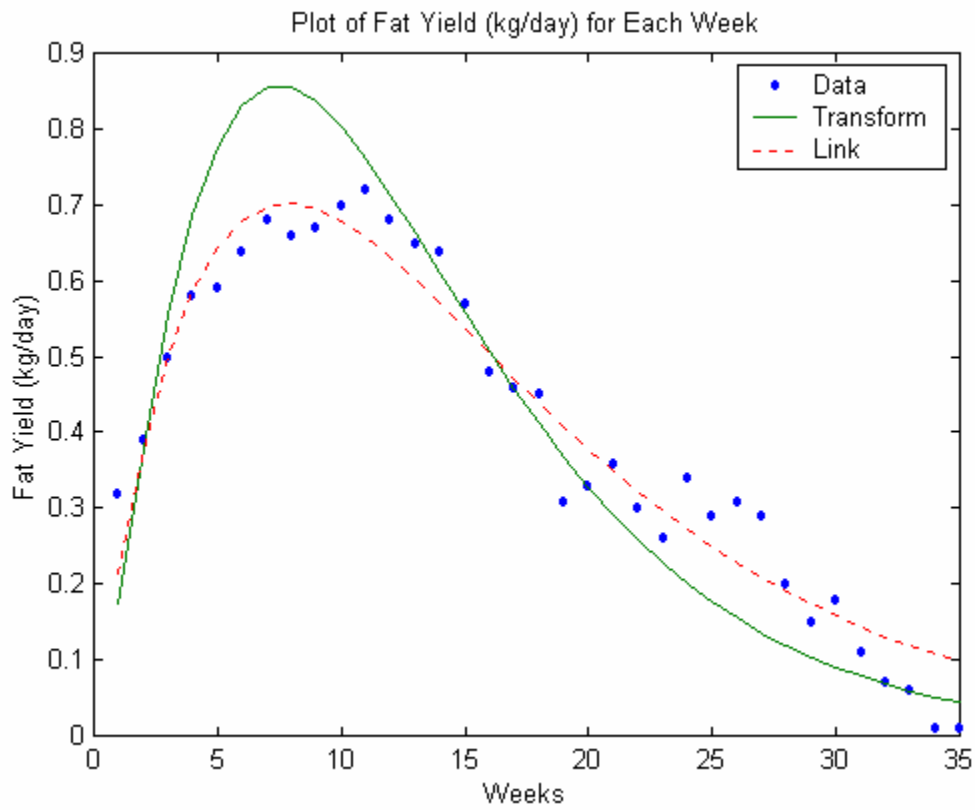
Link:

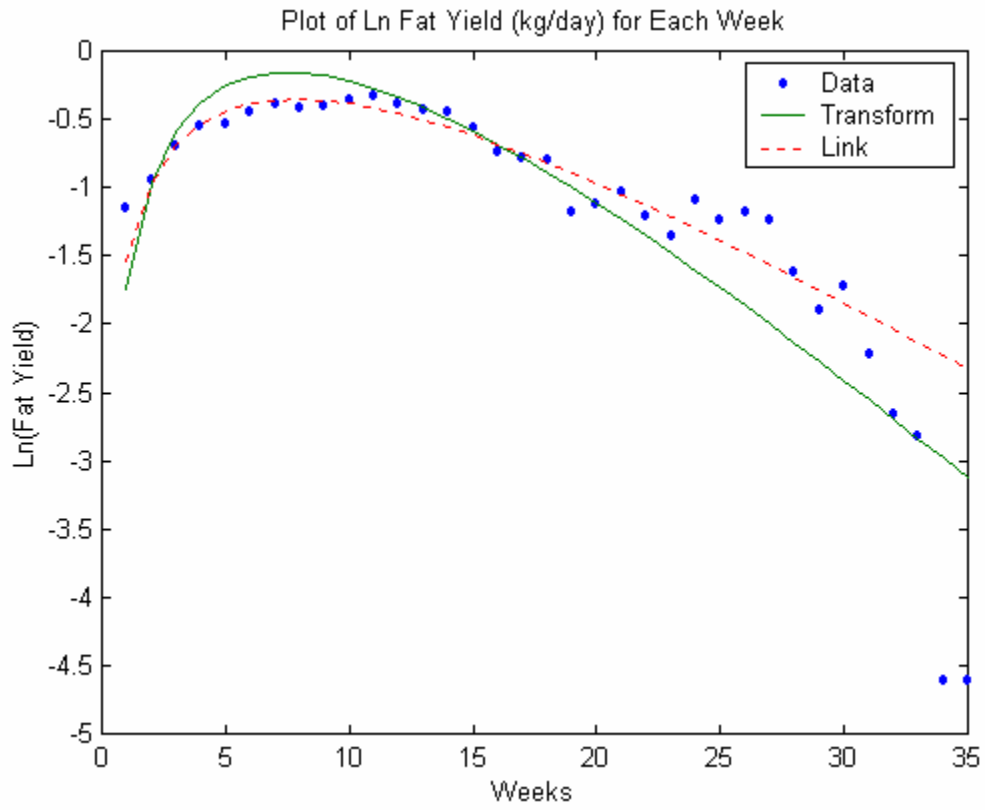
$$Y_i \sim N(\alpha t_i^\beta e^{\gamma t_i}, \tau^2)$$

$$E[Y_i] = \alpha t_i^\beta e^{\gamma t_i}$$

$$\ln(E[Y_i]) = \ln(\alpha) + \beta \ln(t_i) + \gamma t_i$$

$$Y_i = \alpha t_i^\beta e^{\gamma t_i} + \delta_i \quad \delta_i \sim N(0, \tau^2)$$





Homoscedasticity: The GLM analysis assumes a constant variance on the original scale. The transformed analysis assumes a constant variance on the transformed scale.

Trouble with transformations: With Poisson distributed data with zero counts, using a link function avoids the problems of a log transformation and zero counts.

See Ruppert, Cressie, and Carroll (1989) for a discussion.

Estimation and Tests in GLMs

Estimation of parameters by maximum likelihood or maximum quasi-likelihood. Maximize the log of the likelihood or quasi-likelihood.

Quasi-likelihood estimation:

Suppose $\text{Var}(Y) = \sigma^2 V(\mu)$

Define $U = \frac{Y - \mu}{\sigma^2 V(\mu)}$

and $Q(\mu; y) = \int_y^\mu \frac{y-t}{\sigma^2 V(t)} dt$

Note that $E[U] = 0$

$$\text{Var}(U) = \frac{1}{\sigma^2 V(\mu)}$$

$$-E\left[\frac{\partial U}{\partial \mu}\right] = \frac{1}{\sigma^2 V(\mu)}$$

which is similar to the properties of $\frac{\partial \ln f_{Y_i}}{\partial \mu}$.

For maximum likelihood we solve

$$\frac{\partial \ln L}{\partial \mu} = \sum_i \frac{\partial \ln f_{Y_i}}{\partial \mu} = 0.$$

For maximum quasi-likelihood we solve

$$\frac{\partial}{\partial \mu} \sum_i Q(\mu_i, y_i) = 0.$$

Example:

$$\sigma^2 = 1, V(\mu) = \mu$$

$$U = \frac{Y - \mu}{\mu}$$

$$\text{and } Q(\mu; y) = \int_y^\mu \frac{y-t}{t} dt$$

$$= y \int_y^\mu \frac{1}{t} dt - \int_y^\mu \frac{t}{t} dt$$

$$= y \ln(\mu) - y \ln(y) - (\mu - y)$$

So $\sum_i Q(\mu, y_i) = \sum_i y_i \ln(\mu) - n\mu + \text{constant}$.

For a Poisson,

$$\ln L = \sum_i Q(\mu, y_i) = \sum_i y_i \ln(\mu) - n\mu + \text{constant}$$

Measure of model (lack of) fit: deviance or Pearson chi-square statistic.

Deviance = $2(\text{max possible loglikelihood} - \text{loglikelihood of fitted model})$

So large values of deviance indicate a model which fits poorly.

Difference in Deviance for models 1 and 2 =
 $2(\text{loglik model 2} - \text{loglik model 1})$
= likelihood ratio statistic

Example: Potato flour dilutions (continued)

Maximum achievable loglikelihood = -12.597

Model 1: $\text{logit}(p(x_i)) = \alpha + \beta x_i$

ML estimates: $\hat{\alpha} = 4.1737$

$\hat{\beta} = 1.6226$

loglikelihood: -14.214

Deviance: $2(-12.597+14.214)$
 $= 3.234$ with $10-2 = 8$ d.f.

Model 2: $\text{logit}(p(x_i)) = \alpha$ (no slope)

ML estimate: $\hat{\alpha} = 0.4896$

Note: $1/(1+\exp(-0.4896))=.62=\text{ave prop.}$

loglikelihood: -33.203

Deviance: $2(-12.597+33.203)$
 $= 41.212$ with $10-1 = 9$ d.f.

Difference in deviance = $41.212 - 3.234$
 $= 37.978$ with 1 d.f.

Software

Software for LMMs and GLMs is readily available, either through special purpose routines, e.g., for logistic regression, or general routines. The package GLIM was the pioneer of software for GLMs, but other packages, e.g., SAS have caught up and now offer GLMs.

In SAS, PROC MIXED fits linear mixed models with the assumption of normality and PROC GENMOD fits generalized linear models.

Analysis of the potato flour data using GENMOD:

Program:

```
data one;  
set work.potflour;  
lndil=log(dilution);  
run;  
proc genmod descending;  
model nopos/noplate=lndil/dist=bin;  
run;
```

Output:

The GENMOD Procedure

Model Information

Data Set	WORK.ONE	
Distribution	Binomial	
Link Function	Logit	
Response Variable (Events)	Nopos	Nopos
Response Variable (Trials)	Noplate	Noplate
Observations Used	10	
Number Of Events	31	
Number Of Trials	50	

Criteria For Assessing Goodness Of Fit

Criterion	DF	Value	Value/DF
Deviance	8	3.2329	0.4041
Scaled Deviance	8	3.2329	0.4041
Pearson Chi-Square	8	2.7175	0.3397
Scaled Pearson X2	8	2.7175	0.3397
Log Likelihood		-14.2136	

Algorithm converged.

Analysis Of Parameter Estimates

Parameter	DF	Estimate	Standard Error	Wald	95% Confidence Limits	Chi-Square	Pr > ChiSq
Intercept	1	4.1737	1.2522	1.7194	6.6280	11.11	0.0009
lndil	1	1.6226	0.4571	0.7266	2.5185	12.60	0.0004
Scale	0	1.0000	0.0000	1.0000	1.0000		

NOTE: The scale parameter was held fixed.

4. Introduction to Generalized Linear Mixed Models (GLMMs)

Example: Similar to Abu-Libdeh, Turnbull, Clark, (Biometrics, 1990). Effect of selenium on prevention of skin cancer. 770 patients from seven clinics followed for four years.

Recorded:

response: Number of new basal cell epithelioma (BCE) sites found.

predictors: Selenium? (SEL), Sex (SEX), Exposure to the sun (SUN).

[Also?] Age, childhood farm exposure?, smoker?, skin damage, no. of tumors previously, clinic...

Q1: Does selenium decrease the number of BCEs?

Q2: Are some patients more sensitive to sun exposure? If so, which ones?

Features from a modelling viewpoint

Nature of response: count data

How to relate the response to the predictors?

$\lambda = \text{mean}$

$$\ln(\lambda) = \mu + \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}$$

=> Poisson regression

Problems: 1.

2.

A Generalized Linear Mixed Model

Let Y_{ij} be the response for patient i at visit j .

Y_{ij} = Number of new BCE sites

Assume $Y_{ij} \sim \text{Poisson}(\lambda_{ij})$, where λ_{ij} is the mean number of new BCEs for patient i at visit j .

$$\ln(\lambda_{ij}) = \mu_i + \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}$$

$$\mu_i \sim \text{Normal}(\mu, \tau_\mu)$$

assume a distribution for μ_i

$$\text{Cov}(\ln(\lambda_{ij}), \ln(\lambda_{ik})) = \tau_\mu$$

A correlation is induced in the model between observations taken on the same patient.

Other features

1. Assume a distribution on γ :

$$\ln(\lambda_{ij}) = \mu_i + \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma_i \text{SUN}$$

From the previous model:

γ = sun exposure effect (same across patients)

γ_i = sun exposure effect for the *i*th patient

$$\gamma_i \sim \text{Normal}(\gamma, \tau_\gamma)$$

- $\tau_\gamma > 0 \iff$ patients have different responses.
- Extreme values of γ_i indicate sensitive individuals.

2. Assume a distribution on β_2 :

$$\ln(\lambda_{ij}) = \mu_i + \beta_1 \text{SEX} + \beta_{2i} \text{SEL} + \gamma \text{SUN}$$

$$\beta_{2i} \sim \text{Normal}(\beta_2, \tau_\beta)$$

- If SEL is coded 1 for yes and 0 for no, then for the placebo group, the contribution of the $\beta_{2i} \text{SEL}$ term is zero, while for the treatment group it is β_{2i} . If $\tau_\beta > 0$, then the treatment group will have a larger variance.

Specifying GLMMs

1. What is the distribution of the data?
2. What aspects will be modelled?
3. What are the factors?
- * 4. Which factors will be assumed to have a distribution?

GLMMs

<u>general case</u>	<u>logit-normal</u>
$Y \sim$ distribution	$Y \sim$ Bernoulli
$\mu =$ mean of Y	$p =$ mean of Y
$g(\mu) = X\beta + Zu$	$\ln(p/(1-p)) = \beta x + u_i$
link function $g(\cdot)$	logit link
fixed factors $X\beta$	fixed factor x
random factors Zu	random intercepts u_i
$u \sim$ distribution	$u_i \sim$ Normal(μ_u, τ_u)

Prediction in GLMMs

In GLMMs we can adopt the same strategy as in LMMs:

- (1) Calculate $\tilde{\mathbf{u}} = \mathbf{E}[\mathbf{u} | \mathbf{Y}]$
- (2) Estimate any unknown parameters

However, either of these steps may be problematic.

5. Modeling in GLMMs

Example 1: Progabide and seizures (Diggle, Liang and Zeger, 1994).

Epileptics were randomly allocated to a placebo or an anti-seizure drug (Progabide) group. The number of seizures was recorded for a baseline period of 8 weeks and during consecutive two-week periods for four periods after beginning treatment. Is the drug effective at reducing the number of seizures?

Example 2: Cartoons and learning disabilities

This study concerned the comprehension of humor in two groups of adolescents (normal and learning disabled). Each subject was exposed to 24 different cartoons (in three types). There are two response variables whether or not the child got the cartoon and whether or not s/he liked it. The types of cartoon are: visual only, linguistic only, and both visual and linguistic.

Two questions of interest are: Is there a difference between normal and learning disabled children? How consistent are the responses within cartoon type?

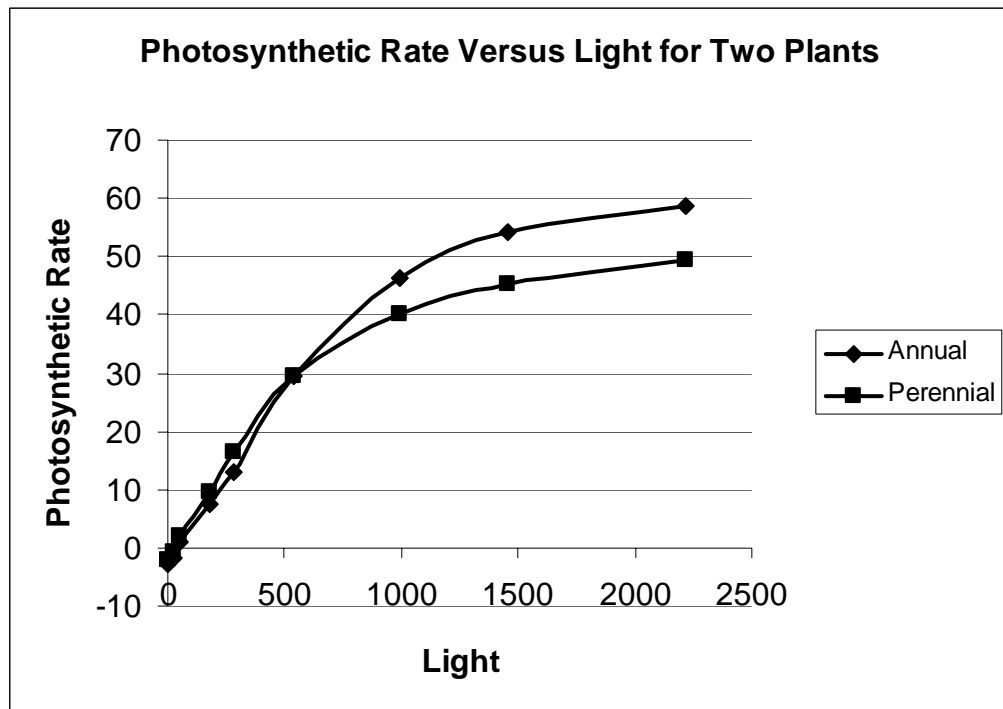
Example 3: Photosynthesis in corn relatives

Two species of corn relatives (an annual and perennial) are being compared with respect to photosynthetic physiology. Seeds from two populations of each species were collected and grown in the greenhouse. The experimental design was a randomized complete block design with four blocks and three seeds from each population in each block (for a total of 12 seeds per block). After 24 days, photosynthesis was recorded at nine different light levels from full sunlight to darkness on one individual from each population in each block (N=16). Measurements on the same 16 plants were repeated after 48 days. From these data, photosynthesis versus irradiance (PAR) response curves reflecting the change in photosynthetic rate with light level were derived.

The traits of interest are the maximum photosynthetic rate, dark respiration, the light compensation point, and the quantum yield. The maximum photosynthetic rate measures the maximum amount of carbon dioxide the plants

are able to assimilate in full sunlight, the dark respiration indicates how much carbon dioxide they respire in the dark, the light compensation point is the light level at which photosynthesis overcomes respiration and carbon assimilation becomes positive, and quantum yield is the efficiency of carbon assimilation at low light levels, or the slope of the light response curve as it crosses the light compensation point.

The main question of interest is to compare the two species with respect to their photosynthetic traits



Example 4: Chestnut Leaf Blight

The American chestnut tree was a predominant hardwood in the forests of the eastern United States, reaching 80-100 feet in height at maturity and providing timber and low-fat, high-protein nutrition for animals and humans in the form of chestnuts. In the early 1900's an imported fungal pathogen, which causes chestnut leaf blight, was introduced into the United States. The pathogen spread from infected trees in the New York City area and, by 1950, had killed over 3 billion trees and virtually eliminated the chestnut tree in the United States. Economic losses in both timber and nut production have been estimated in the hundreds of billions of dollars. As well, there are ecological impacts of eliminating a dominant species.

Attempts to restore this tree to the U.S. forests include

- development of blight resistant varieties
- weakening of the fungus by infecting it with a virus which reduces the fungus' virulence.

I'll describe the latter in more detail. The basic idea is to release hypovirulent isolates of chestnut blight fungus and let the viruses infect the natural populations of the fungus, thereby allowing chestnuts trees to survive.

Viruses spread between fungal individuals when they come in contact and fuse together. A major obstacle in spreading this virus and thus controlling the disease is that different isolates of the fungus cannot necessarily transfer the virus to one another.

Michael Milgroom - Cornell Plant Pathology, and his colleague, Paolo Cortesi - from the University of Milan, have worked with six incompatibility genes, which may block the transmission of this virus between isolates of the fungus.

To estimate the effects of these genes, they have paired numerous isolates which differ on the first gene only, the second gene only, the first and the second gene, etc. For each combination of isolates they have averaged about 30 attempts and record a binary response of whether or not the attempt succeeded in transmitting the virus.

Questions of interest include whether pre-identified genes actually do have an influence on transmission of the virus (and if so, to what degree), whether there are other, as yet unidentified, genes which might affect transmission, and whether transmission is symmetric. By symmetry of transmission we mean the following: suppose the infected fungus is type b at the locus for the first gene and the non-infected isolate (which we are trying to infect) is type B. The two isolates are the same at the other five loci. Is the probability of transmission the same as when using a type B to try to infect a type b?

Example 5: Combat vehicle design

Army combat vehicles of the future will likely locate crew stations deep within the vehicle, to achieve lower silhouettes and increased crew protection against ballistic and directed energy threats. This will require indirect vision systems such as liquid crystal displays.

In “Indirect Vision Driving With Fixed Flat Panel Displays for Near-Unity, Wide, and Extended Fields of Camera View” (Smyth, Gombash, Burcham, ARL-TR-2511, 2001) eight drivers tested each of four vision systems: direct and three types of indirect (with three different fields of view -- unity, wide and extended). Outcomes included speed to traverse the course, number of barrels knocked over and severe motion sickness (yes/no).

Example 6: Troponin and hemorrhage

Heart damage in patients experiencing brain hemorrhage has historically been attributed to pre-existing conditions. However, more recent evidence suggests that the hemorrhage itself can cause heart damage through the release of norepinephrine following the hemorrhage. To study this, researchers at UCSF measured cardiac troponin levels, an enzyme released following heart damage, at up to three occasions after patients were admitted to the hospital for a specific type of brain hemorrhage (subarachnoid hemorrhage or SAH).

The primary question was whether severity of injury from the hemorrhage was a predictor of troponin levels, as this would support the hypothesis that the SAH caused the cardiac injury. To make a more convincing argument in this observational study, we would like to show that severity of injury is an independent predictor, over and above other circulatory and clinical factors that would predispose the patient to higher troponin levels.

Possible clinical predictors included age, gender, history of heart failure, heart rate, whether the person was a smoker, diabetic or had high cholesterol levels. Circulatory status was described using systolic blood pressure, history of hypertension (yes/no) and left ventricular ejection fraction a measure of heart function. The severity of neurological injury was graded using a subject's Hunt-Hess score on admission. This score is an ordered categorical variable ranging from 1 (little or no symptoms) to 5 (severe symptoms such as deep coma).

The study involved 175 subjects with at least one troponin measurement and between 1 and 3 visits per subject.

6. Features of GLMMs

6 a) Consequences of model assumptions

What impact does this have on the distribution of Y ? Here are some calculations for the skin cancer example.

$$\begin{aligned} E[Y_{ij}] &= E[E[Y_{ij} | \mu_i]] \\ &= E[\exp\{\mu_i + \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}\}] \\ &= \exp\{\beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}\} E[\exp\{\mu_i\}] \end{aligned}$$

So $\log E[Y_{ij}] = \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN} + \log M_{\mu}(1)$, where $M_{\mu}(t)$ is the moment generating function of μ_i .

$$\begin{aligned} \text{Var}(Y_{ij}) &= \text{Var}[E[Y_{ij} | \mu_i]] + E[\text{Var}(Y_{ij} | \mu_i)] \\ &= \text{Var}(E[\exp\{\mu_i + \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}\}]) \\ &\quad + \exp\{\beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}\} E[\exp\{\mu_i\}] \\ &= \text{Var}(E[\exp\{\mu_i + \beta_1 \text{SEX} + \beta_2 \text{SEL} + \gamma \text{SUN}\}]) \\ &\quad + E[E[Y_{ij} | \mu_i]] \\ &> E[Y_{ij}] \end{aligned}$$

Marginal distribution for Probit models

$$Y_{ij} \sim \text{Bernoulli}(\Phi[\mu + a_i + \beta x_{ij}])$$

$$a_i \sim \text{Normal}(0, \tau_a).$$

What is the marginal distribution?

Persistence of links

Which other links “persist” like this?

Log link:

$$\begin{aligned} E[Y_i] &= E[E[Y_i | \mathbf{u}]] \\ &= E[\exp\{\mathbf{x}'_i \boldsymbol{\beta} + \mathbf{z}'_i \mathbf{u}\}] \\ &= \exp\{\mathbf{x}'_i \boldsymbol{\beta}\} E[\exp\{\mathbf{z}'_i \mathbf{u}\}] \end{aligned}$$

So $\log E[Y_i] = \mathbf{x}'_i \boldsymbol{\beta} + \log E[\exp\{\mathbf{z}'_i \mathbf{u}\}]$

(partial)

Logit link and other links do not persist.

6 b) Marginal versus conditional models

Illustration of difference of conditional and marginal approaches:

$Y_{ij} = 1$ if the i th woman miscarries during her j th pregnancy and is 0 otherwise.

$x_{ij} = j =$ pregnancy number.

Model:

$$E[Y_{ij}|u_i] = \Phi(\mu + \beta x_{ij} + u_i)$$

which gives

$$E[Y_{ij}] = \Phi\left(\frac{\mu + \beta x_{ij}}{\sqrt{1 + \sigma_u^2}}\right)$$

Interpretation of β ?

Advantages/disadvantages of the approaches

Because of computational problems with conditionally specified GLMMs there are many alternative methods (e.g., GEEs) for clustered data that focus on models for the marginal expectation of the response, $E[Y_{ij}]$.

Marginal models have the following advantages:

- Marginal models avoid the specification of the conditional structure, so misspecification of this portion of the model can be avoided.
- For example, when the underlying random effects distribution is heteroscedastic, assuming it is homoscedastic and using a conditional approach can lead to biased estimators (Heagerty and Kurland, 2001)
- When paired with a GEE approach to estimation, estimates of the marginal parameters are consistent, even under misspecification of the association structure.

Major drawbacks of the marginal approach include:

- Often does not measure covariate effects of primary scientific interest.
- In extreme circumstances, features of scientific interest present in every conditional model may not be present in the marginal model.
- Marginal quantities can be calculated from a conditional model but the converse is not typically true.
- Marginal modeling approaches are susceptible to Simpson's paradox and the Ecological Fallacy, potentially giving misleading results.
- If the question of interest is based on the marginal distribution, a longitudinal design may not be the most appropriate.

For a more detailed critique of marginal modeling see Lindsey and Lambert (1998).

7. Inference for GLMMs

Estimation: Maximum likelihood (or variants) based on normality assumptions are relatively standard for linear mixed models. For example, SAS PROC MIXED using ML or REML.

For many GLMs, maximum likelihood is also standard, e.g., logistic regression or Poisson regression.

What about GLMMs?

A simple GLMM

A logit-normal model:

$$Y_{ij} | u \sim \text{Bernoulli}(p_{ij}),$$

$$i=1,2, \dots, n; j=1,2, \dots, q.$$

q clusters, n observations per cluster.

$$\ln(p_{ij}/(1-p_{ij})) = \beta x_{ij} + u_j$$

logit link

one fixed and one random factor

$$u_j \sim \text{Normal}(0, \sigma^2)$$

Scenario:

$Y_{ij} = 1$ if blood pressure on day i on individual j decreases after using medicine at dose x_{ij} , 0 otherwise.

q individuals, n days of measurement on each.

u_j is the individual specific propensity to increase or decrease blood pressure.

ML Estimation?

$$\begin{aligned}
 \text{Likelihood} &= P\{Y=y|\beta, \sigma^2\} \\
 &= \int P\{Y=y|\beta, \sigma^2, u\}f(u)du \\
 &= \int P\{Y=y|\beta, u\}f(u)du \\
 &= \int \prod_{i,j} P\{Y_{ij}=y_{ij}|\beta, u\} f(u)du \\
 &= \prod_j \int \prod_i P\{Y_{ij}=y_{ij}|\beta, u_j\}f(u_j)du_j \\
 &= \\
 &\prod_j \int \exp\{\beta \sum_i Y_{ij}x_{ij} + Y_{+j}u_j\} \prod_i (1 + \exp\{\beta x_{ij} + u_j\})^{-1} \times \\
 &\quad \exp\{-u_j^2/2\sigma^2\} / (2\pi\sigma^2)^{1/2} du_j.
 \end{aligned}$$

Cannot be evaluated in closed form but is not too hard to do numerically for this example.

Brute force ML

When the model has a single random effect or two nested random effects, it is relatively easy to evaluate the integrals in the likelihood. For example, with a single random factor we have seen that the likelihood is a product of one-dimensional integrals.

One can then maximize the likelihood numerically to find ML estimates and to perform likelihood ratio tests.

Numerical evaluation of the likelihood

When there is a single, normally distributed random effect, the likelihood can be written as a product of integrals of the form:

$$\int_{-\infty}^{+\infty} g(x) \exp\{-x^2\} dx$$

These can be accurately evaluated using Gauss-Hermite quadrature:

$$\int_{-\infty}^{+\infty} g(x) \exp\{-x^2\} dx \approx \sum_i w_i g(x_i)$$

The weights, w_i , and the evaluation points, x_i , can be found in books on numerical integration, e.g., Abramowitz and Stegun (1964).

In general, however, the evaluation of the likelihood can be quite difficult. For the general case,

$$\int \dots \int_{\text{dim of } u} \exp(\sum_i Y_i (x_i' \beta + z_i' u)) \prod_i (1 + \exp(x_i' \beta + z_i' u))^{-1} dF(u).$$

The dimension of u can get large quickly. For example, in the leaf blight data, the dimension of the integral is larger than 250!

What to do?

Other approaches to ML

Simulation approximations

Monte Carlo EM

Monte Carlo Newton-Raphson

Stochastic approximation

Importance sampling

Inference using ML would proceed using the usual asymptotic approximations:

ML estimates are asymptotically normal, with SEs coming from second derivatives of the log likelihood.

Tests would be based on the likelihood ratio test, comparing $-2\log\text{likelihood}$ for nested models.

Best predicted values would be estimated by calculating $E[\text{random effect}|\text{data}]$ and plugging in ML or REML estimates. In general, the conditional expected values can't be evaluated in closed form either.

Tests on variances of random effects The usual asymptotic theory breaks down when testing whether the variance components are equal to zero just as with LMMs. For example, in testing whether a single variance component is zero, the large-sample distribution under H_0 is a 50:50 mixture of a χ_1^2 and 0.

Summary: ML

- + Known large sample properties
- + Likelihood ratio tests
- Hard to compute for many GLMMs
- Small sample performance needs to be assessed for any particular model.

Conditional Inference

A very different approach to random effects is to treat them as nuisance parameters and condition them away.

Classic situation: Matched pairs binary logistic regression.

Example: Do cancer patients get more effective treatment in a major cancer center or a community hospital? Can't directly compare rates. Patients are matched on treatment date, treatment, protocol and other factors. The response is whether or not there is a large shrinkage in their tumor within 90 days.

Data: (1=shrinkage, 0=no shrinkage).

Pair	Cancer Center	Community Hospital
1	1	1
2	1	0
3	1	1
.		
.		
.		
936	0	1

A model:

$Y_{ij} = 1$ for shrinkage and 0 otherwise. i indexes pairs ($i=1,2,\dots,N$) and j indexes treatment (with j being 1 for a comm hosp. and 2 for a cancer center).

$$Y_{ij} \sim \text{Bernoulli}(p_{ij})$$

$$\text{logit}(p_{ij}) = \mu_i + \beta x_{ij},$$

where $x_{ij} = 0$ for $j=1$ and 1 for $j=2$ (cancer center or “treatment” indicator).

μ_i treated as fixed parameters

Maximum likelihood gives

$$\hat{\beta} = 2 \log \frac{N_{01}}{N_{10}},$$

where N_{10} is the number of pairs with $Y_{i1}=1$ and $Y_{i2}=0$ and N_{01} is the number of pairs with $Y_{i1}=0$ and $Y_{i2}=1$.

This is perhaps easiest to visualize in a 2×2 format:

	Treatment	
Control	Failure	Success
Failure	N_{00}	N_{01}
Success	N_{10}	N_{11}

The ML estimator is twice the sensible answer.

Remedy? A commonly used approach is that of conditional likelihood.

Basic idea: Derive the sufficient statistics for the μ_i and work with the conditional distribution given those sufficient statistics.

From the form of the density it is clear that the sufficient statistic is $(S_1, S_2, \dots, S_N, T) = (Y_{1.}, Y_{2.}, \dots, Y_{N.}, Y_{.2})$. Since the distribution is discrete, to find the distribution of \mathbf{S} we merely sum over the appropriate values of \mathbf{Y} :

$$f_{\mathbf{S},T}(\mathbf{s},t) = \sum_{\mathbf{y}: s_i = y_{i.}, t = y_{.2}} f_{\mathbf{Y}}(\mathbf{y})$$

$$= C(\mathbf{s},t) \frac{e^{\sum \mu_i s_i + \beta t}}{d},$$

where $C(\mathbf{s},t)$ represents the number of combinations of values of \mathbf{y} that satisfy the constraints.

From this it is straightforward to get the marginal distribution of \mathbf{S} :

$$\begin{aligned} f_{\mathbf{S}}(\mathbf{s}) &= \sum_z f_{\mathbf{S},T}(\mathbf{s},z) \\ &= \sum_z C(\mathbf{s},z) \frac{e^{\sum \mu_i s_i + \beta z}}{d} \end{aligned}$$

and the conditional distribution of T given \mathbf{S} :

$$\begin{aligned} f_{T|\mathbf{S}}(t|\mathbf{s}) &= f_{\mathbf{S},T}(\mathbf{s},t) / f_{\mathbf{S}} \\ &= \frac{C(\mathbf{s},t)e^{\beta t}}{\sum_z C(\mathbf{s},z)e^{\beta z}} \end{aligned}$$

None of the μ_i remain, as expected. This conditional likelihood can thus be used to estimate β or to form tests or confidence intervals.

For the matched pairs situation the combinatorial coefficient is straightforward to evaluate.

Conditional on $S_i = 0$ we know $Y_{i1}=0$ and $Y_{i2}=0$.
Conditional on $S_i = 2$ we know $Y_{i1}=1$ and $Y_{i2}=1$.

The only remaining randomness involves those pairs for which $S_i = 1$.

Using $r = t - N_{00} - N_{11}$ = number of successes in the discordant pairs, is equivalent to using t .

Then it isn't hard to show that

$C(\mathbf{s},t)$ = number of ways the successes in the N_{10} and N_{01} pairs can be distributed

$$= \binom{N_{10} + N_{01}}{r} = \binom{N'}{r}$$

Illustration: The conditional approach discards the $132+501 = 633$ responses which are concordant and bases the analysis on the 303 remaining.

$$\text{p-value} = 2 \times \Pr\{X \leq 146\}$$

where $X \sim \text{Binomial}(303, 1/2)$.

So $\text{p-value} = 2(0.283) = 0.566$.

Drawbacks to the conditional approach

Recover information from concordant pairs?

Inferences about random effects?

Between versus within “subjects.”?

Generalized Estimating Equations (GEEs)

GEEs are a computationally less demanding method than ML estimation. They are applicable (mainly) to longitudinal data.

Longitudinal data = data collected on a subject on two or more occasions.

Number of occasions is small compared to the number of subjects.

Longitudinal Data

Begin by considering longitudinal data with linear models under normality.

(1) Separate effects that are constant across subjects (β) from those which vary across subjects (\mathbf{u}_i).

(2) For the i th individual write a linear model conditional on the value of \mathbf{u}_i :

$$\mathbf{Y}_i = \mathbf{X}_i\beta + \mathbf{Z}_i\mathbf{u}_i + \varepsilon_i$$

$$\varepsilon_i \sim N(\mathbf{0}, \mathbf{R}_i)$$

(3) Incorporate subject-to-subject variability by assigning a distribution to \mathbf{u}_i :

$$\mathbf{u}_i \sim N(\mathbf{0}, \mathbf{D}).$$

Result: $\mathbf{Y}_i \sim \text{indep } N(\mathbf{X}_i\beta, \mathbf{Z}_i\mathbf{D}\mathbf{Z}_i' + \mathbf{R}_i)$

Longitudinal Data

Example: (Diggle, Liang and Zeger, 1994).
Milk was collected from 79 cows on one of three diets: barley, lupins, and a mixture of both. Protein content of the milk was recorded weekly for 19 weeks after the earliest calving.

Constant effects: diet, time

Effects that vary across animals: intercepts

Model for the i th cow on diet j , at time t

$$Y_{ijt} = \mu + c_{i(j)} + \alpha_j + f(t) + e_{ijt}$$

$$\mathbf{e}_{ij} \sim N(\mathbf{0}, \mathbf{R}_{i(j)})$$

$$\mathbf{R}_{i(j)}: \text{cov}(\mathbf{e}_{ijt}, \mathbf{e}_{ijt'}) = \sigma_e^2 \exp(-\phi|t-t'|)$$

$$c_{i(j)} \sim N(0, \sigma_c^2)$$

(More generally for awhile):

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

$$\mathbf{u} \sim \mathbf{N}(\mathbf{O}, \mathbf{D})$$

$$\mathbf{e} \sim \mathbf{N}(\mathbf{O}, \mathbf{R})$$

So $\mathbf{Y} \sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}=\mathbf{Z}\mathbf{D}\mathbf{Z}' + \mathbf{R})$.

What about using $\hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$?

$\hat{\beta}_{OLS}$ is unbiased.

$$\begin{aligned} E[\hat{\beta}_{OLS}] &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E[\mathbf{Y}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta = \beta \end{aligned}$$

$\hat{\beta}_{OLS}$ is usually fairly efficient.

$$\text{Var}(\hat{\beta}_{OLS}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$$

(As compared to $\text{Var}(\hat{\beta}_{GLS}) = (\mathbf{X}'\mathbf{V}\mathbf{X})^{-1}$)

In fact, with balanced designs, $\hat{\beta}_{OLS} = \hat{\beta}_{GLS}$.

So why not just use $\hat{\beta}_{OLS}$ and standard software?

$$\text{Var}(\hat{\beta}_{OLS}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$$

but, using standard software,

$\hat{\text{Var}}(\hat{\beta}_{OLS}) = (\mathbf{X}'\mathbf{X})^{-1}\hat{\sigma}^2$, which will often be very wrong. That is, the OLS estimate isn't so bad, but the usual variance estimate is way off.

Going back to the longitudinal data setting the basic idea is, with $Y_i \sim$ independently, to use the “replication” across subjects to get an empirical estimate of the variance. For the longitudinal data setting,

$$\hat{\beta}_{OLS} = (\sum_i \mathbf{X}'_i \mathbf{X}_i)^{-1} (\sum_i \mathbf{X}'_i \mathbf{Y}_i)$$

$$\text{Var}(\hat{\beta}_{OLS}) = (\sum_i \mathbf{X}'_i \mathbf{X}_i)^{-1} (\sum_i \mathbf{X}'_i \mathbf{V}_i \mathbf{X}_i) (\sum_i \mathbf{X}'_i \mathbf{X}_i)^{-1}$$

which can be estimated by

$$(\sum_i \mathbf{X}'_i \mathbf{X}_i)^{-1} (\sum_i \mathbf{X}'_i (\mathbf{Y}_i - \hat{\mu}_i) (\mathbf{Y}_i - \hat{\mu}_i)' \mathbf{X}_i) (\sum_i \mathbf{X}'_i \mathbf{X}_i)^{-1}$$

For the milk protein data from Diggle, Liang and Zeger (1994), if all the animals had all 19 weeks of data we could just get empirical estimates from the multivariate observations.

With some missing data the previous formula can still be used.

Non-normal data?

GEEs work most easily for models specified on the unconditional distribution. In contrast, we have been specifying models which are conditional on the random effects, u .

For example, for binary data, we could specify:

$$\begin{aligned} E[Y_{ij}] &= p_{ij} \\ \text{logit}(p_{ij}) &= \mathbf{X}_i\boldsymbol{\beta}. \end{aligned}$$

Obtain $\hat{\boldsymbol{\beta}}$ by solving the GEE:

$$\sum_{i=1}^n \left(\frac{\partial \mathbf{p}_i}{\partial \boldsymbol{\beta}} \right)' W \text{Var}(\mathbf{Y}_i)^{-1} (\mathbf{Y}_i - \mathbf{p}_i) = 0,$$

where $W \text{Var}$ indicates a “working” or assumed covariance structure, possibly dependent on unknown parameters.

This has properties similar to the estimating equations for the LMM:

$$\sum_{i=1}^n \mathbf{X}_i' \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) = 0$$

A big advantage of the GEE approach is the ability to use a “robust” variance estimate.

In such a case the inferences about the mean structure are asymptotically valid, even when the working variance is incorrect.

This offers a useful tool for inference or, at least, model checking.

GEEs are most naturally adapted to marginal models, not the conditional random effects models of GLMMs. But see Zeger, Liang and Albert (1988) see for some results in this direction.

In addition to the drawbacks above relating to marginal models, the GEE approach in particular also has the following drawbacks compared to GLMMs:

- GEEs by themselves do not help to separate out different sources of variation.
- GEEs are not directly a technology for best prediction of random effects. But see Waclawiw and Liang (1993) and Heagerty (1999).
- GEEs are not the best technique for other-than-longitudinal (but correlated) data, either crossed or nested random factors.
- GEEs may be inefficient when the goal is estimation of the variance covariance structure.

Summary: GEEs

Mainly for longitudinal data.

Easiest for marginal models, not random effects models: GEEs by themselves do not help to separate out sources of variation that may be present and do not provide predicted values.

Robust standard errors:

- + Robust
- + Often relatively efficient
- Estimates many parameters
- Does not work well when the number of time points is large compared to the number of subjects
- Does not work well with missing data

Penalized Quasi-likelihood (PQL)

$\mathbf{Y} \sim$ exponential family with mean $\boldsymbol{\mu}$

$$\mathbf{g}(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} \quad \mathbf{u} \sim \mathbf{N}(\mathbf{0}, \mathbf{D})$$

$$\mathbf{g}(\mathbf{y}) \approx \mathbf{g}(\boldsymbol{\mu}) + (\mathbf{y} - \boldsymbol{\mu})\mathbf{g}'(\boldsymbol{\mu}) \equiv \mathbf{z}$$

$$= \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + (\mathbf{y} - \boldsymbol{\mu})\mathbf{g}'(\boldsymbol{\mu})$$

$$= \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}\mathbf{g}'(\boldsymbol{\mu})$$

Idea: treat \mathbf{z} as a LMM with

$$\text{Var}(\mathbf{z}) = \mathbf{Z}\mathbf{D}\mathbf{Z}' + \mathbf{R}(\mathbf{g}'(\boldsymbol{\mu}))^2$$

Use the Mixed Model Equations iteratively to find both

$\hat{\boldsymbol{\beta}}$ and the BLUP of \mathbf{u}

Schall (1991) also suggests ways to get approximate SEs.

Summary: PQL

- + Computationally fairly easy
- + Works well when the data are approximately normal to start with.
- Does not work well for highly non-normal data (e.g., binary).
- Only for $\mathbf{u} \sim \text{Normal}$.

Why PQL? See Breslow and Clayton (1993).

Other Approaches

1. Models for specific situations.
 - Beta-binomial (Crowder, 1978)
 - Poisson-gamma (Abu-Libdeh, et al, 1990)
 - Other (Conaway, 1990)

3. Other marginal models
 - Liang, Zeger and Qaqish (1992)

4. BLUP estimators
 - Engel and Keen (1994)
 - McGilchrist (1994,1995)

5. Maximum hierarchical likelihood.
 - Lee and Nelder (1996)

More on the beta-binomial

Scenario: A potentially toxic chemical is administered to pregnant rats in the treatment group(s). There is also a control group. The response we record is the presence or absence of a birth defect in animal k from litter j in group i .

$$Y_{ijk} \mid p_{ij} \sim \text{indep. Bernoulli}(p_{ij})$$

$$p_{ij} \sim \text{indep. Beta}(\alpha_i, \beta_i)$$

Hence $Y_{ijk} \sim \text{Bernoulli}(\mu_i)$, where μ_i is given by $E[p_{ij}] = \alpha_i / (\alpha_i + \beta_i)$.

The joint density of \mathbf{Y} is given by

$$f_{\mathbf{Y}} = \prod_{i,j} f_{\mathbf{Y}_{ij}},$$

where $\mathbf{Y}_{ij} = (Y_{ij1}, Y_{ij2}, \dots, Y_{ijn_{ij}})'$. Dropping the i and j subscripts,

$$\begin{aligned} f_{\mathbf{Y}} &= \int_0^1 \prod_k p^{Y_k} (1-p)^{(1-Y_k)} p^{\alpha-1} (1-p)^{\beta-1} / B(\alpha, \beta) dp \\ &= \int_0^1 p^{\alpha+Y_{\cdot}-1} (1-p)^{n-Y_{\cdot}+\beta-1} / B(\alpha, \beta) dp \\ &= \frac{B(\alpha+Y_{\cdot}, \beta+n-Y_{\cdot})}{B(\alpha, \beta)} \end{aligned}$$

Therefore, the likelihood is given by

$$L = \prod_{ij} \frac{B(\alpha_i + Y_{ij\cdot}, \beta_i + n_{ij} - Y_{ij\cdot})}{B(\alpha_i, \beta_i)}.$$

Extensions? E.g., different doses of the toxic chemical?

Software

Maximum likelihood

Linear Normal Mixed Models: SAS PROC MIXED or SPSS.

Linear Normal Nested Models: MlwiN (<http://multilevel.ioe.ac.uk>) and HLM (<http://www.ssicentral.com/hlm/hlm.htm>) fit hierarchical models, using maximum likelihood for normal data and penalized quasi-likelihood for binary and binomial data (see below).

Logit/Probit normal, Ordinal logit: MIXOR program runs on PCs available free from Don Hedeker via the WWW at <http://www.uic.edu/~hedeker/mix.html>.

Nonlinear normal mixed models: S-Plus functions free from Pinheiro, Bates and Lindstrom at: <http://www.stat.wisc.edu/p/stat/ftp/src/NLME/>

SAS NLMIXED (new in Version 7.0) can handle random effects for the longitudinal data situation (i.e., data are in clusters).

GEE software

SAS GENMOD allows GEE estimation through its REPEATED statement. SUDAAN and STATA allow GEE estimation for a variety of statistical methods including multinomial logistic regression.

PQL software

GLIMMIX macros available from SAS at

<http://ftp.sas.com/techsup/download/stat/glmm800.html>

For nested models MlwiN and HLN use PQL and improvements of PQL.

Bayes software:

BUGS fits a wide variety of Bayesian models and allows the incorporation of distributions for the parameters. A description of BUGS (and it can be downloaded from)

<http://www.mrc-bsu.cam.ac.uk/bugs/>

Case Studies

Case study 1: Breeding Bird Survey. (James, et al, 1996). Counts of number of birds “sighted” has been made each June at thousands of locations across the U.S. and Canada. Many of the locations have been surveyed since the mid 1960s. Responses are summarized by estimating whether the trend in population size is positive within a stratum.

response: increase (yes/no) for species i in stratum j .

distribution: Bernoulli *link:* probit

predictors: species (fixed), stratum (random).

Question: Is destruction of overwintering habitat causing the decline of neo-tropical migrant bird populations on a continent-wide basis?

Breeding Bird Survey (cont)

Model: Y_{ij} = (1/0) increase for species i in stratum j ? (Probit-normal)

$$Y_{ij} \sim \text{Bernoulli}(p_{ij}) \quad i=1,\dots,26; j=1,\dots,37$$

$$p_{ij} = \Phi(\mu_i + s_j), \quad s_j \sim N(0, \sigma_s^2)$$

Data layout

Species	Stratum							
	1	...	10	11	12	13	...	37
1			0			0		
2								
3								
4	0							
.					0	0		
26				1	1	1	0	0

blank = species not present (about 2/3)

1 = increase

0 = decrease

ML Estimates:

$$\hat{\mu}_1 = -0.66, \hat{\mu}_2 = 0.26, \dots, \hat{\mu}_{26} = -1.28$$

$$\hat{\sigma}_s^2 = 0.45$$

Interpretations:

$$\Phi(0.26) = 0.60$$

$$\hat{\sigma}_s^2 / (\hat{\sigma}_s^2 + 1) = 0.31$$

Test of $\sigma_s^2 = 0$:

$$\text{diff in } -2\log\text{lik} = 11.88$$

$$\text{compare to a } \frac{1}{2}\chi_1^2$$

Estimated best predicted values: $E[s_j|Y]$

$$\text{e.g., } E[s_{23}|Y] = -1.10$$

$$\Phi(-1.10) = 0.14$$

Case Study 2: Progabide and Seizures (Diggle, Liang and Zeger, 1994). Epileptics were randomly allocated to a placebo group or an drug (Progabide) group. The number of seizures was recorded for a baseline period of 8 weeks and during consecutive two-week periods for 4 periods after beginning treatment. Is the drug effective at reducing the number of seizures?

Patient	Number of seizures					Trt
	Base -line	Period 1	Period 2	Period 3	Period 4	
1	11	5	3	3	3	0
2	11	3	5	3	3	0
3	6	2	4	0	5	0
4	8	4	4	1	4	0
.
57	13	0	0	0	0	1
58	12	1	4	3	2	1

response: number of seizures for individual i at
time $j=1,2,3,4,5$

distribution: Poisson

predictors: period, treatment (both fixed),
individual, individual \times treatment (?) (both
random).

The baseline period is 8 weeks long, whereas the
observation periods are only 2 weeks long.

Question: Does Progabide reduce the frequency
of seizures?

Model:

Y_{ij} = count for subject i at time j

t_{ij} = time (in weeks) for the observation period for subject i at time j (either 8 or 2 weeks).

$$Y_{ij} | \lambda_{ij} \sim \text{indep. Poisson}(\lambda_{ij})$$

$$\ln(\lambda_{ij}) = \mu + s_i + \beta_1 \text{TIME}_{ij} + \beta_2 \text{TRT}_{ij} + \beta_3 \text{TIME}_{ij} \times \text{TRT}_i + \ln(t_{ij})$$

$$s_i \sim N(0, \sigma_s^2)$$

$\text{TIME}_{ij} = 1$ if the observation is post baseline and 0 otherwise.

Mainly interested in β_3 .

How to estimate this model?

SAS Programs for the Progabide data

```

data thall;
input id y visit trt bline age;
cards;
104 5 1 0 11 31
104 3 2 0 11 31
104 3 3 0 11 31
103 0 4 1 19 20
...
232 0 4 1 13 36
236 1 1 1 12 37
236 4 2 1 12 37
236 3 3 1 12 37
236 2 4 1 12 37
;

data new;
  set thall (drop=age);
  output;
  if visit=1 then do; y=bline; visit=0; output; end;
run;

proc sort;
  by id visit;
run;

data new3;
  set new;
  if id ne 207;
  if visit=0 then do; time=0; ltime=log(8); end;
  else do; time=1; ltime=log(2); end;
run;

proc nlmixed data=new3 qpoints=20;
  parms mu=1 bl=0 b2=0 b3=0 sig1=0.1;
  eta=mu+bl*time+b2*trt+b3*time*trt+u1+ltime;
  lam=exp(eta);
  model y~Poisson(lam);
  random u1~Normal(0,sig1) subject=id;
run;

proc nlmixed data=new3 qpoints=20;
  parms mu=1 bl=0 b2=0 b3=0 sig1=0.1 cov=0.05 sig2=0.1;
  eta=mu+bl*time+b2*trt+b3*time*trt+u1+u2*time+ltime;
  lam=exp(eta);
  model y~Poisson(lam);
  random u1 u2~Normal([0, 0],[sig1, cov, sig2]) subject=id;
run;

proc genmod data=new3;
  class id;
  model y= time trt time*trt / d=poisson offset=ltime;
  repeated subject=id / corrw covb type=exch;
run;

```

SAS output

The NLMIXED Procedure

Specifications

Data Set	WORK.NEW3
Dependent Variable	y
Distribution for Dependent Variable	Poisson
Random Effects	u1
Distribution for Random Effects	Normal
Subject Variable	id
Optimization Technique	Dual Quasi-Newton
Estimation Method	Adaptive Gaussian Quadrature

Dimensions

Observations Used	290
Observations Not Used	0
Total Observations	290
Subjects	58

Max Obs Per Subject	5
Parameters	5
Quadrature Points	20

Parameters					
mu	b1	b2	b3	sig1	NegLogLike
1	0	0	0	0.1	1015.04066

Iterations					
Iter	Calls	NegLogLike	Diff	MaxGrad	Slope
1	3	955.8817	59.15896	57.97191	-7021.37
2	5	952.178211	3.70349	55.02552	-34.165
3	6	948.800206	3.378005	13.73927	-29.1799
4	7	948.525761	0.274445	2.674937	-0.54631
5	9	948.511099	0.014661	1.814969	-0.00427
6	10	948.486503	0.024596	0.610316	-0.02313
7	12	948.483431	0.003072	0.070732	-0.00642
8	14	948.483277	0.000154	0.052537	-0.00008
9	16	948.483246	0.000031	0.002884	-0.00005
10	18	948.483246	3.612E-8	0.000061	-7.42E-8

The NLMIXED Procedure

NOTE: GCONV convergence criterion satisfied.

Fitting Information

-2 Log Likelihood	1897.0
AIC (smaller is better)	1907.0
BIC (smaller is better)	1917.3
Log Likelihood	-948.5
AIC (larger is better)	-953.5
BIC (larger is better)	-958.6

Parameter Estimates

Parameter	Estimate	Standard Error	DF	t Value	Pr > t	Alpha	Lower	Upper	Gradient
mu	1.0359	0.1415	57	7.32	<.0001	0.05	0.7526	1.3192	0.000061
b1	0.1108	0.04689	57	2.36	0.0216	0.05	0.01691	0.2047	-0.00002
b2	-0.01049	0.1968	57	-0.05	0.9577	0.05	-0.4047	0.3837	0.000044
b3	-0.3016	0.06975	57	-4.32	<.0001	0.05	-0.4413	-0.1619	0.000041
sig1	0.5167	0.1013	57	5.10	<.0001	0.05	0.3139	0.7196	-0.00001

The NLMIXED Procedure

Specifications

Data Set	WORK.NEW3
Dependent Variable	Y
Distribution for Dependent Variable	Poisson
Random Effects	u1 u2
Distribution for Random Effects	Normal
Subject Variable	id
Optimization Technique	Dual Quasi-Newton
Estimation Method	Adaptive Gaussian Quadrature

Dimensions

Observations Used	290
Observations Not Used	0
Total Observations	290
Subjects	58
Max Obs Per Subject	5
Parameters	7
Quadrature Points	20

Parameters

mu	b1	b2	b3	sig1	cov	sig2	NegLogLike
1	0	0	0	0.1	0.05	0.1	952.625769

Iterations

Iter	Calls	NegLogLike	Diff	MaxGrad	Slope
1	3	912.039756	40.58601	264.2311	-7563.74
2	4	898.775177	13.26458	18.03295	-979.652
3	5	896.722486	2.052691	13.29285	-4.27047
4	7	895.336636	1.38585	12.48232	-1.57163
5	8	894.589275	0.747361	10.88548	-2.00732
6	9	894.365239	0.224036	12.92942	-0.83035
7	11	893.837266	0.527973	8.323319	-1.0606
8	13	893.468515	0.36875	10.17836	-0.16147
9	15	893.346311	0.122204	10.68736	-0.15449
10	16	893.139226	0.207085	5.499622	-0.11613
11	18	893.053203	0.086023	0.895064	-0.18358
12	20	893.046781	0.006421	0.235226	-0.01462
13	22	893.045961	0.000821	0.13539	-0.00128

The NLMIXED Procedure

Iterations

Iter	Calls	NegLogLike	Diff	MaxGrad	Slope
14	24	893.045484	0.000477	0.121494	-0.00047
15	26	893.045332	0.000152	0.029424	-0.00026
16	28	893.04533	1.603E-6	0.00084	-3.17E-6

NOTE: GCONV convergence criterion satisfied.

Fitting Information

-2 Log Likelihood	1786.1
AIC (smaller is better)	1800.1
BIC (smaller is better)	1814.5
Log Likelihood	-893.0
AIC (larger is better)	-900.0
BIC (larger is better)	-907.3

Parameter Estimates

Parameter	Estimate	Standard Error	DF	t		Alpha	Lower	Upper	Gradient
				Value	Pr > t				
mu	1.0696	0.1343	56	7.96	<.0001	0.05	0.8005	1.3387	-0.0006
b1	0.005870	0.1070	56	0.05	0.9564	0.05	-0.2085	0.2202	0.000209
b2	-0.00970	0.1860	56	-0.05	0.9586	0.05	-0.3823	0.3629	-0.0005
b3	-0.3471	0.1489	56	-2.33	0.0233	0.05	-0.6453	-0.04888	-0.00024
sig1	0.4528	0.09354	56	4.84	<.0001	0.05	0.2654	0.6402	0.000059
cov	0.01725	0.05287	56	0.33	0.7455	0.05	-0.08867	0.1232	-0.00084
sig2	0.2161	0.05864	56	3.69	0.0005	0.05	0.09862	0.3336	-0.00047

The GENMOD Procedure

Model Information

Data Set	WORK.NEW3
Distribution	Poisson
Link Function	Log
Dependent Variable	y
Offset Variable	ltime
Observations Used	290

Class Level Information

Class	Levels	Values
id	58	101 102 103 104 106 107 108 110 111 112 113 114 116 117 118 121 122 123 124 126 128 129 130 135 137 139 141 143 145 147 201 202 203 204 205 206 208 209 210 211 213 214 215 217 218 219 220 221 222 225 226 227 228 230 232 234 236 238

Parameter Information

Parameter	Effect
Prm1	Intercept
Prm2	time
Prm3	trt
Prm4	time*trt

Criteria For Assessing Goodness Of Fit

Criterion	DF	Value	Value/DF
Deviance	286	2413.0245	8.4371
Scaled Deviance	286	2413.0245	8.4371
Pearson Chi-Square	286	3015.1555	10.5425
Scaled Pearson X2	286	3015.1555	10.5425
Log Likelihood		5631.7547	

Algorithm converged.

The GENMOD Procedure

Analysis Of Initial Parameter Estimates

Parameter	DF	Estimate	Standard Error	Wald 95% Confidence Limits		Chi-Square	Pr > ChiSq
				Lower	Upper		
Intercept	1	1.3476	0.0341	1.2809	1.4144	1565.44	<.0001
time	1	0.1108	0.0469	0.0189	0.2027	5.58	0.0181
trt	1	-0.1080	0.0486	-0.2034	-0.0127	4.93	0.0264
time*trt	1	-0.3016	0.0697	-0.4383	-0.1649	18.70	<.0001
Scale	0	1.0000	0.0000	1.0000	1.0000		

NOTE: The scale parameter was held fixed.

GEE Model Information

Correlation Structure	Exchangeable
Subject Effect	id (58 levels)
Number of Clusters	58
Correlation Matrix Dimension	5
Maximum Cluster Size	5
Minimum Cluster Size	5

Covariance Matrix (Model-Based)

	Prm1	Prm2	Prm3	Prm4
Prm1	0.01223	0.001520	-0.01223	-0.001520
Prm2	0.001520	0.01519	-0.001520	-0.01519
Prm3	-0.01223	-0.001520	0.02495	0.005427
Prm4	-0.001520	-0.01519	0.005427	0.03748

Covariance Matrix (Empirical)

	Prm1	Prm2	Prm3	Prm4
Prm1	0.02476	-0.001152	-0.02476	0.001152
Prm2	-0.001152	0.01348	0.001152	-0.01348
Prm3	-0.02476	0.001152	0.03751	-0.002999
Prm4	0.001152	-0.01348	-0.002999	0.02931

Algorithm converged.

The GENMOD Procedure

Working Correlation Matrix

	Col1	Col2	Col3	Col4	Col5
Row1	1.0000	0.5941	0.5941	0.5941	0.5941
Row2	0.5941	1.0000	0.5941	0.5941	0.5941
Row3	0.5941	0.5941	1.0000	0.5941	0.5941
Row4	0.5941	0.5941	0.5941	1.0000	0.5941
Row5	0.5941	0.5941	0.5941	0.5941	1.0000

Analysis Of GEE Parameter Estimates
Empirical Standard Error Estimates

Parameter	Estimate	Standard Error	95% Confidence Limits		Z	Pr > Z
			Lower	Upper		
Intercept	1.3476	0.1574	1.0392	1.6560	8.56	<.0001
time	0.1108	0.1161	-0.1168	0.3383	0.95	0.3399
trt	-0.1080	0.1937	-0.4876	0.2716	-0.56	0.5770
time*trt	-0.3016	0.1712	-0.6371	0.0339	-1.76	0.0781

Parameter estimates and SEs (in parenthesis) for the Progabide data.

Variable	Estimation Method		
	MLE ¹	PQL ²	GEE ³
Intercept	1.03 (0.14)	1.00 (0.14)	1.35 (0.16)
TRT	-0.01 (0.20)	-0.009 (0.19)	-0.11 (0.19)
TIME	0.11 (0.05)	0.11 (0.05)	0.11 (0.12)
TIME×TRT	-0.30 (0.07)	-0.30 (0.07)	-0.30 (0.17)
σ_s^2	$\hat{\sigma}_s^2=0.52$ (0.10)	$\hat{\sigma}_s^2=0.53$ (0.10)	$\hat{\rho}=0.60$

¹SAS Proc NLMIXED

²From Diggle, Liang, and Zeger (1994, p.188)

³SAS Proc GENMOD

Case Study 3: Potomac River Fever in Horses: (Atwill, et al, 1996) Potomac River Fever (equine monocytic ehrlichiosis) is a blood-borne rickettsial disease whose transmission mechanism is unknown. Both arthropod (e.g. blackfly) and direct oral transmission have been suspected but not verified. Identification of risk factors of horses in New York State might give clues to the spread of this disease and help with reducing its frequency.

511 farms were studied, each with several social groups of horses, for a total of 2,587 horses.

response: seropositive (yes/no) response for horse k in social group j at farm i .

distribution: Bernoulli *link:* logit

predictors: Frequency stall cleaned, Frequency fly spray applied, Breed, Sex, ...(fixed), Farm and Social group (farm) (random).

Questions: Transmission mechanism of Potomac River Fever?

Model: Y_{ijk} = infection for horse k in social group j on farm i .

$$Y_{ijk} \sim \text{Bernoulli}(p_{ijk})$$

$$\text{logit}(p_{ijk}) = \mu + s_{j(i)} + f_i + \text{fixed effects},$$

$$s_{j(i)} \sim N(0, \sigma_{\text{group}(farm)}^2)$$

$$f_i \sim N(0, \sigma_{farm}^2)$$

Analysis: Focus on the random factors. The estimated variances of the random effects were:

$$\hat{\sigma}_{farm}^2 = 1.26$$

$$\hat{\sigma}_{\text{group}(farm)}^2 = 0$$

So the difference in loglikelihood for testing $\sigma_{\text{group}(farm)}^2 = 0$ is zero and hence not statistically significant when compared to a $\frac{1}{2} \chi_1^2$.

Implications: There is a strong correlation among horses within a farm on the logit scale ($0.32 = \hat{\sigma}_{farm}^2 / (\hat{\sigma}_{farm}^2 + \sigma_{logistic}^2)$), but no correlation within social groups. This suggests the disease is not transmitted directly from horse to horse, but instead is related to environmental or management factors operating at a farm scale.

Case study 4: Chestnut Leaf Blight. Recall the situation: Viruses spread between fungal individuals when they come in contact and fuse together. A major obstacle in spreading this virus and thus controlling the disease is that different isolates of the fungus cannot necessarily transfer the virus to one another.

To estimate the effects of these genes, they have paired numerous isolates which differ on the first gene only, the second gene only, the first and the second gene, etc. For each combination of isolates they have averaged about 30 attempts and record a binary response of whether or not the attempt succeeded in transmitting the virus.

Questions of interest include whether pre-identified genes actually do have an influence on transmission of the virus (and if so, to what degree), whether there are other, as yet unidentified, genes which might affect transmission, and whether transmission is symmetric. By symmetry of transmission we mean the following: suppose the infected fungus is type b at the locus for the first gene and the non-infected isolate (which we are trying to infect) is type B. The two isolates are the same at the other five loci. Is the probability of transmission the same as when using a type B to try to infect a type b?

Model:

$Y_i = 1$ if virus is transmitted, 0 otherwise

$Y_i \sim \text{indep. Bernoulli}(p_i)$

$$p_i = \Phi(\mu + \sum_j \beta_j MCH_{ij} + \sum_j \gamma_j ASY_{ij}),$$

where $MCH_{ij} = 1$ if there is a mismatch at locus j for pairing i and 0 otherwise, and $ASY_{ij} = 1/2$ if there is a mismatch at locus j in pairing i with a b donor, $-1/2$ if there is a mismatch at locus j pairing i with a B donor and 0 if there is no mismatch.

β_j = effect of a mismatch on gene j .

γ_j = asymmetry effect.

= difference between a mismatch with a donor type b and type B .

Question: Is there asymmetric transmission?

maximized log likelihood of the model:

$$\log l = -955.303$$

with 13 parameters

maximized log likelihood of the model with all the γ_i set equal to zero:

$$\log l = -1116.639$$

with 7 parameters

Likelihood ratio test:

Difference is $1116.639 - 955.303 = 161.336$

$$\text{p-value} = P\{\chi_6^2 \geq 2*161.336\} \approx 0$$

Threshold model

A common model in genetics for describing the presence or absence of a trait is the threshold model. This arises from assuming that a large number of genes each have a small and additive effect and when the cumulative effect exceeds a threshold of zero the trait is present in an individual.

$Y = 1$ if trait is present, and 0 otherwise.

$\mathbf{x}'\boldsymbol{\beta}$ = either genetic or non-genetic fixed effects.

ε = the genetic effect not captured in $\mathbf{x}'\boldsymbol{\beta}$.

Appealing to the central limit theorem gives the probit model:

$$\begin{aligned} P\{Y=1\} &= P\{ \mathbf{x}'\boldsymbol{\beta} + \varepsilon > 0 \} \\ &= P\{ -\varepsilon < \mathbf{x}'\boldsymbol{\beta} \} = \Phi(\mathbf{x}'\boldsymbol{\beta}) \end{aligned}$$

Different isolates of the fungus are used which may differ with regard to genes other than the six pre-identified.

We might model their effects as being selected from a normal distribution.

Y_{ijk} = k th observation from an attempted infection from the i th isolate (the donor) to the j th isolate (the recipient).

\mathbf{x}_{ijk} = vector of covariates for Y_{ijk}

A reasonable model might then be:

$$P\{Y_{ijk} = 1 | \mathbf{u}\} = P\{ \mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j} + \varepsilon_{ijk} > 0 \},$$

where u_{1i} represent the (random) effects of the donor isolate and u_{2j} represent the (random) effects of the recipient isolate.

This gives

$$P\{Y_{ijk} = 1 | \mathbf{u}\} = \Phi(\mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j})$$

Consequences of introducing random factors

On the mean

$$E\left[Y_{ijk}\right] = E\left[E\left[Y_{ijk} \mid \mathbf{u}\right]\right] = \Phi(\mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j}),$$

or, using the threshold representation:

$$\begin{aligned} E\left[Y_{ijk}\right] &= E\left[P\{\mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j} + \varepsilon_{ijk} > 0 \mid \mathbf{u}\}\right] \\ &= P\{\mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j} + \varepsilon_{ijk} > 0\} \\ &= P\{-(u_{1i} + u_{2j} + \varepsilon_{ijk}) < \mathbf{x}'_{ijk} \boldsymbol{\beta}\} \\ &= P\{W < \mathbf{x}'_{ijk} \boldsymbol{\beta}\}, \end{aligned}$$

where $W \sim N(0, 1 + \sigma_1^2 + \sigma_2^2)$. So

$$\begin{aligned} E\left[Y_{ijk}\right] &= \Phi\left(\mathbf{x}'_{ijk} \boldsymbol{\beta} / \sqrt{1 + \sigma_1^2 + \sigma_2^2}\right) \\ &= \Phi\left(\mathbf{x}'_{ijk} \boldsymbol{\beta}^*\right) \end{aligned}$$

On the variance-covariance structure

For example, for two observations with the same donor and recipient isolate:

$$E\left[Y_{ijk} Y_{ijl} \right] = \int_{-\infty}^{+\infty} \Phi\left(\mathbf{x}'_{ijk} \boldsymbol{\beta} + \sigma z \right) \Phi\left(\mathbf{x}'_{ijl} \boldsymbol{\beta} + \sigma z \right) \exp(-z^2 / 2) / \sqrt{2\pi} dz,$$

where $\sigma = \sqrt{\sigma_1^2 + \sigma_2^2}$.

A correlation is therefore induced between responses which share one or more random effects.

Likelihood:

The conditional density of \mathbf{Y} given \mathbf{u} is

$$f_{\mathbf{Y}|\mathbf{u}} =$$

$$\prod \Phi(\mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j})^{y_{ijk}} [1 - \Phi(\mathbf{x}'_{ijk} \boldsymbol{\beta} + u_{1i} + u_{2j})]^{1 - y_{ijk}},$$

so that the likelihood is given by

$$L = \int \cdots \int f_{\mathbf{Y}|\mathbf{u}} f_{\mathbf{u}} d\mathbf{u},$$

which, for this example, is a 256-dimensional integral.

Question of interest: Are there other genes causing incompatibility?

If there are no other genes affecting the transmission of the virus, then all isolates with a given set of fixed effects will behave the same.

$$\Rightarrow H_0: \sigma_1^2 = 0, \sigma_2^2 = 0$$

Suppose we reject H_0 . How could we go about finding the genes that control incompatibility? We might look at the isolates that have the most extreme values of u_{1i} or u_{2j} .

Extreme values of u_{1i} or u_{2j} : Want predicted values of the u_{1i} and u_{2j} .

$$\text{best predictor} = \tilde{u}_{1i} = E[u_{1i} | \mathbf{Y}]$$

Two problems

- Depends on unknown parameters
- $E[u_{1i} | \mathbf{Y}] = \int u_{1i} f_{\mathbf{u} | \mathbf{Y}} d\mathbf{u}$ and

$$f_{\mathbf{u} | \mathbf{Y}} = f_{\mathbf{Y}, \mathbf{u}} / f_{\mathbf{Y}}$$

Research

Some selected topics in research.

1. Computing maximum likelihood estimates.

McCulloch (1994) - uses Gibbs sampler

McCulloch (1997) - uses Metropolis

Booth and Hobert (1999) – Indep. sampler

Geyer (1994) - Simulated ML

Geyer and Thompson (1992) - Simulated ML

Econometrics literature (Borsch-Supan and
Hajivassiliou, 1993)

Casella and Berger (1995) - Another method of
simulating to find ML estimates

Ruppert, et al (1984) - Stochastic approximation

2. PQL, Laplace approximations

Gilmour, Anderson and Rae (1984)

Schall (1991)

Breslow and Clayton (1994)

Breslow and Lin (1995)

Lin and Breslow (1996)

Wolfinger (1994)

3. Bayes estimates

Gilks, et al (1993)

Zeger and Karim (1991)

(But) Natarajan and McCulloch (1995)

4. GEEs

Zeger and Liang (1986)

Liang and Zeger (1986)

(But) Fitzmaurice (1995), Lipsitz, et al (1994)

5. Other

Engel and Keen (1994)

Kuk (1995)

McGilchrist (1994, 1995)

Heagerty and Lele (1998)

Drum and McCullagh (1993)

(1999)

SUMMARY

The Good News:

GLMMs can handle

- Non-normal data

- Nonlinear responses

- Random effects covariance structure

Can be used to

- Incorporate correlations in models

- Model the correlation structure

- Identify sensitive subjects

- Handle heterogeneous variances

Modelling process

1. Distribution of the data?

2. What is to be modelled?

3. Factors?

4. Fixed or random?

Software is available for linear and nonlinear normal models, some GLMs with normal random effects and for GEE estimation.

The not-so-Good News:

Computing methods for much of the class of GLMMs is an area of active research. Advances are being made in ML estimation, PQL, GEEs and Bayes methods.

General purpose software is still developing.

Tests and confidence intervals are asymptotic and approximate.

DNA Microarrays in Medicine

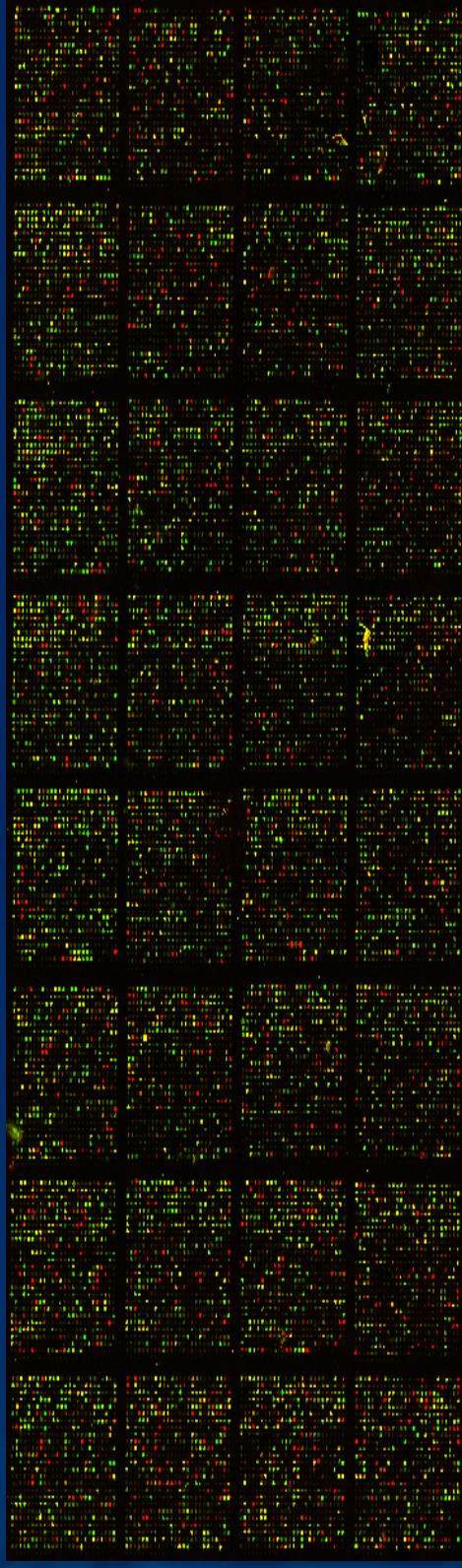
Expressing yourself one pixel at a time

J. Patrick Vandersluis, Ph.D.
HealthRx Corporation
www.HealthRx.com

Patrick@HealthRx.com

Agenda

- The Data of Biological Systems
- Review of Life Sciences Notions
- DNA Microarrays:
- The Bioinformatics Approach



The Data of Biological Systems

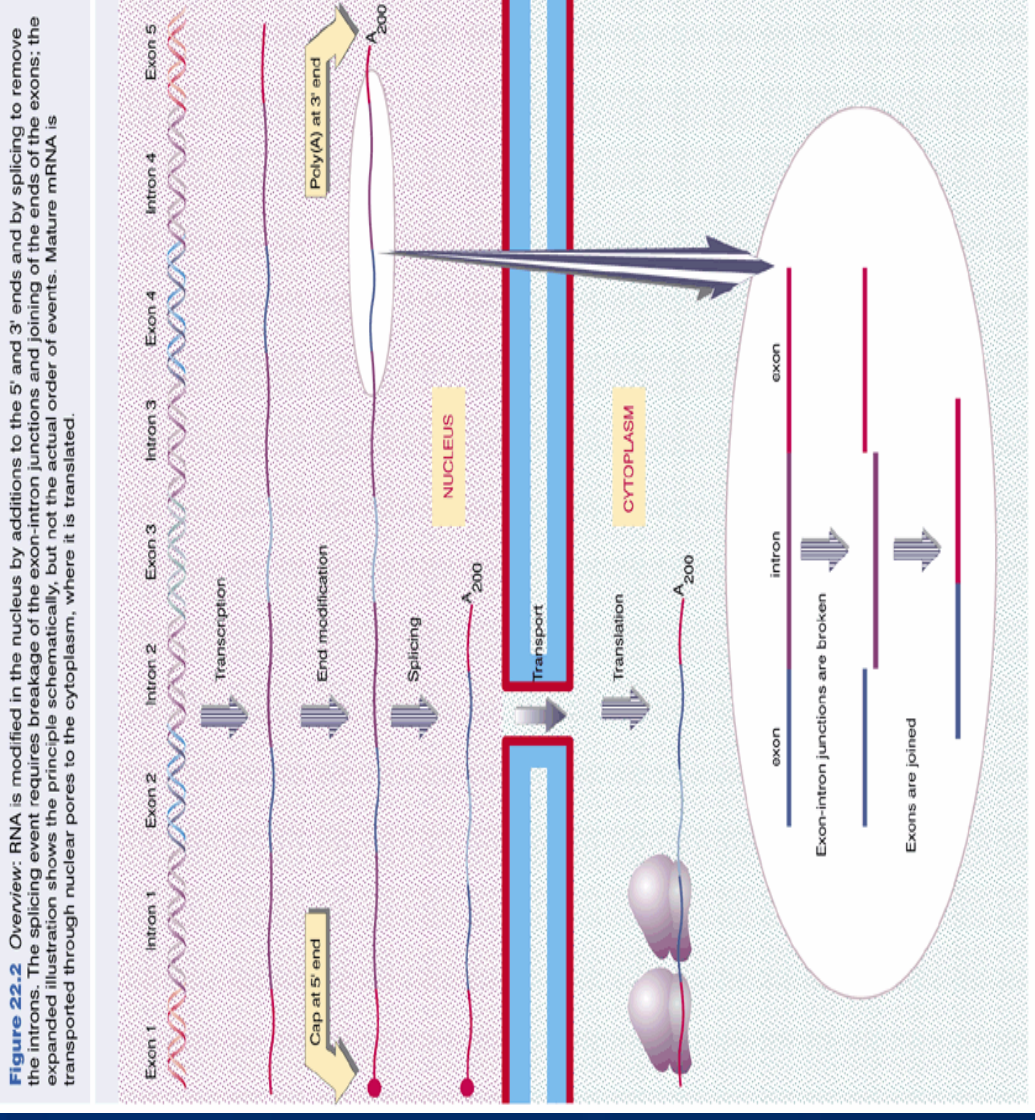
- Genomic-centric view of bioinformatics
 - Sequence data are signals
 - Proteins as Expression
 - Metabolic Pathways
- Biomedical signals
 - Signals data are sequences
 - ECG, EEG, EMG
- Both gene expression data and signals have significant temporal thread
- In combination, huge discoveries are waiting to be made

Review of Life Sciences Notions

- Sequences don't tell the whole story
- In same genome, expression is cell-type specific, tissue specific, and environment specific
- Disease, stress, metabolism cause variability of gene expression
- There is some “distance” between sequence and protein
- Just how do we get from genes to function?
- DNA -> mRNA -> Protein -> Function

RNA Processing

- Occurs in nucleus
- Pre-mRNA → mRNA
- Introns removed
- Mature mRNA transported through nuclear membrane
- Translated in cytoplasm
- Why look at mRNA instead of protein?



Protein, mRNA Profiles Differ

- Temporal differences between gene expression and protein accumulation
- Differential stability of mRNA and protein
 - Difference between rate of synthesis and amount of product
- Spatial differences, compartments, transport of mRNA and protein
- Differential processing of mRNA yielding various protein products
- Protein post-translational modifications

Protein Post-Translational Modifications: Examples

- Phosphorylation
- Glycosylation
- Proteolytic cleavage
- Disulfide bond formation
- Methylation
- Prenylation
- Adenylation
- Association with cofactors
- Etc.



One Gene - Multiple Proteins

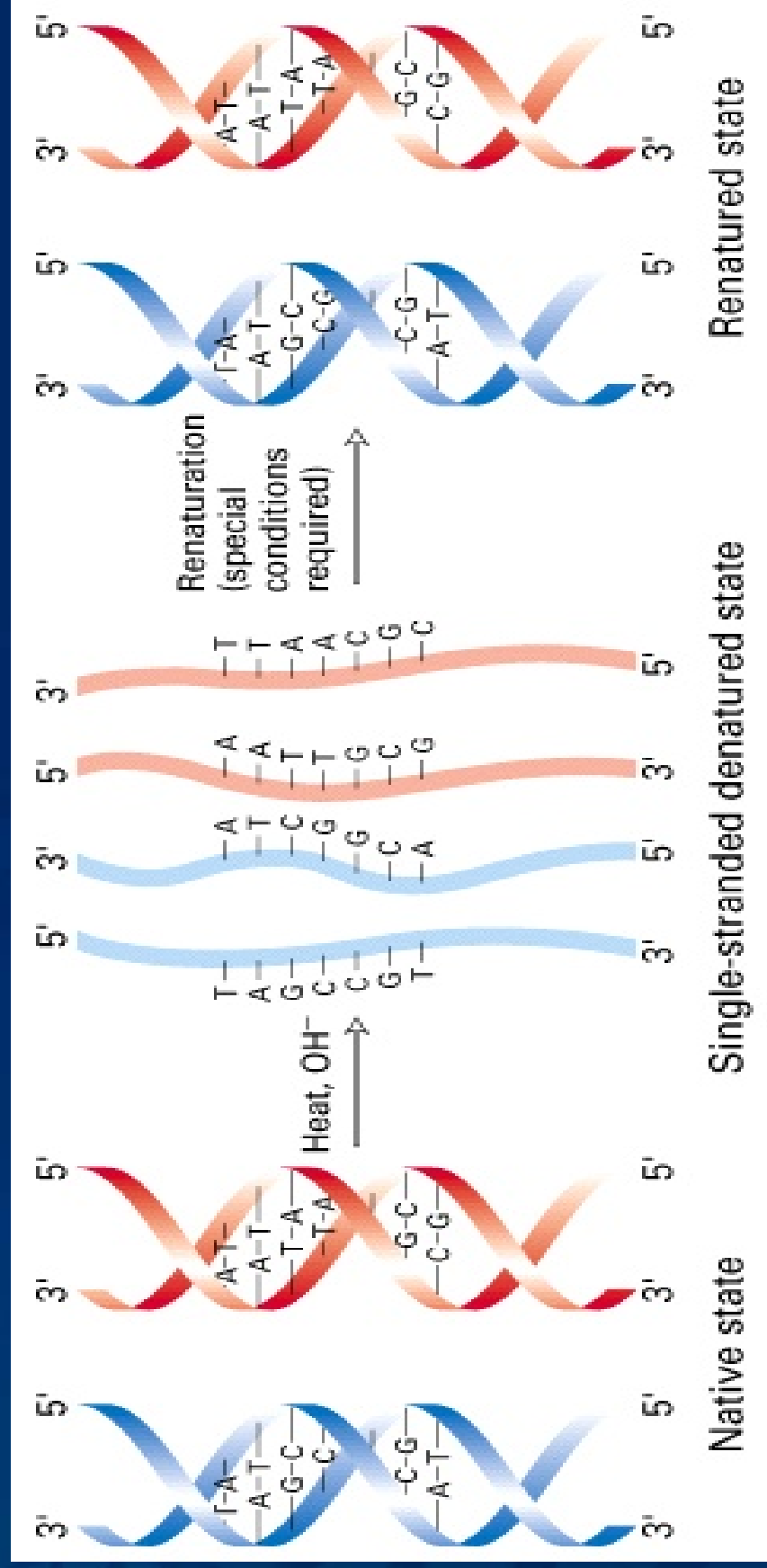
Insulin as Example

- Translated as preproinsulin, followed by membrane transport and processing.
- Cleavage yielding proinsulin and peptide C, then cleavage of proinsulin into 2 peptides.
- Ultimately three peptides: two (A, B) combine via disulfide bonds to form insulin.
- The third peptide C has unique biological activity of its own: cardiovascular effects.
- Summary: 1 gene, 3 polypeptides, 2 proteins, with 2 functions

What is the Problem to be Solved?

- Need strategy for understanding why certain proteins are produced under specific conditions (function)
- Understanding that mRNA is a measure of gene activity, need effective technique for measuring it through time
- Need tool to compare entire experimental transcriptome with a reference genome through time and varying environmental conditions to give greater understanding of complex gene interactions

Reminder - Nucleic Acid Hybridization is Specific



DNA Microarrays

- DNA Chips
- Massively parallel measurements
- Allow simultaneous measurement of the level of transcription for every gene over time
- Provide means for collecting data about differential gene expression over time in changing environment
- Robotic manufacture, standardized
- Automated readers, computer interpreted

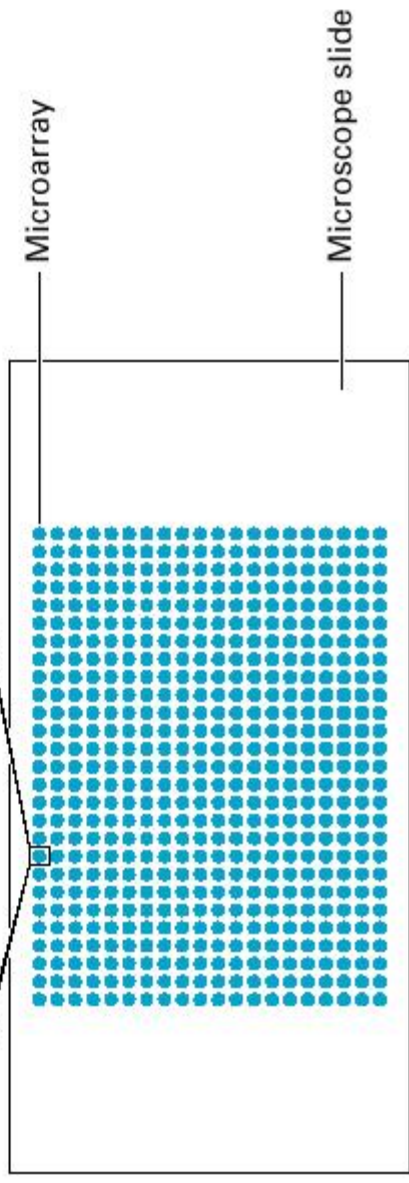
Making DNA Microarrays

- Known genes are amplified using PCR
- PCR products are verified, purified, and spotted onto ordinary glass slide
- Spot is about the size of a printed period
- Typical size is 6,200 spots, made robotically at about 12 slides per hour
- Once spotted, DNA is denatured, linked to slide with covalent bonds
- Stored at room temperature, very stable

GeneChip® in Schematic View

Sequence of one gene

```
TCC TTTCCGG AACGGTTGGC GTCTGCGCAC GCGGTGTGG GGCATGACAT
GCCGCCCCAG GAACAACCC C GACACGGCTT TAAGCCTCTC AAATCGCTGT
AGACATCATC TTTACGTGCT TGCCACCATT TGCCACCATT AGGGCTGTTC
CCGCGACGAC TCGCCATTCA ACCTCAGTCC TTCGGGTTGA GCGAGTGGGT
CGCGCGCAAG GTGCGAATGG GTCGCGCGCA AAGTGTGCG CTGGCTGTAT
TATATGCTGC CTATAGCGAG ACTAACGACC CACACTTCA CACAAGGATT
TCCCCTAAT GGTACCTCG CGTCAGGACC TTGACGCAAG CGCCCTTCG
GTTGGCCCA AGCTTGCTAG GACTACTTAT CTTGAGCTCA TTTAACATCC
CGCGCCCTCT CCGGGAGCGG TCGTCGCGAA GAAGTCAAAC CCGGAACGGC
GTTGACAAAG CGTGGAGACA TCGATACCTC TGTGT CAGCG GCCACAATC
```



GeneChip® Probe Array

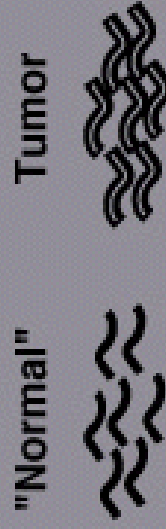


DNA Microarray Processing

- mRNAs harvested from two populations of cells
- For each population, mRNA copied to cDNA with reverse transcriptase, using nucleotides labeled with fluorescent dye
- Resulting cDNA mixed in equal amounts, microarray flooded and incubated overnight and allowed to hybridize
- Microarray washed, scanned with laser, and images captured
- Green and red scans are done separately, merged computationally. Yellow represents exactly 1:1 ratio of control/reference and experimental – an unusual situation
- Image processed to record color and intensity of spots

Labeled Cell cDNA Binds to EST or Gene DNA on Microarray

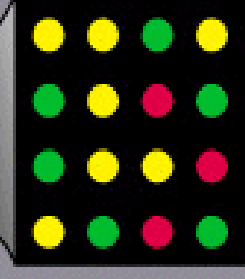
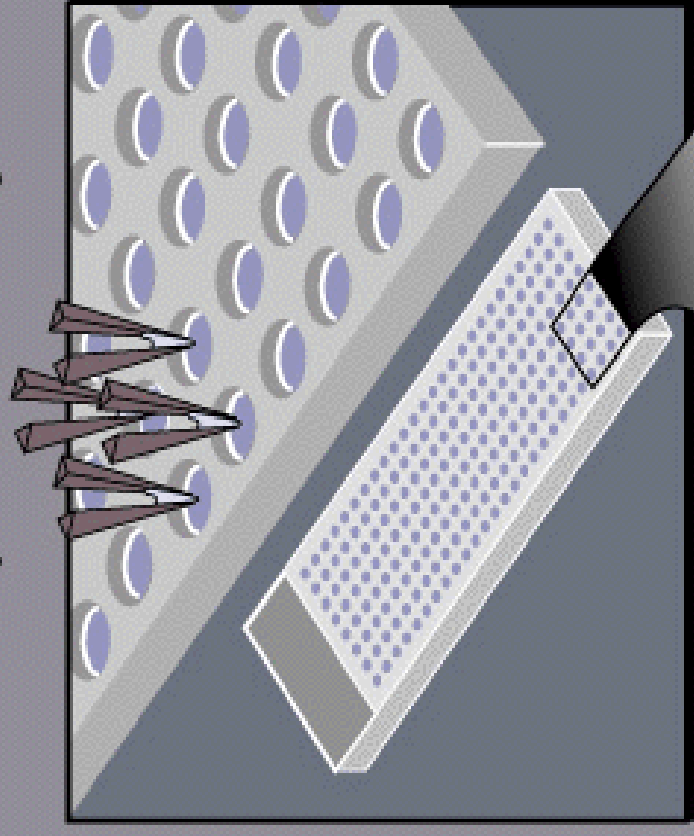
Prepare cDNA Probe



Combine Equal Amounts

Hybridize probe to microarray

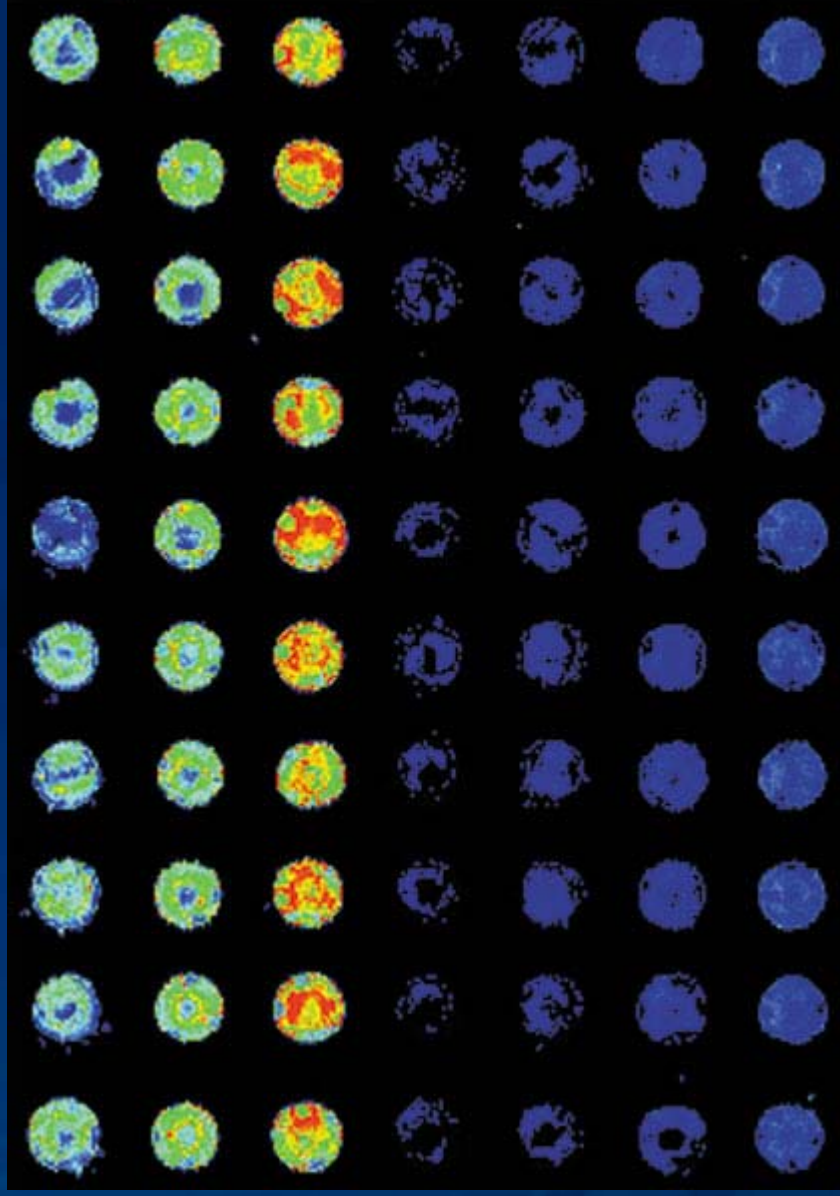
Prepare Microarray



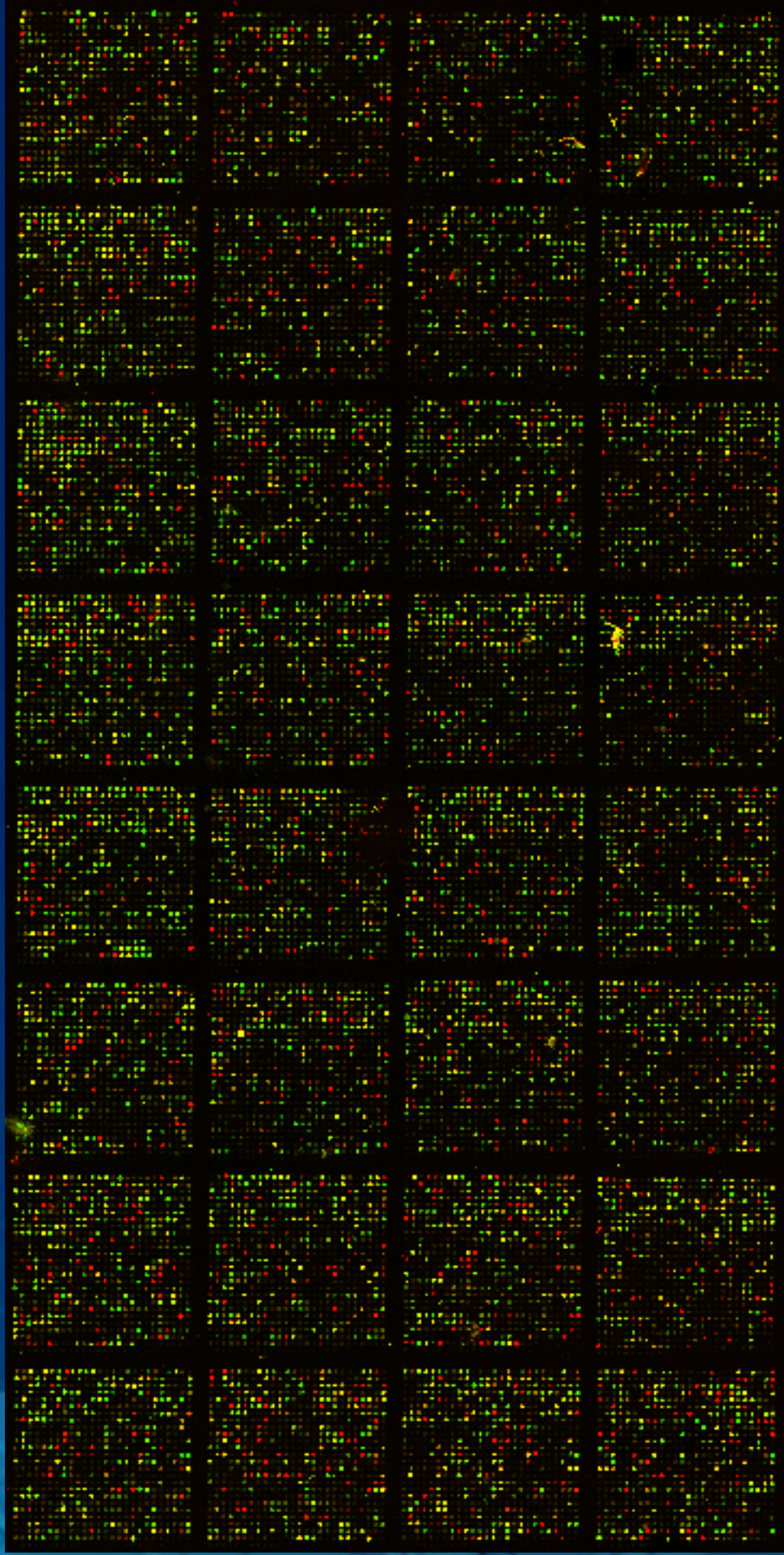
SCAN

Microarray Technology

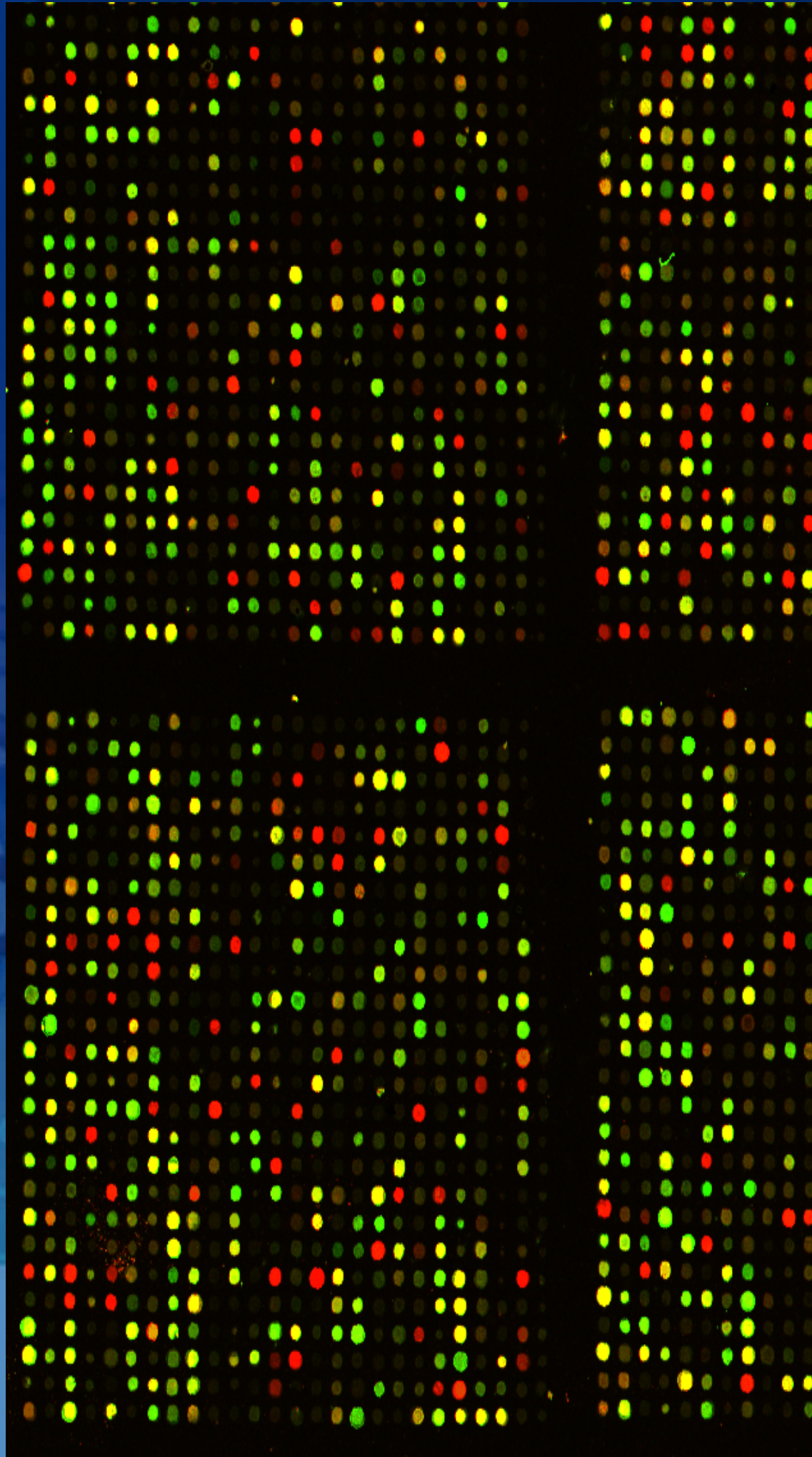
Sample of Affymetrix Arrayer Output



“All” 19,353 Genes of *C. elegans*

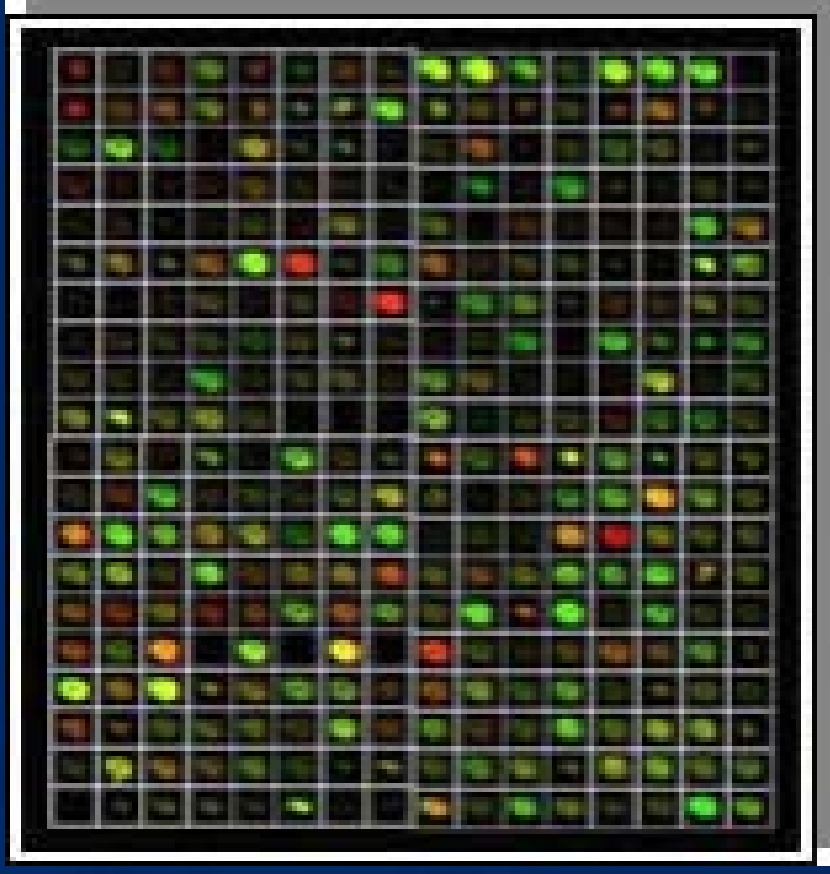


Zoomed View



Fluorescent Microarray Data

- Microarray raw data:
Color and intensity correlate with expression level
- Each spot represents hybridization of probe to a different defined DNA: EST or cDNA representing a gene
- Absolute intensity has no meaning, only relative intensity



cDNA Arrays

- In a population of cDNAs, each cDNA sequence and number of copies tell us which gene was expressed by how much, but exon regions only
- Compare with genomic (reference) DNA to find introns, promoters, intergenic DNA
- cDNAs represent the entire set of genes expressed in a particular cell or tissue type at a particular time
 - Compare normal vs. disease, tissue types, metabolic stress, differing environment
- Progress from structural genomics to functional genomics

Now What?

- We now know how to collect plenty of expression data
- We can create a large database with quantitative information
- How do we go from data to understanding?

The Bioinformatics Approach

- Deals with biological data as collections, in very large numbers
- Typically sequence-oriented data
- Algorithms seek to parallelize
- Computationally intensive
- Born of readily-available massive computing resources
- Data leads to hypothesis rather than reverse

Summary

- Understanding biological systems requires more than defining sequence data
- From sequence to function is a long jump
- Microarray data are fodder for tools that give insight into expression under varying conditions
- Combining microarray data with other clinical data along common temporal thread will provide insight into both normal and pathophysiology
- With your help, great discoveries will be made in this space!

SYSTEM TEST TIME BASED ON LINDSTROM/MADDEN APPROACH FOR CONTINUOUS DATA WITH WEIGHTED SUBSYSTEMS

Thomas R. (Tom) Walker
US Army Evaluation Center
4120 Susquehanna Avenue
APG, MD 21005-3013
DSN 458-0473, (410) 306-0473
FAX DSN 458-0476
tom.walker@usaec.army.mil

OBJECTIVES

- DESCRIBE THE METHOD
- SHOW AN EXAMPLE
- EXPLAIN SOME EXCEL ROUTINES
- SHOW HOW THIS METHOD WAS USED TO DETERMINE THE SAMPLE SIZE

DESCRIPTION

Method for computing the Lower Confidence Bound (LCB) on system reliability for independent discrete subsystems in series

LINDSTROM-METHOD FOR COMPUTING LCB'S

- **Step 1:** For each subsystem, calculate the Maximum Likelihood Estimate (MLE) of reliability, $\hat{\mathbf{R}}$, by dividing the number of successes by the number of trials, \mathbf{N}
- **Step 2:** Calculate the MLE of the system reliability by taking the product of all the MLEs of subsystem reliabilities
- **Step 3:** Obtain the equivalent number of trials (\mathbf{N}) for the system by choosing the minimum “number of trials” for each individual subsystem
- **Step 4:** Calculate the equivalent number of successes and failures for the system by multiplying \mathbf{N} by $\hat{\mathbf{R}}$ and \mathbf{N} by $(1 - \hat{\mathbf{R}})$
- **Step 5:** Obtain the $(1 - \alpha) * 100\%$ LCB for the series system as the $(1 - \alpha) * 100\%$ binomial LCB for a single component with F failures in \mathbf{N} tests

“EXAMPLES”
LINDSTROM/MADDEN METHOD
(DISCRETE DATA)

**OBJECTIVE IS TO DETERMINE THE SYSTEM RELIABILITY 90% LCB
WHEN THERE ARE TWO
OR MORE SUBSYSTEMS IN SERIES**

EXAMPLE 1	EXAMPLE 2
$R1_{HAT} = 20 / 20 = 1.00$ $R2_{HAT} = 45 / 50 = 0.90$	$R1_{HAT} = 9 / 10 = 0.90$ $R2_{HAT} = 40 / 50 = 0.80$ $R3_{HAT} = 5 / 5 = 1.00$
$R_{SYSTEM} = 1.00 * 0.90 = 0.90$	$R_{SYSTEM} = 0.90 * 0.80 * 1.00 = 0.72$
SYSTEM N $MIN(20, 50) = 20$	SYSTEM N $MIN(10, 50, 5) = 5$
$SYSTEM S = 20 * 0.90 = 18.0$	$SYSTEM S = 5 * 0.72 = 3.6$
90% LCB (18, 20) FOR THE SYSTEM = 0.75523	90% LCB (3.6, 5) FOR THE SYSTEM = 0.34394

BACKGROUND INFORMATION

EXCEL ROUTINE FOR DETERMINING BINOMIAL LOWER AND UPPER CONFIDENCE BOUNDS

EXAMPLES:

EXAMPLES	NUMBER OF SUCCESSES	NUMBER OF TRIALS	90% LOWER CONF BOUND	90% UPPER CONF BOUND
#1	18	20	0.75523	0.97309
#2	3.6	5	0.34394	0.94921

LOWER CONFIDENCE BOUND

= betainv (1 – conf level, # of successes, # of failures + 1)

= betainv (.1, 18, 3) = **0.75523 (Example #1)**

= betainv (.1, 3.6, 2.4) = **0.34394 (Example #2)**

UPPER CONFIDENCE BOUND

= betainv (conf level, # of successes +1, # of failures)

= betainv (.9, 19, 2) = **0.97309 (Example #1)**

= betainv (.9, 4.6, 1.4) = **0.94921 (Example #2)**

“EXAMPLE”

LINDSTROM/MADDEN APPROACH
(CONTINUOUS DATA)

OBJECTIVE IS TO DETERMINE THE SYSTEM RELIABILITY
80% LCB WHEN THERE ARE TWO OR MORE SUBSYSTEMS
IN SERIES

Lindstrom-Madden Method									
(Fixed Configurations)									
			$\alpha = 0.2$	yellow = input					
			$(1 - \alpha)$	LCL = 8.09					
			Sys Point Estimate = 10.00						
			Number of Subsystems = 4						
n sys =	λ sys * min =	20.00							
λ sys =	$\Sigma w \lambda =$	0.1			Min (Test Times / weights) =	200			
Utilization	Factor	λ hat	Weight	times	Failures	Test Time	Test Time	Weight	divided by
1.0000		0.05000	0.05000		20	400	400	400	400
0.5000		0.04000	0.02000		4	100	100	200	200
0.5000		0.02000	0.01000		4	200	200	400	400
0.2000		0.10000	0.02000		5	50	50	250	250

SYSTEM POINT ESTIMATE = 200 / 20 = 10.00

80% LCB = (2 * 200) / CHINV (.2, 42) = 400 / 49.45596 = 8.09

ORIGINAL TEST REQUIREMENT

DEMONSTRATE SYSTEM REQUIREMENT OF 66 HOURS AND
SUBSYSTEM REQUIREMENTS WITH 80% CONFIDENCE BASED
ON THE FOLLOWING ASSUMED USAGE RATES

SUBSYSTEM	ASSUMED USAGE RATE	MTBF REQUIREMENT (HOURS)
A	1.00 (39/39)	170
B	.31 (12/39)	100
C	.15 (6/39)	60
D	.18 (7/39)	50

METHOD USED TO DETERMINE THE NUMBER OF SYSTEM TEST HOURS BASED ON THE ORIGINAL TEST REQUIREMENT

Requirement	Weight	Allowable # of Failures	Subsystem Test Time Required to Demonstrate 80% LCB	System Test Time Required to Demonstrate Subsystem 80% LCB
170	1.0000	2	727	727
100	0.3077	2	428	1391
60	0.1538	2	257	1669
50	0.1795	2	214	1192

The minimum system test time required to demonstrate each subsystem requirement with at least 80% confidence is **1669** hours

$$\begin{aligned}
 \text{Requirement} &= (2 * \text{Subsystem Test Time}) / (\text{CHINV} (.2, 6)) \\
 \text{Subsystem Test Time} &= (\text{CHINV} (.2, 6) * \text{Requirement}) / 2 \\
 &= (8.558058 * 100) / 2 \\
 &= 428
 \end{aligned}$$

USING A TEST TIME OF 1669 HOURS, THE MTBF REQUIREMENTS FOR THE SYSTEM AND SUBSYSTEMS WILL BE MET WITH AT LEAST 80% CONFIDENCE WHEN THE NUMBER OF FAILURES IS LESS THAN OR EQUAL TO 2

Lindstrom-Madden Method (Fixed Configurations)		
$\alpha =$	0.2	
$(1 - \alpha)$ LCL =	146.66	
Sys Point Estimate =	208.63	
Number of Subsystems =	4	
$n_{sys} =$	$\lambda_{sys} * min =$ 8.00	
$\lambda_{sys} =$	$\sum w\lambda =$ 0.0048	
Utilization	Min (Test Times / weights) = 1669	
Factor	Weight	
(Weight)	Test Time	
	Weight	
	Requirement	
1.0000	1669	170
0.3077	514	100
0.1538	257	60
0.1795	300	50

IF THE SUBSYSTEM TEST TIMES ARE IN ACCORDANCE WITH THEIR UTILIZATION FACTORS (WEIGHTS), THE LINDSTROM MADDEN METHOD IS NOT NECESSARY FOR DETERMINING THE LCB.

REVISION TO THE ORIGINAL TEST REQUIREMENT FOR THE SUBSYSTEMS

- **DEMONSTRATE EACH SUBSYSTEM REQUIREMENT AS A POINT ESTIMATE WHEN THE TEST TIME IS AT LEAST TWICE ITS REQUIREMENT (I.E., THERE ARE NO MORE THAN TWO (2) FAILURES FOR EACH SUBSYSTEM)**

AND

- **DEMONSTRATE THE SYSTEM REQUIREMENT OF 66 HOURS WITH AT LEAST 80% CONFIDENCE.**

METHOD USED TO DETERMINE THE NUMBER OF SYSTEM TEST HOURS BASED ON THE REVISION TO THE ORIGINAL TEST REQUIREMENT

Subsystem	Requirement	Weight	Allowable Number of Failures	Subsystem Test Time Required to Demonstrate the Point Estimate	System Test Time Required to Demonstrate the Point Estimate
A	170	1.0000	2	340	340
B	100	0.3077	2	200	650
C	60	0.1538	2	120	780
D	50	0.1795	2	100	557

The minimum system test time required to demonstrate each subsystem’s requirement as a point estimate would be equal to **780** hours.

USING A SYSTEM TEST TIME OF 780 HOURS AND NO MORE THAN TWO FAILURES FOR EACH SUBSYSTEM, THE MTBF POINT ESTIMATE FOR EACH SUBSYSTEM IS MET AND THE SYSTEM MTBF REQUIREMENT (66 HOURS) IS MET WITH AT LEAST 80% CONFIDENCE.

Linstrom-Madden Method									
(Fixed Configurations)									
	$\alpha =$	0.2							
	$(1 - \alpha)$ LCL =	68.54							
	Sys Point Estimate =	97.50							
	Number of Subsystems =	4							
	$n_{sys} = \lambda_{sys} * min =$	8.00							
$\lambda_{sys} =$	$\Sigma w\lambda =$	0.0103	Min (Test Times / weights) =	780					
Utilization	Weight			Test Time					
Factor	times			divided by					
(Weight)	λ_{hat}	λ_{hat}	Failures	Test Time	Weight	Requirement			
1.0000	0.00256	0.00256	2	780	780	170			
0.3077	0.00833	0.00256	2	240	780	100			
0.1538	0.01667	0.00256	2	120	780	60			
0.1795	0.01429	0.00256	2	140	780	50			

CONCLUSIONS

- **TOTAL SYSTEM TEST TIME REQUIRED IS 780 HOURS**
- **WITH NO MORE THAN 2 FAILURES FOR EACH SUBSYSTEM, THE 80% LCB FOR THE SYSTEM MTBF WILL BE GREATER THAN ITS REQUIREMENT (66 HOURS) AND THE POINT ESTIMATE FOR EACH SUBSYSTEM WILL MEET ITS REQUIREMENT.**
- **THE LINDSTROM/MADDEN APPROACH CAN BE USED TO DETERMINE SYSTEM TEST TIME WHEN REQUIRED TO DEMONSTRATE THE SYSTEM REQUIREMENT WITH CONFIDENCE AND THE SUBSYSTEM REQUIREMENTS WITH CONFIDENCE OR AS A POINT ESTIMATE.**

BoM

System Level Analysis vs Subsystem Level Analysis
When Testing According to OMS/MP

No. Subsystems = k
 Total System Test Time = T
 Subsystem Utilization Factor (Weight) = w_i
 Subsystem Test Time = $w_i T = T_i$
 Subsystem # failures = f_i
 System # failures = $\sum_{i=1}^k f_i = f$

System Level Analysis: $\hat{MTBF}_{SYS} = \frac{T}{f}$

$Y\% \text{ LCB on MTBF} = \frac{2T}{\chi^2_{r, 2f+2}}$

Subsystem Level Analysis

Subsystem failure rate MLE = $\hat{\lambda}_i = \frac{f_i}{T_i}$
~~System~~ failure rate MLE = $\hat{\lambda}_{SYS} = \sum w_i \hat{\lambda}_i$

$\hat{MTBF}_{SYS} = \frac{1}{\hat{\lambda}_{SYS}} = \frac{1}{\sum_{i=1}^k w_i \hat{\lambda}_i} = \frac{1}{\sum_{i=1}^k w_i (\frac{f_i}{T_i})} = \frac{T}{f}$ (Same as one System Level Analysis)

Y% LCB on MTBF (LINDSTROM/MADDEN APPROACH)

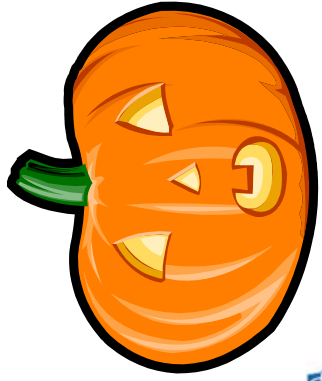
Equivalent System Time = $\min(\frac{T_i}{w_i}) = \min(\frac{w_i T}{w_i}) = T$
 Equivalent # System Failure = $\hat{\lambda}_{SYS} T = \sum_{i=1}^k w_i \hat{\lambda}_i T = \sum_{i=1}^k \frac{T_i}{T} \frac{f_i}{T_i} T = \sum_{i=1}^k f_i = f$

$\therefore \hat{MTBF}_{SYS} \text{ Y\% LCB} = \frac{2T}{\chi^2_{r, 2f+2}}$ (Same as one System Level Analysis)

INFORMATION INTEGRATION FOR STOCKPILE SURVEILLANCE

Alyson Wilson
Statistical Sciences Group
Los Alamos National Laboratory
agw@lanl.gov

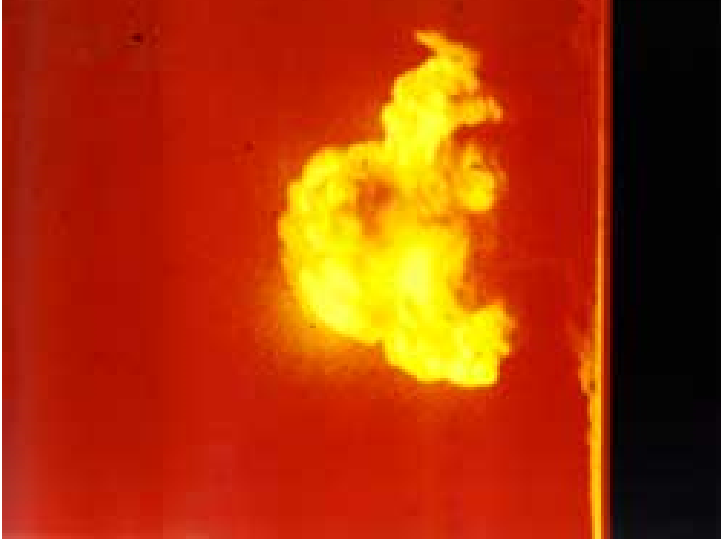
October 31, 2003



Acknowledgments

- Laura McNamara, Sandia National Laboratories
- Kristi O’Grady, graduate student at North Carolina State University and LANL summer student
- Shane Reese, Brigham Young University and LANL visiting faculty
- Todd Graves, Nick Hengartner, and Andrew Koehler, LANL
- Greg Anderson, Kevin Deal, and George Lopez, MCPD

Science-Based Stockpile Stewardship



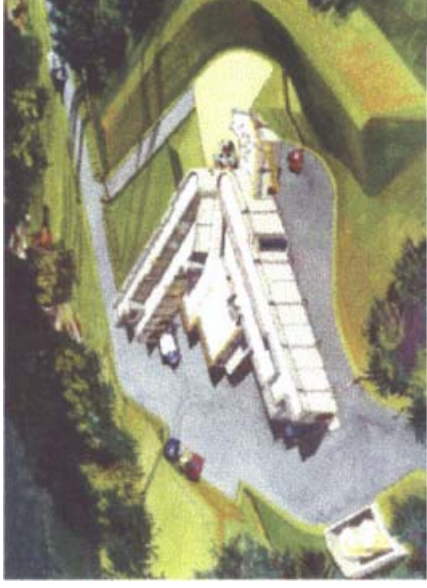
In the absence of full-system testing, how do we understand the stockpile and integrate various sources of information to get a quantitative estimate, with uncertainties, of system reliability and performance?

Surveillance

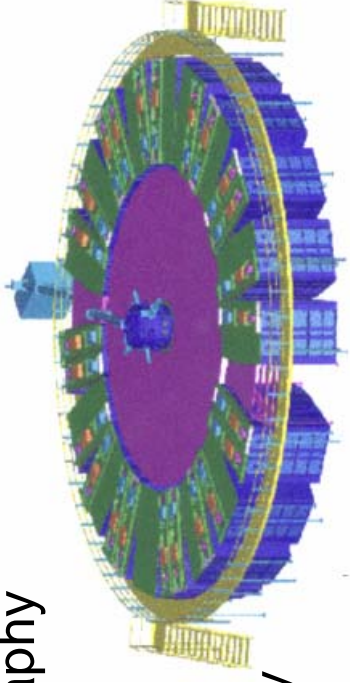
Continuous monitoring of “X”
to ensure the health of “Y”

- Detect and respond through:
- Planned “**data**” collection
 - Simulation, experimental, field, database, text, images, expert judgment, ...
 - Maintenance
 - Life extension programs, special investigations, ...

Science-Based Stockpile Stewardship



- Large Scale Computing
- Advanced Radiography
- Materials Science
 - Pu
 - High explosives
- High-energy density experiments
- Advanced manufacturing
- Information integration



Outline

- These two problems started off feeling like a reliability assessment or a PRA, but ended up somewhere rather different.
- Model development, Bayesian network
- Can we do better than “x/n”?
- Both have relevance back to LANL stockpile surveillance

Example 1: Missile Defense Agency

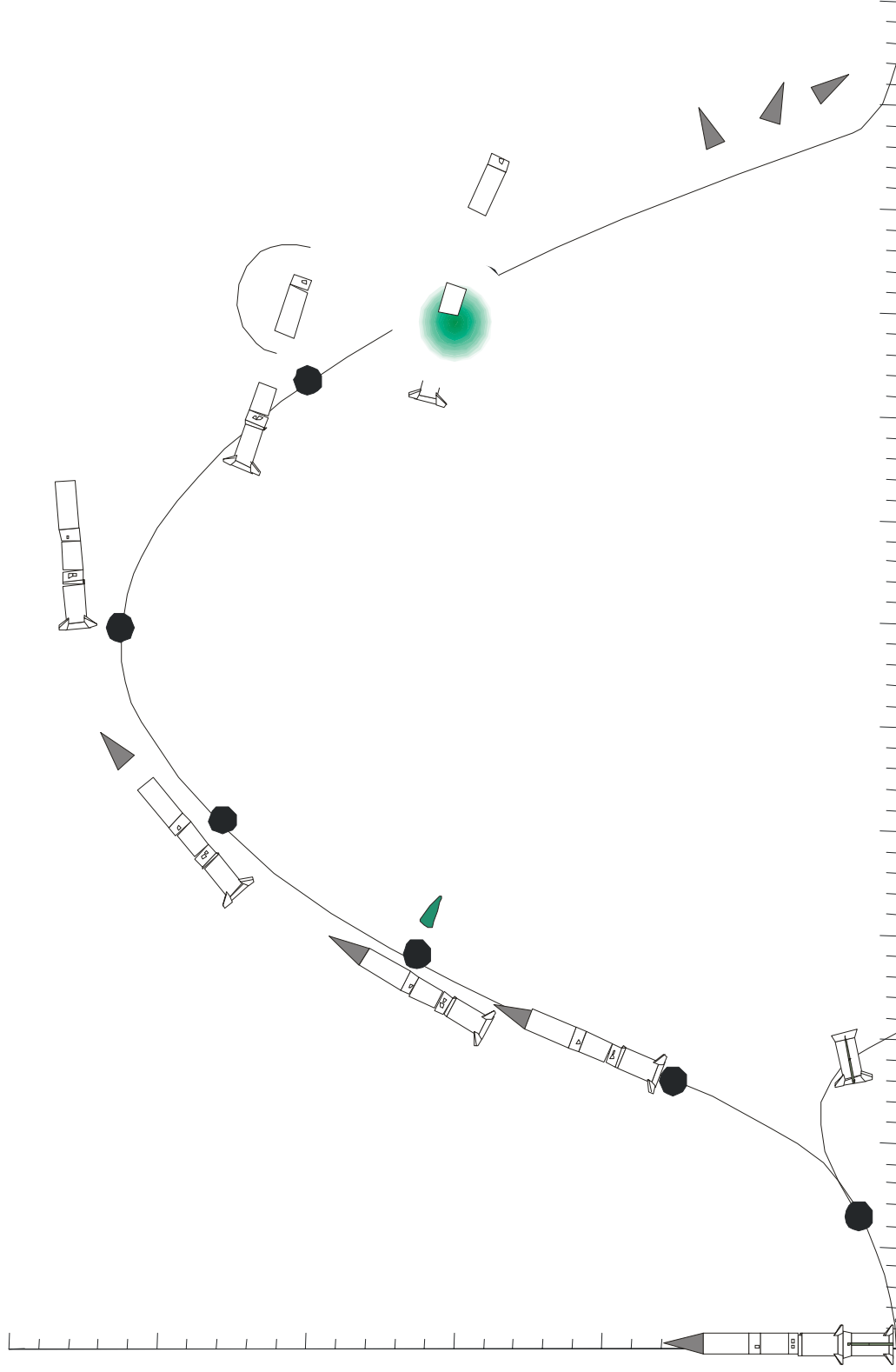
PROGRAM: Fly a high-fidelity, threat-representative missile system for Theater Missile Defense data collection and interoperability exercise

GOAL: “Quantify the probability of mission success” and identify “areas of unacceptable risk” to the program

ISSUES:

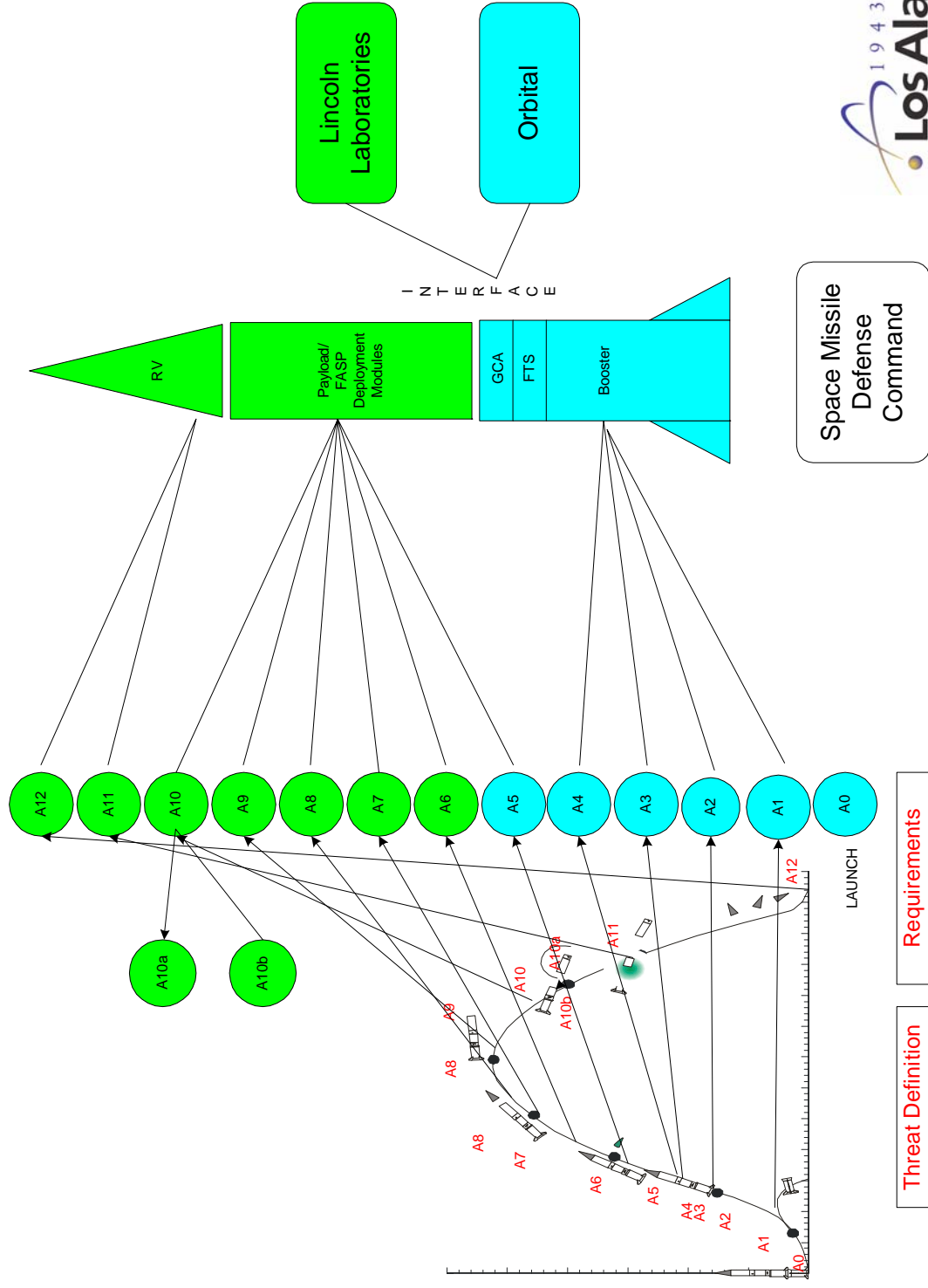
- Multiple partners and contractors
- High reliability demanded
- Full system testing not an option
- System requirements dynamic
- Diverse data sources

Notional Trajectory



Events to System

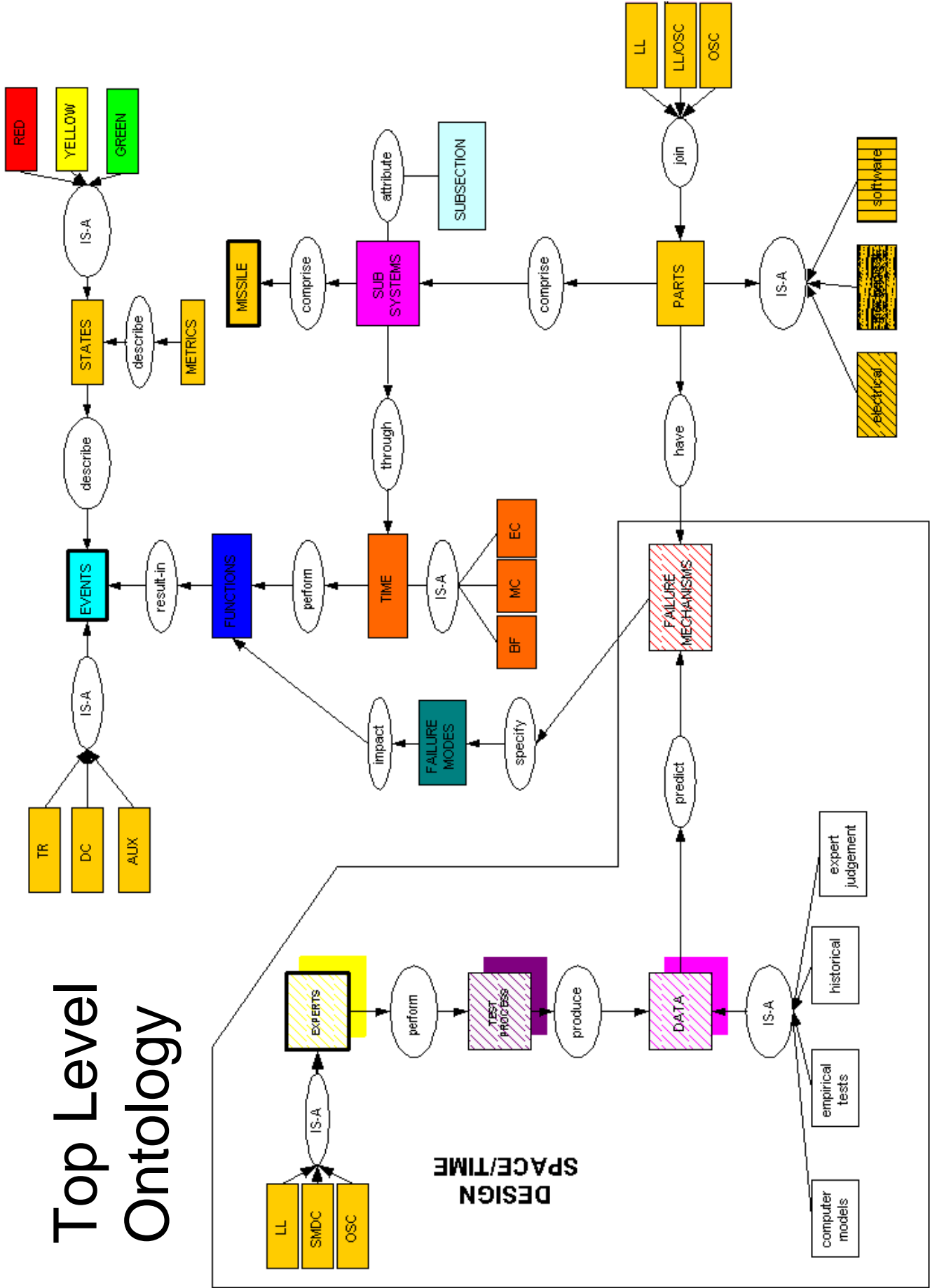
EVENTS IN SCENARIO



Threat Definition Requirements

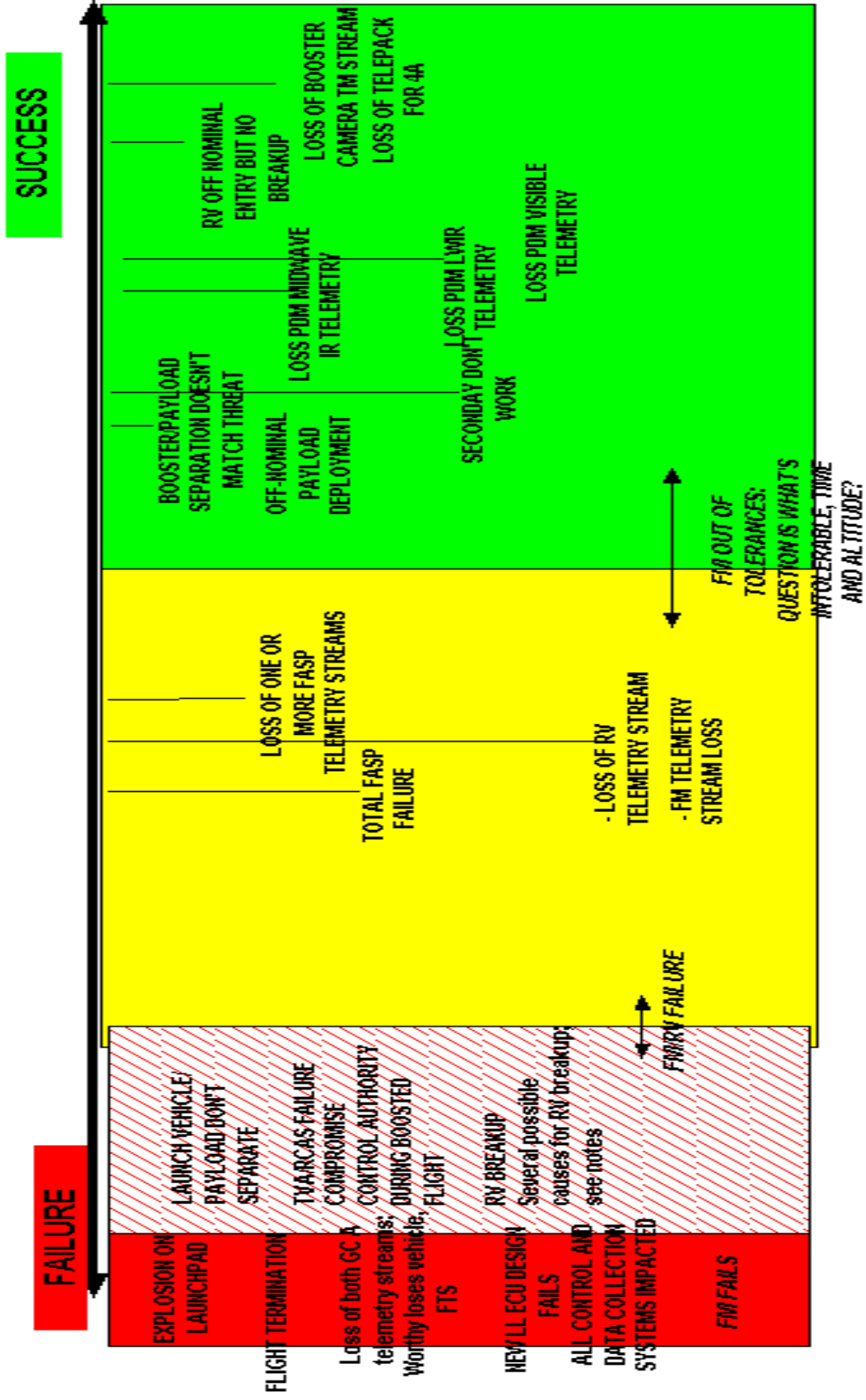
Space Missile Defense Command

RUN SPACE/TIME

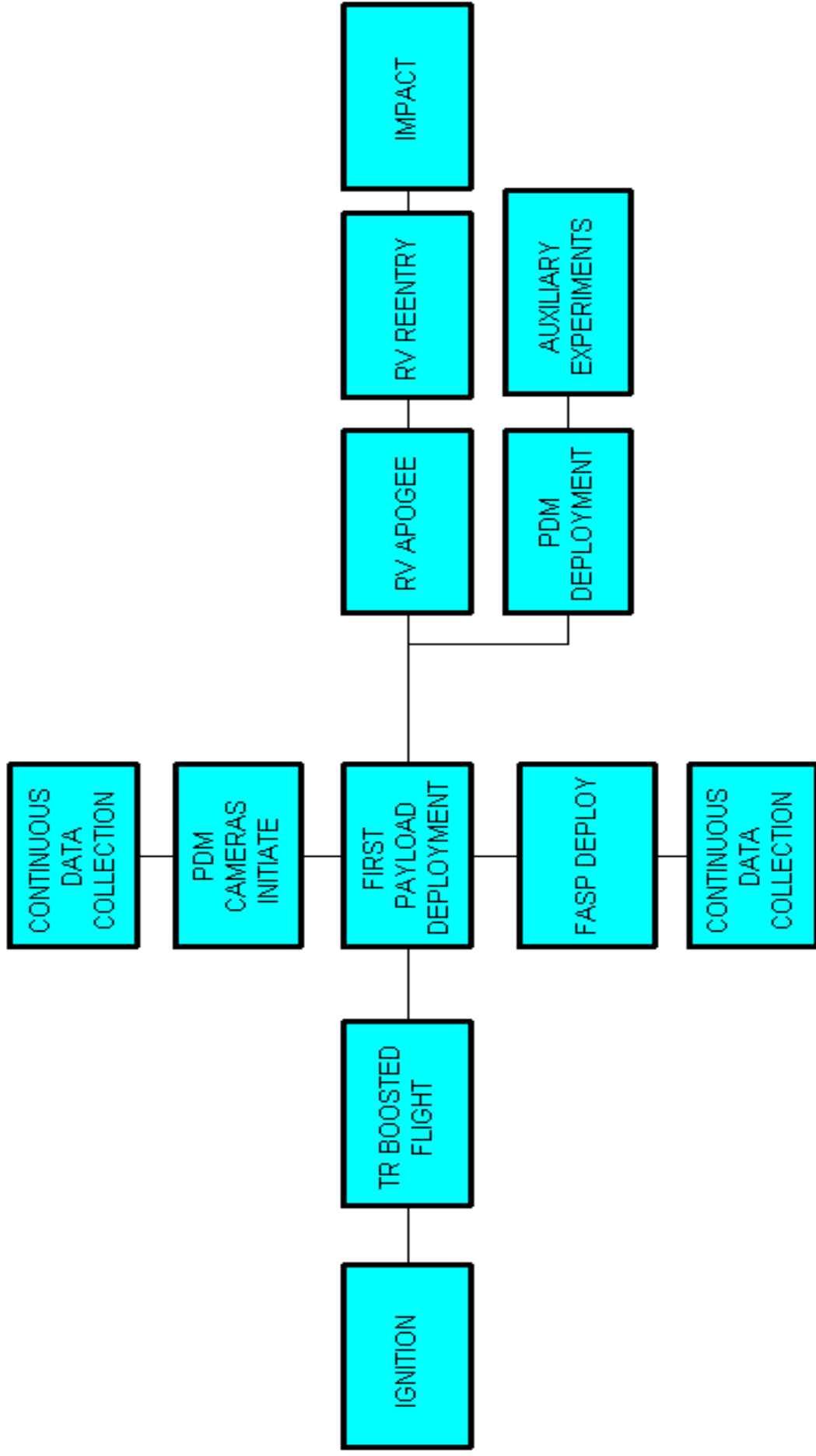


Top Level Ontology

Mission Success



Event Diagram



TR BOOSTED FLIGHT/ TRAJECTORY

TR FLIGHT

VEHICLE GUIDANCE, NAVIGATION, CONTROL

ATTRIBUTE CONTROL

VEHICLE STABILITY

THERMAL PROTECTION

ENV. PROTECTION

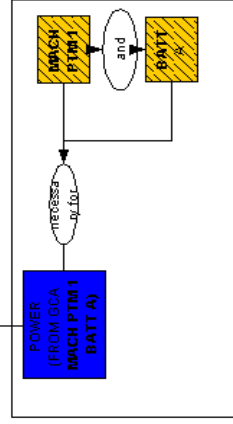
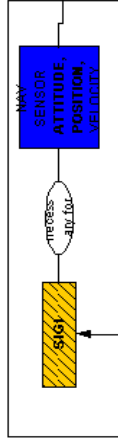
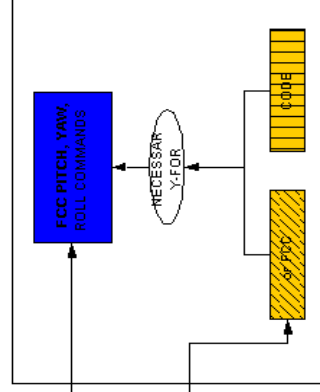
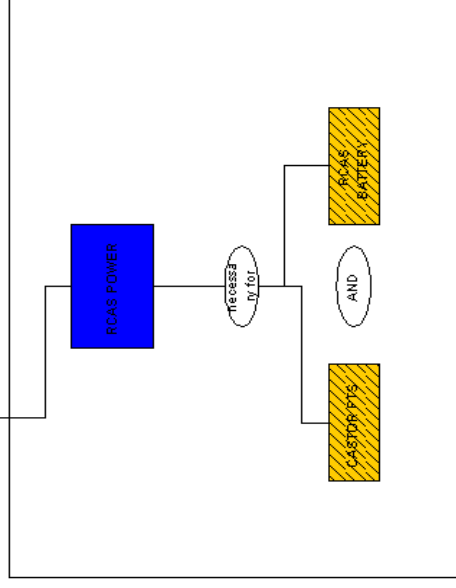
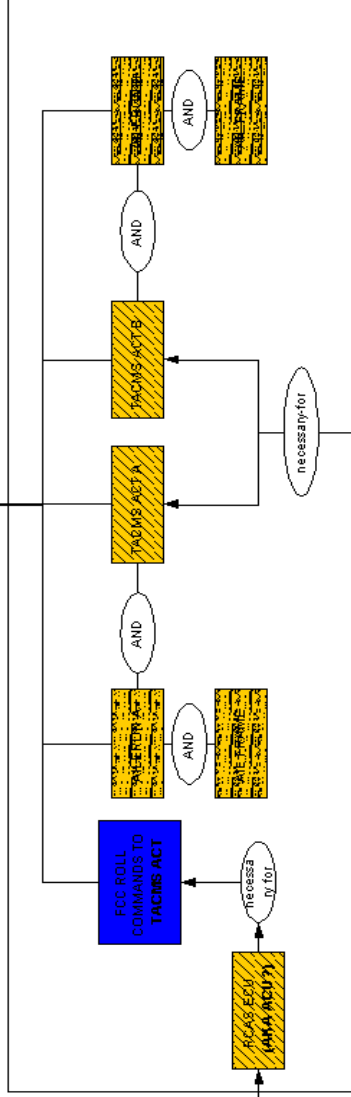
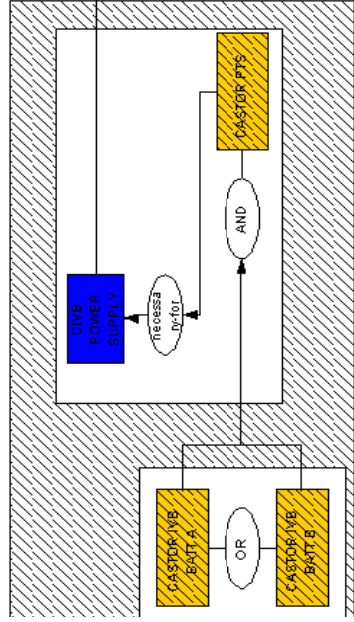
SF ROLL CONTROL

Event

Dependency

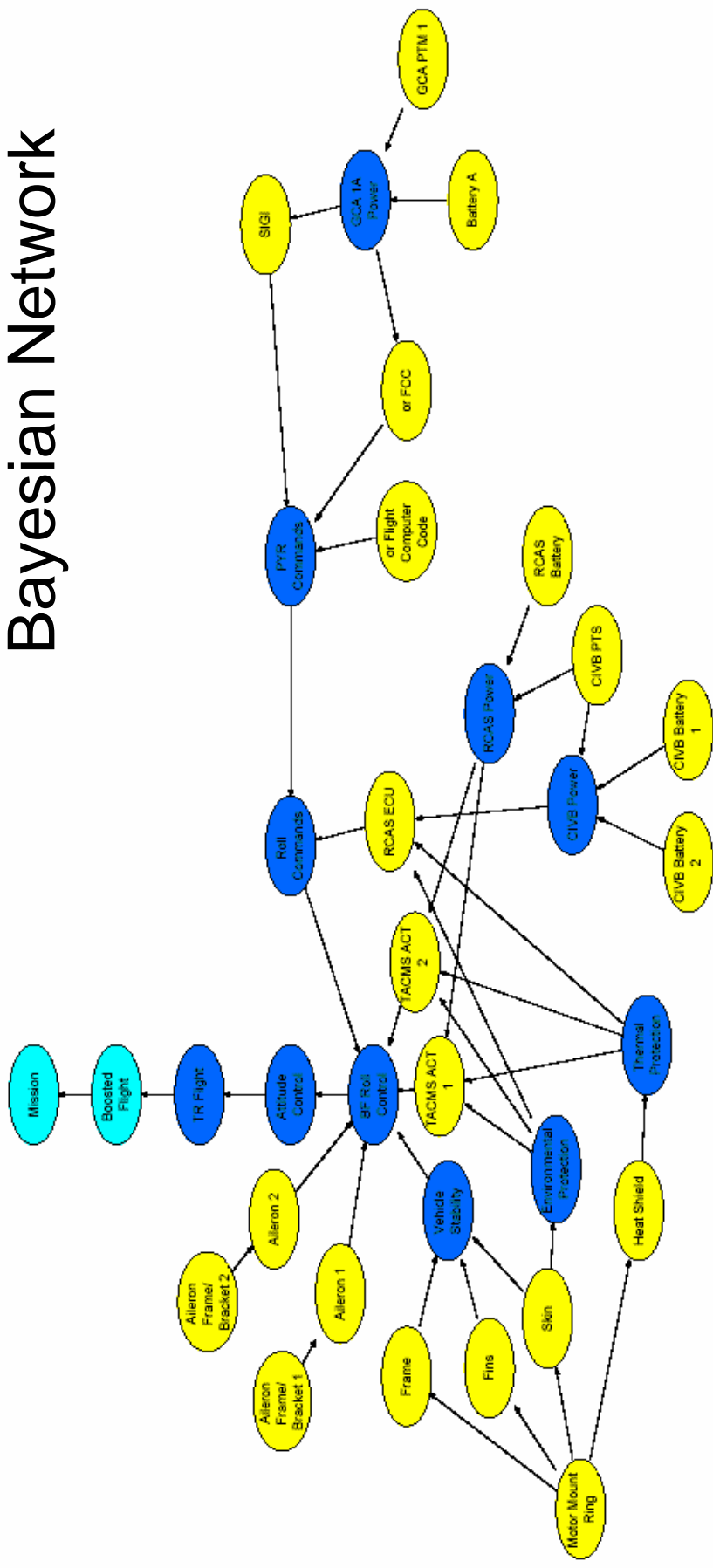
Diagram

VEHICLE STABILITY



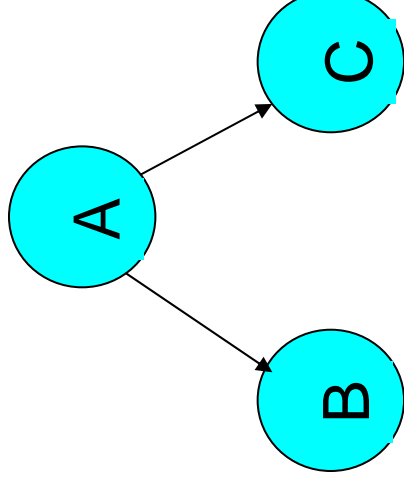
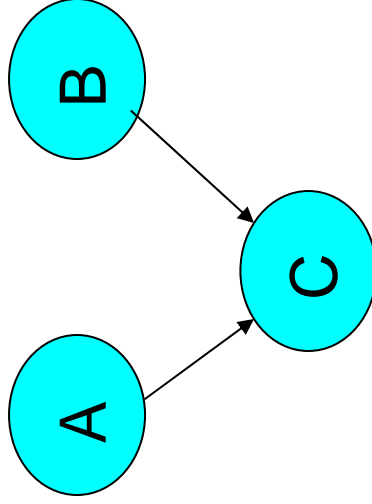
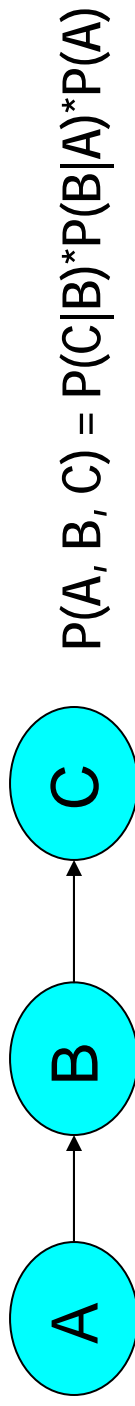
Statistical Model Representation

Boosted Flight Roll Control Bayesian Network



Bayesian Network Calculation

- Local conditional structure (like the elicited data)
- $P(A_1, \dots, A_{599}) = \prod P(A_i | \text{parents})$
- Three structures: serial, converging, diverging



Data

Engineering Judgment

- The probability of the motor mount ring failing catastrophically is under 1%.
- If the motor mount ring fails catastrophically, then the fins and frame fall off the vehicle.
- There is somewhere between a 5% and 10% chance that the skin will peel back.
- If the fins or frame are missing, then the vehicle is unstable.
- If the skin peels back, then the vehicle is unstable.
- If the fins warp, then vehicle stability is compromised.

Experimental Data

- There is about a ten percent chance that the fins will warp during flight.
- The frame will not fail if loads do not exceed 5000 psi.

Computer Model

- Our simulations indicate that there is a 15% chance that flight loads exceed 5000 psi.

Notional Mission Success

Estimates of mission success (full distributions available)

- Mission **yellow** is most likely (50% \pm 10%)
- Mission **red** is second (35% \pm 5%)
- Mission **green** is third (15% \pm 5%)

Decompose these estimates into parts, subsystems, and functions that contribute to size and variability of estimates.

Munitions Example

Two sets of test data:

(Z = success/failure, X = covariates)

(S = spec measurement, X = covariates)

Measuring success/failure is expensive, so it would be useful to figure out how to use the spec data as a surrogate for measuring success/failure.

The ultimate aim is to predict reliability as a function of age, $P(Z = 1|age)$.

Assumptions

- For this example, the probability of success increases (monotonically) with the (unobserved) spec measurement.
 - We do not have data that lets us verify this.
 - We generally choose the functional form using engineering judgment.
 - No restriction on the functional form
- For this example, the spec measurements relate to the covariates through a linear regression.
 - Nick Hengartner has developed a nice way to do the estimation semi-parametrically that does not require the specification of this functional form
 - Accelerated testing

Munitions Example

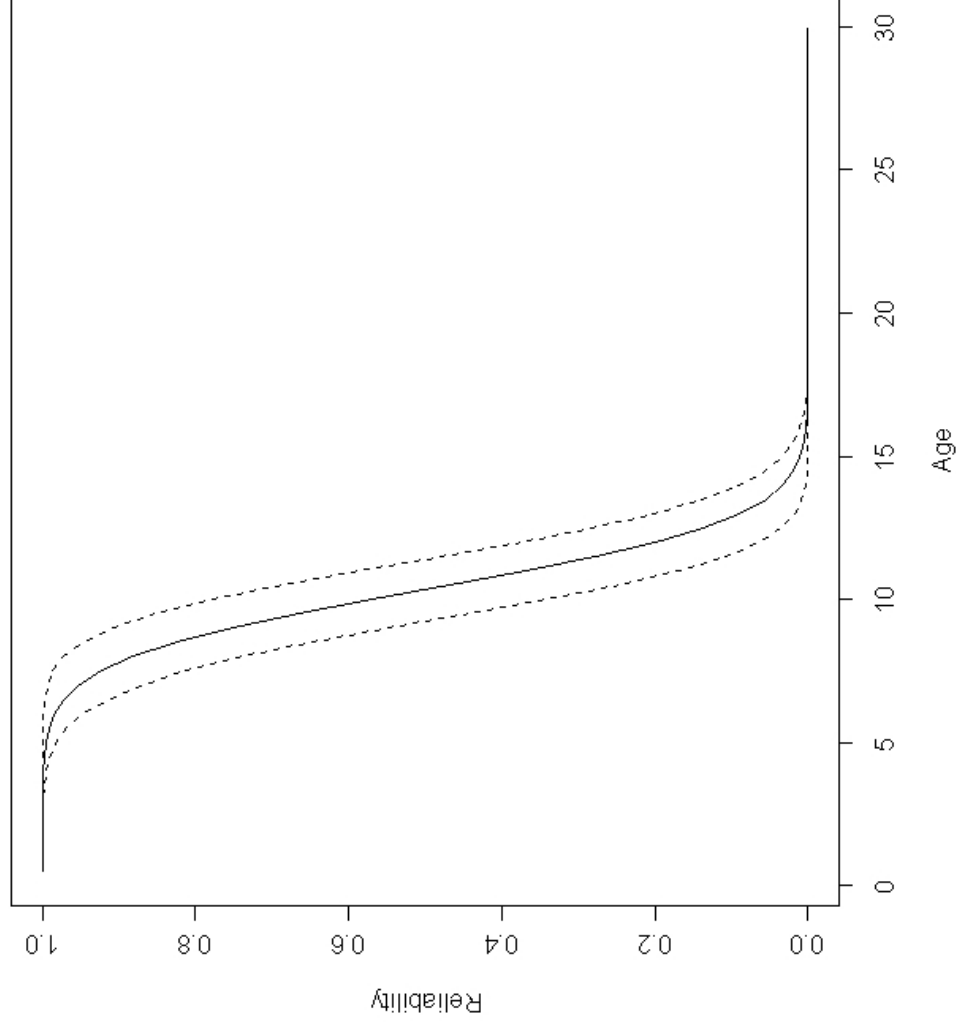
$$Z_i \sim \text{Bernoulli} \left(\Phi \left(\frac{S_i - \theta}{\sigma} \right) \right) \quad (\text{"surrogacy assumption"})$$

$$S_k \sim N(X\alpha, \gamma^2 \mathbf{I})$$

Can integrate out the unobserved S_i and get

$$Z_i \sim \text{Bernoulli} \left(\Phi \left(\frac{X\beta - \theta}{\sqrt{\gamma^2 + \sigma^2}} \right) \right)$$

Munitions Example





**Application of
Fisher's Combined Probability Test
To the Validation of the AIM-9X Missile Model**

***U.S. Army Conference on Applied Statistics
Napa Valley, CA
October 2003***

**Arthur Fries
Institute for Defense Analyses
afries@ida.org**

ACKNOWLEDGMENTS & DISCLAIMER



Portions of Arthur Fries' research were undertaken at the Institute for Defense Analyses (IDA) under IDA Central Research Project C9016, under IDA Professional Development Funding, and as part of various tasks sponsored by the Office of the Director of Operational Test and Evaluation (ODOT&E) in the Office of the Secretary of Defense (OSD) within the U.S. Department of Defense (DoD). The author gratefully acknowledges helpful technical review comments and inputs provided by A. Rex Rivolo and Robert R. Soule, IDA, numerous supporting analyses conducted by LeAnna Guerin, IDA summer intern, and briefing material graciously provided by Denise Beattie and Leonard Vance from Raytheon.

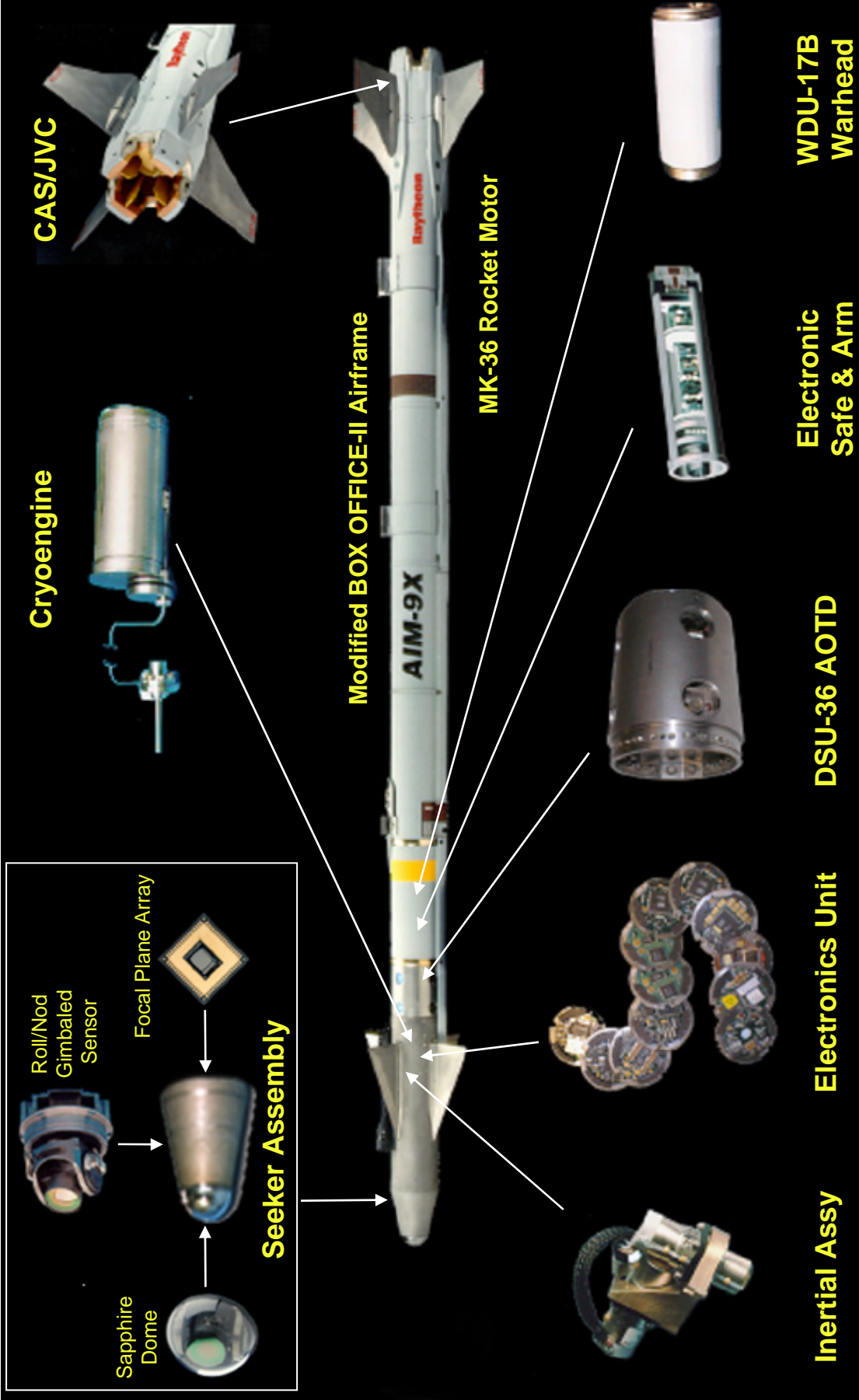
The contents of this presentation and the views therein are solely those of the author. No official endorsement by IDA, DoD, or ODOT&E is intended or should be inferred.



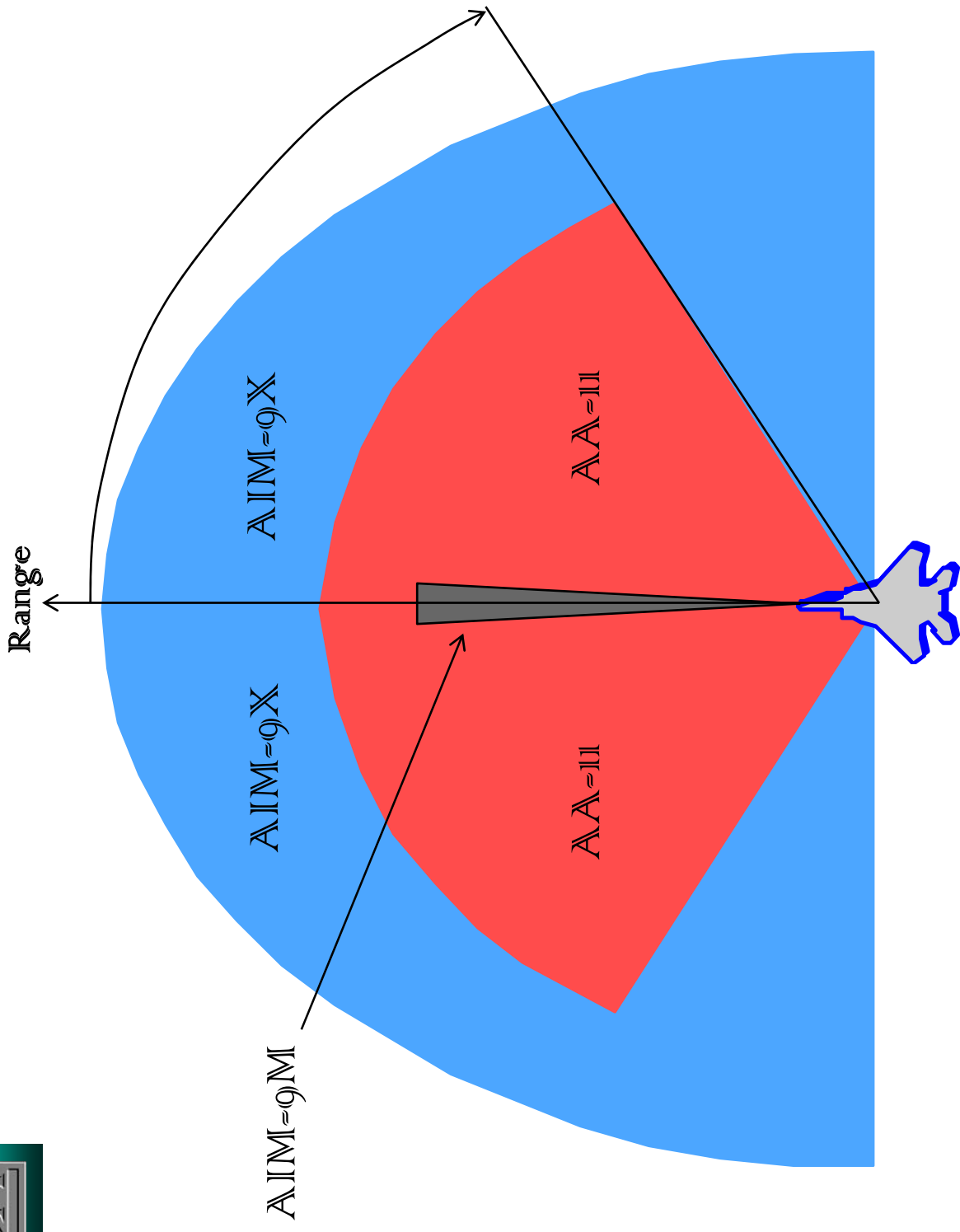
OUTLINE

- **AIM-9X**
 - System Description
 - VV&A Considerations
 - » How much to test?
 - » Where to test?
 - » How to analyze?
- **Fisher's Combined Probability Test**
 - Connection to M&S
 - Methodological considerations
 - » Accommodating measurement errors
 - » 1-sided or 2-sided testing?
 - » Other meta-analysis approaches
 - » Goodness-of-fit procedures
 - Statistical power
- **Application to AIM-9X**
- **Other Challenges**

AIM-9X: System Description



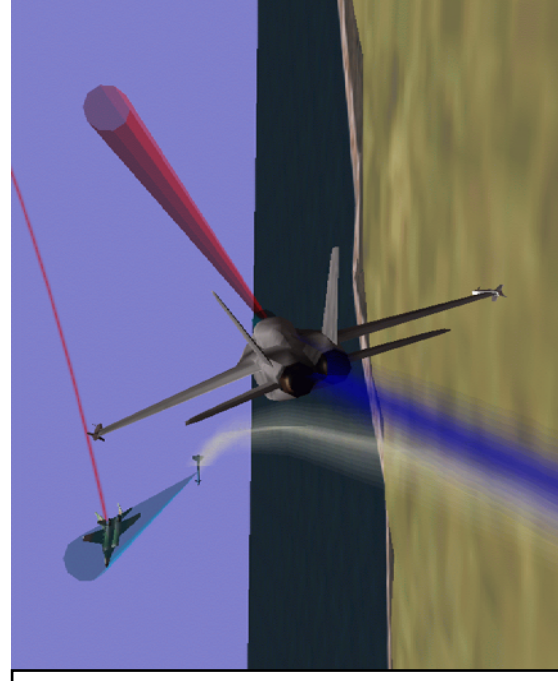
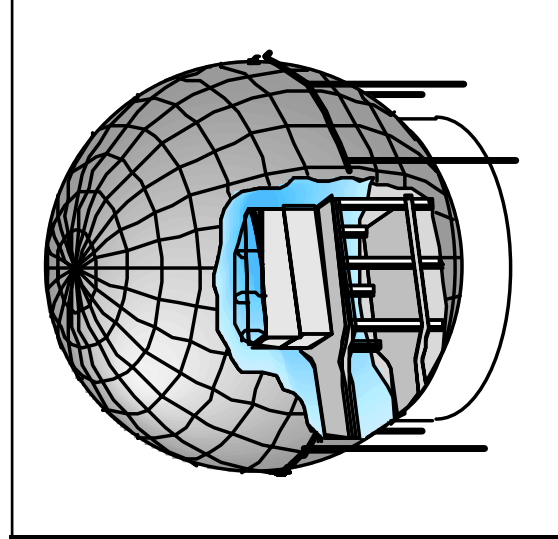
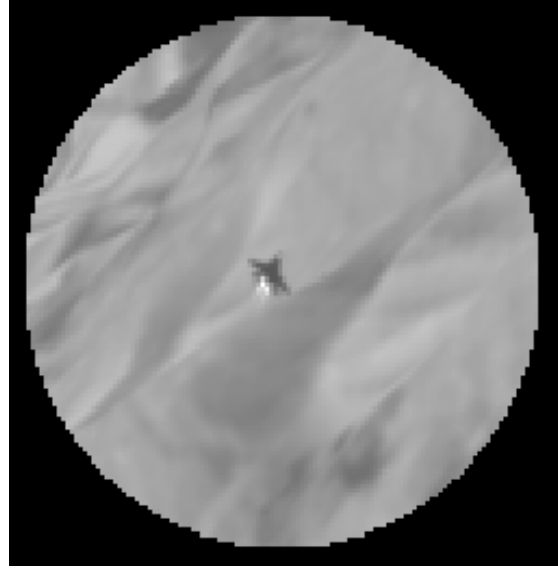
AIM-9X Launch Envelope (notional)





Modeling & Simulation Based Acquisition Program

- **AIM-9X Relies on Extensive Modeling & Simulation**
 - Improve the quality of flight tests
 - Reduce costs
- **Simulation Is Used to Predict All Levels of System Performance**
 - Weapon System, Missile System, Seeker, Kinematic
 - » Pk Performance Verified by Simulation Only
 - » Guided Flight Tests Used for Simulation Validation
 - Integrated Plan With Acquisition, DT, & OT Communities





Why Simulation Based Acquisition?

Simulation Based Acquisition (SBA) can lead to dramatic reductions in weapon system testing

Example: AMRAAM vs AIM-9X Full Scale Development

AMRAAM (1985-1990)	
Control Test Vehicle Flights*	29
Developed Guided Flights	79
Operational Test Flights	17
Total Guided Flights	96

AIM-9X (1999-2002)	
Seperation Control Test Vehicle Flights*	22
Developed Guided Flights	16
Operational Assessment Flights	5
Operational Test Flights	22
Total Guided Flights	43



Key Elements of Simulation Credibility

- **Software Accuracy**
 - Correctness of implementation
 - Accuracy of any data manipulations
 - » Unit conversions, coordinate transformations, etc.
 - Accuracy of software documentation

Measured by “verification”
- **Output validity**
 - Comparison with the “real world” (properly defined)
 - Accuracy of simulation inputs

Measured by “validation”
- **Software Stability**
 - Adequate management and control of M&S software, embedded data, and software documentation

“Configuration Management”
- **Data Accuracy**
 - Accuracy of any embedded data
 - » Physical constants, system parameters, etc.
 - Accuracy of any input data

“Data Quality”

Accreditation activities ensure robust assessments of simulation credibility in all these dimensions



WHAT IS “VALIDATION”?

- **Broad VV&A Context (AR 5-11, Management of Army Models and Simulations, 1997)**
 - “Validation is the process of determining the extent to which the M&S adequately represents the real-world from the perspectives of its intended use.”
- **Limited Statistical Context**
 - Quantify consistency between observed data (real world or test) and predicted outcomes
 - Provoke detailed follow-on discussions and analyses
- **No “One-Time” Validation**
 - Same model for different purposes
 - Continuous comprehensive evaluation



AIM-9X Flight Test Matrix

System Simulation Aspect Validated	DTIIB/C										OTI/A					DTII B/C					DTIID				
	e1	e2	e3	e4	e1	e2	e3	e4	o1	o2	o3	o4	o5	p1	p2	p3	p4	p5	p6	p7	p8	p9			
Target Signature (Low, (-)med, Hi)	-	-	-	-	L	-	-	-	-	-	-	-	-	-	H	-	-	-	-	-	-	-			
Target Terminal Aspect (Nose, Tail, Planform, Beam)	B	P	P	P	N	B	N	B	N	B	B	P	P	B	P	B	P	B	P	N	P	T			
Closing Velocity (Low, (-)med, Hi)	L	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	L	-	H	-	L	L			
OBA ((-)Low, High(vignetting))	-	-	H	H	H	-	H	-	H	-	H	-	H	-	H	H	-	-	-	-	-	-			
CM (midcourse,terminal w/ A,B,C,AC or (-) none	-	-	-	mB	-	-	mB	mB	-	-	mB	mB	B	B	B	-	A	AC	-	mA	C	-			
CM - Lag(+), Lead (+), (-)n/a	-	-	-	+	-	-	+	+	-	-	+	+	+	+	(-)	-	(-)	+	-	(-)	+	+			
Altitude (Low, (-)Med, High)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	L	H	H	L	H			
Background type: (blue sky [bs], desert, cloud, sea, horizon)	bs	d	c	d	bs	bs	bs	d	bs	bs	d	d	d	bs	h	d	s	bs	bs	d	bs	bs			
Contrast: positive(+), washout(0), neg washout(-)	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+			
Range (Low, (-)Med, High)	-	-	-	-	H	L	L	-	-	-	-	-	-	-	L	-	H	H	H	-	-	H			
Target Maneuver ((-)None, In-plane, Out of plane)	-	-	I	I	-	I	O	I	-	I	-	I	-	I	-	I	O	O	I	-	-	-			
Initialization (Radar, Helmet)	R	R	R	H	H	R	H	H	H	R	H	H	R	H	H	H	R	R	R	R	H	R			
Dome heating ((-)Moderate, High)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	H	H	H	H	-			
Afterburner (On, On Midgame, On Terminal, (-)None)	-	-	M	-	-	-	-	M	-	-	-	M	-	M	O	M	O	O	-	-	T	T			
	CL	WS	WS	CL	CL	EG	EG	CL	CL	CL	CL	CL	CL	CL	CL	CL	CL	CL	EG	EG	CL	CL			
	30 Jun99	1 Sep99	16 Dec99	31 Mar00	21 Apr00	23 May00	25 May00	28 Jun00	21 Jul00	12 Jun00	26 Sep00	17 Nov00	21 Dec00	5 Apr01	May-01	Aug-01	Sep-01	Oct-01							

EDM DT shot
 EDM OT shot
 PRM DT shot (DTIIB/C)
 PRM DT shot (DTIID)
 PRM DT Assist

IMPLEMENTATION CONSIDERATIONS



- **How Much To Test?**
- **Where To Test?**
- **How To Analyze?**

HOW MUCH TO TEST?



- **Whatever You Can Get**
- **20 – 40 ??**
- **Simulation-Based Investigation**

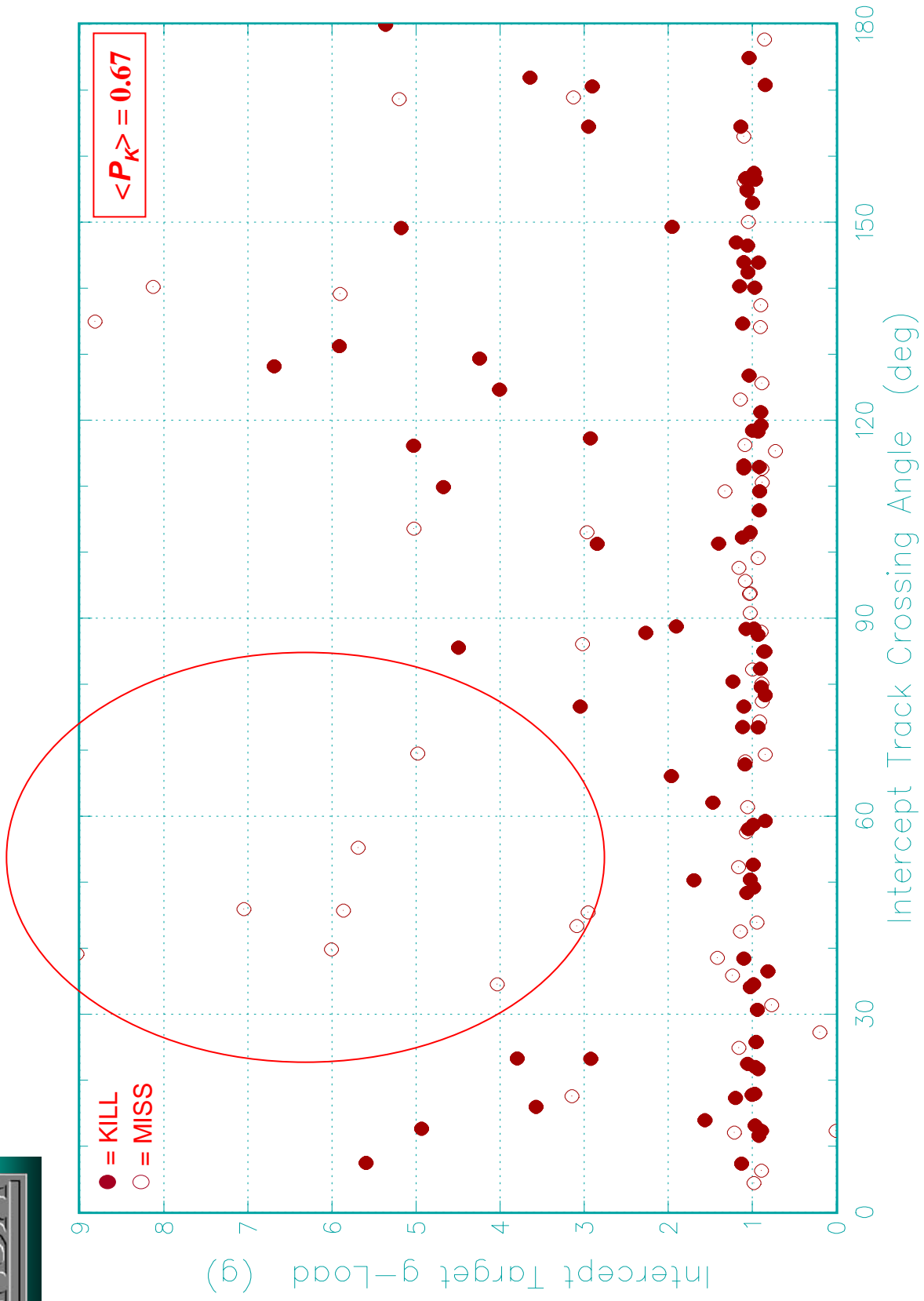


WHERE TO TEST?

- **One Point Replicated**
- **Few Points & Less Replication**
- **“Many” Points & No Replication**
 - Uniformly distributed
 - Concentrate on disparate “pockets”
 - » Based on M&S insights
 - » Based on physics insights
 - Include specific combinations
 - » Likely
 - » Extreme

M&S ALONE MAY BE INADEQUATE

(hypothetical example, not AIM-9X)



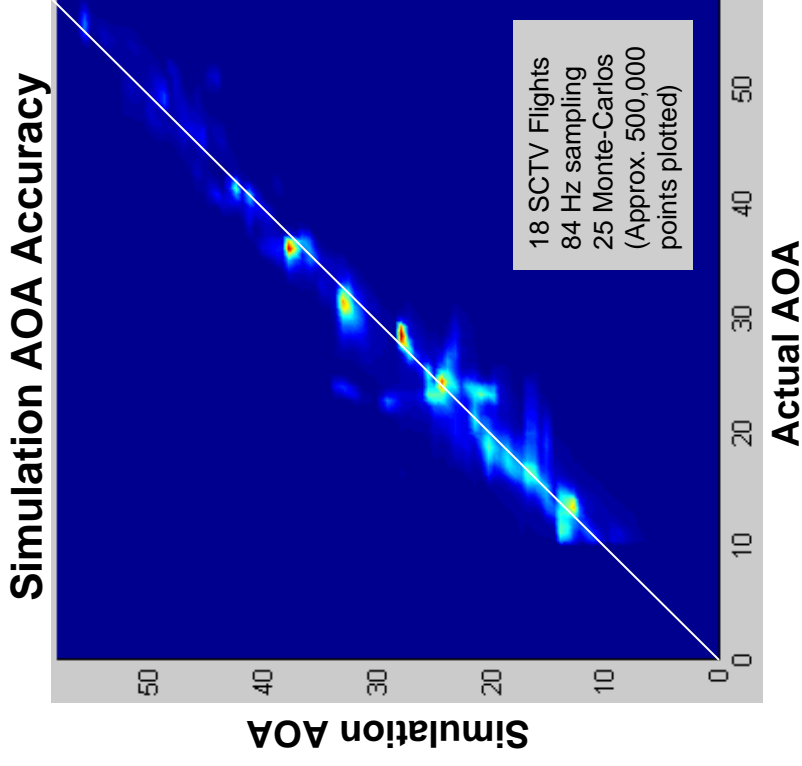
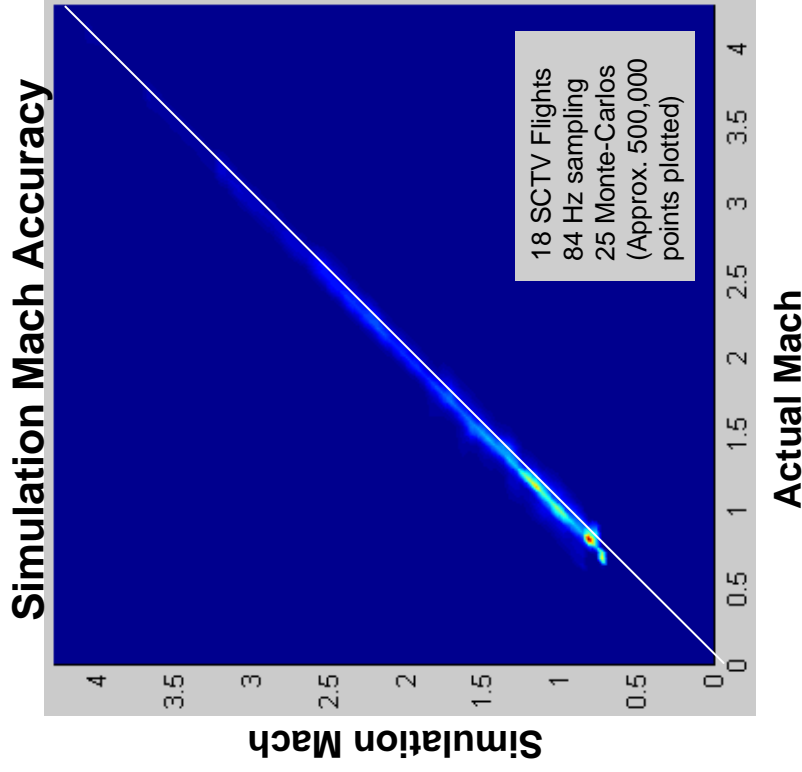
HOW TO ANALYZE?



- **Focus On End-Game Miss Distances**
 - *Old Method*: Look at the actual vs predicted miss distances and use our judgment
 - *Better Method*: Determine if actual miss distances are statistically plausible when compared to simulation predictions

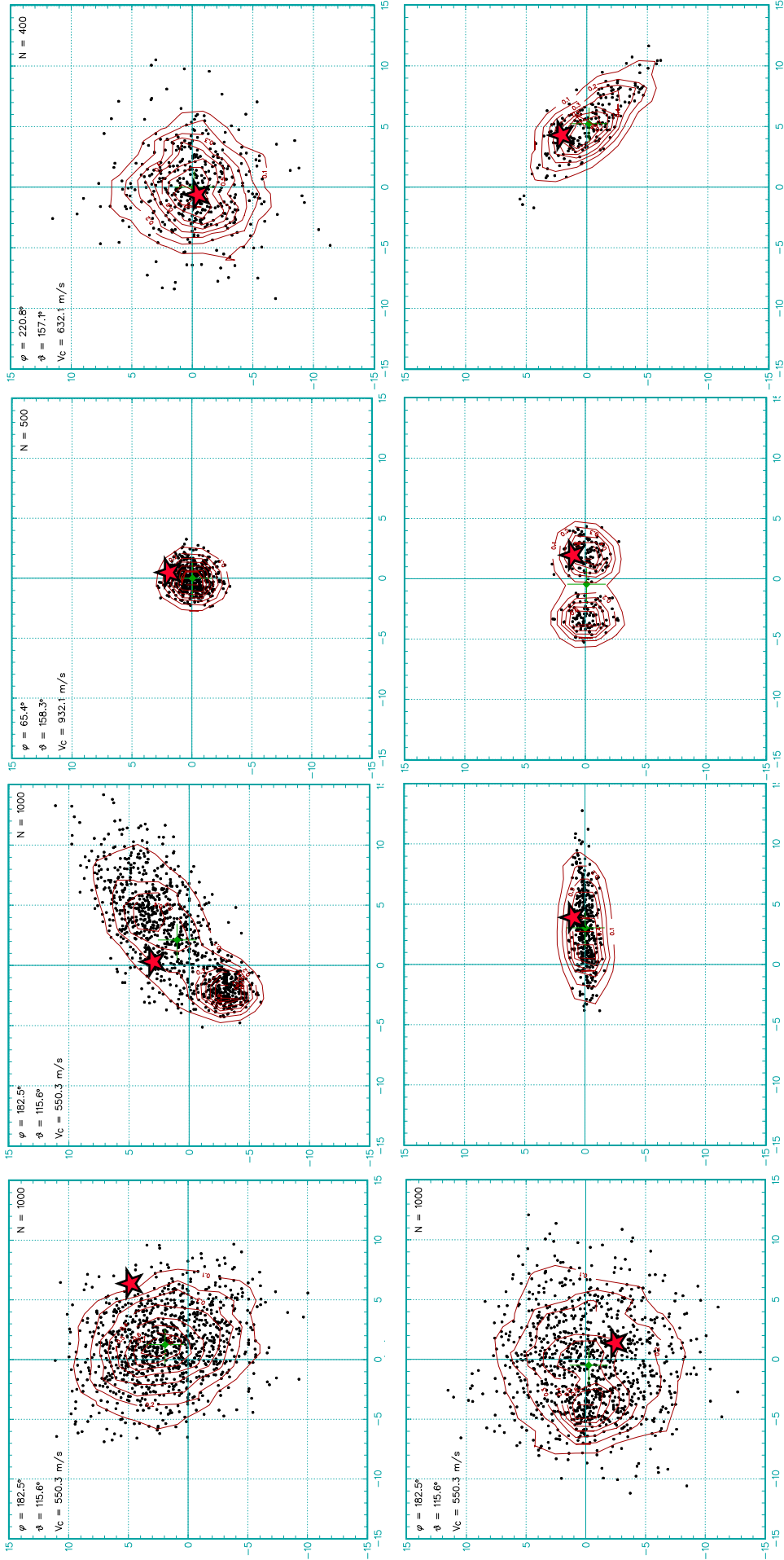


Excellent Kinematic Match



- Actual vs simulation values of Mach and AOA are plotted for all 18 SCTV flights
- Excellent overall 1:1 correlation is evident

HOW TO ANALYZE MISS DISTANCES?





- **Meta-Analysis Technique**
- **Combines Information From Different “Experiments”**
 - Same H_0
 - Different settings
 - Small sample sizes
 - Limited statistical significance



FISHER'S STATISTIC

$$F = - \sum_{i=1}^N 2 \log_e(P_i)$$

Single Experiment P-Value

H_0 : P_i is uniformly distributed in the interval [0 - 1.0]

↓
Implies



F is χ^2 distributed with $2N$ Degrees of Freedom

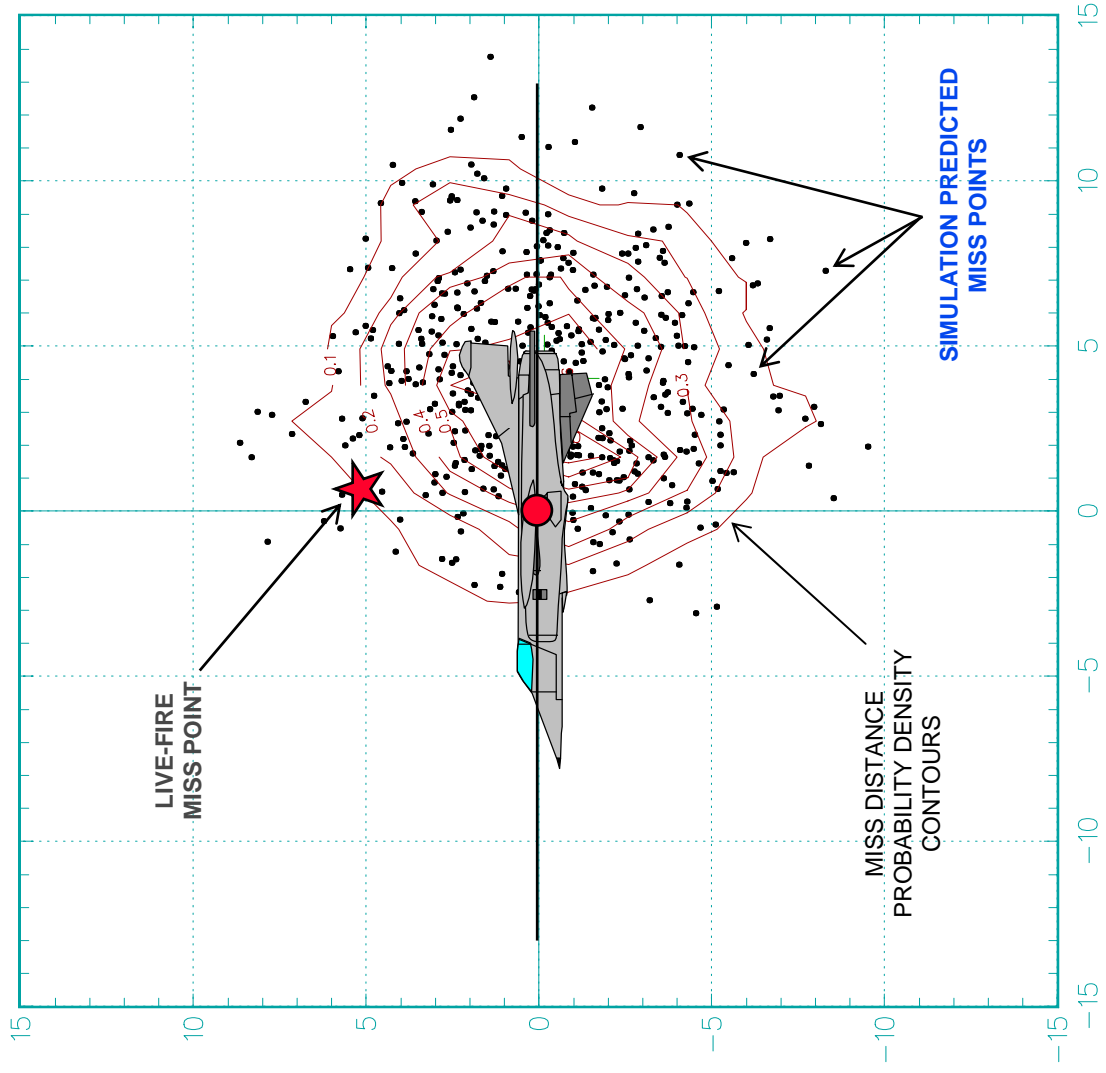


MISS DISTANCE MEASURE OF AGREEMENT

Definition: P_i (Tail Probability).
Fraction of prediction points
outside the tangent point
contour.

$$P_i = 0.16$$

Plane of Closest Approach





FISHER'S STATISTIC

$$F = - \sum_{i=1}^N 2 \log_e(P_i)$$

Single Tail Probability

H_0 : P_i is uniformly distributed in the interval [0 - 1.0]



Implies



F is χ^2 distributed with $2N$ Degrees of Freedom

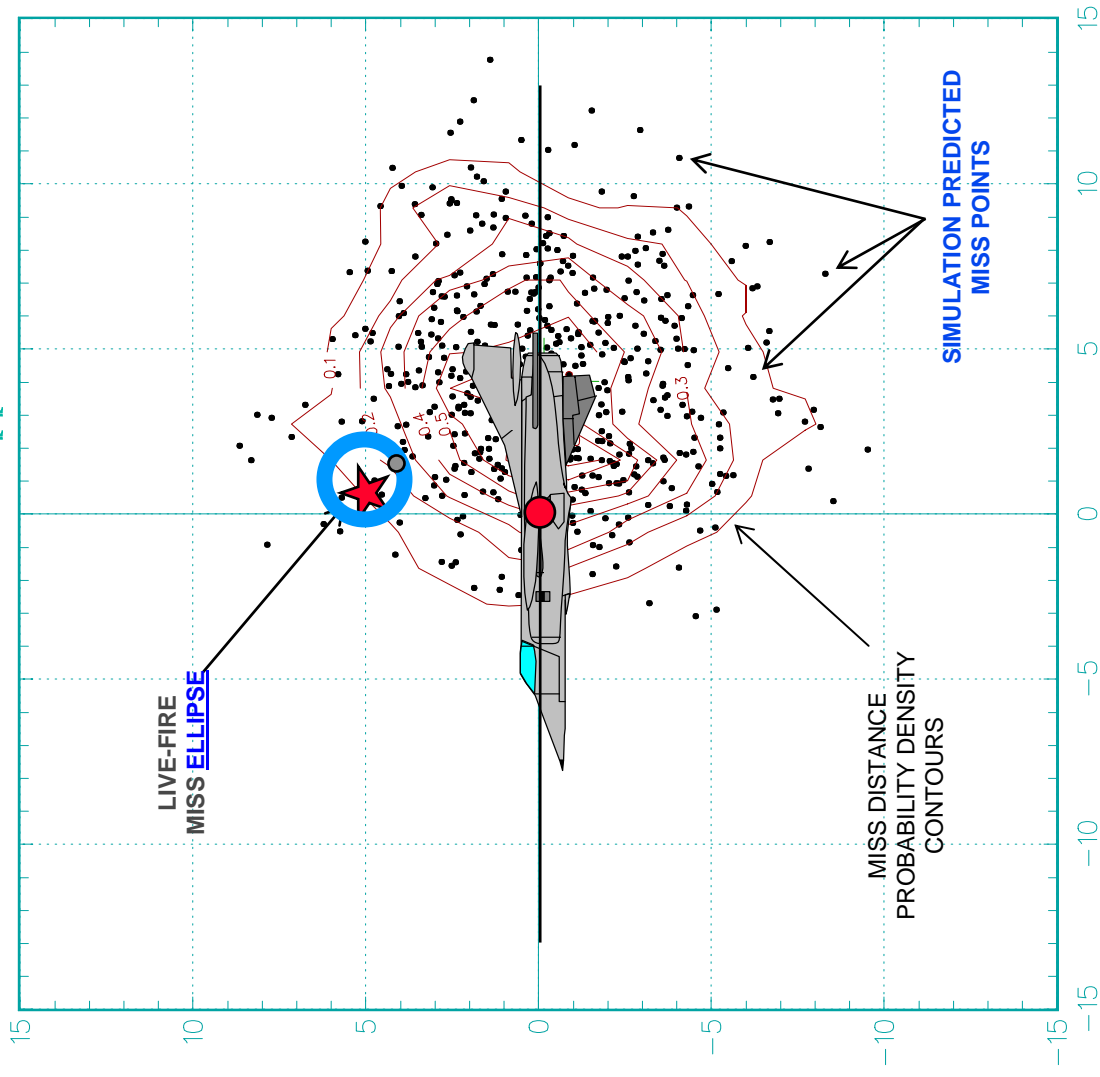
MISS DISTANCE MEASUREMENT ERRORS



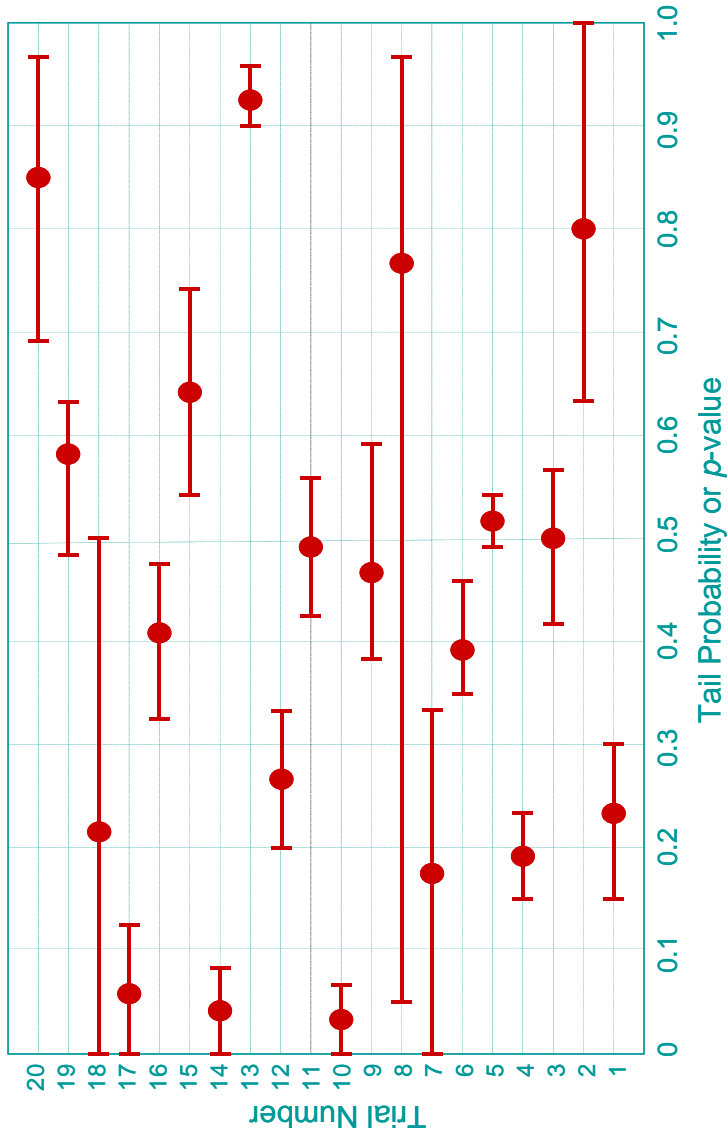
Definition: P_i (Tail Probability).
Fraction of prediction points
outside the tangent point
contour.

$$P_i = 0.02 - 0.28$$

Plane of Closest Approach

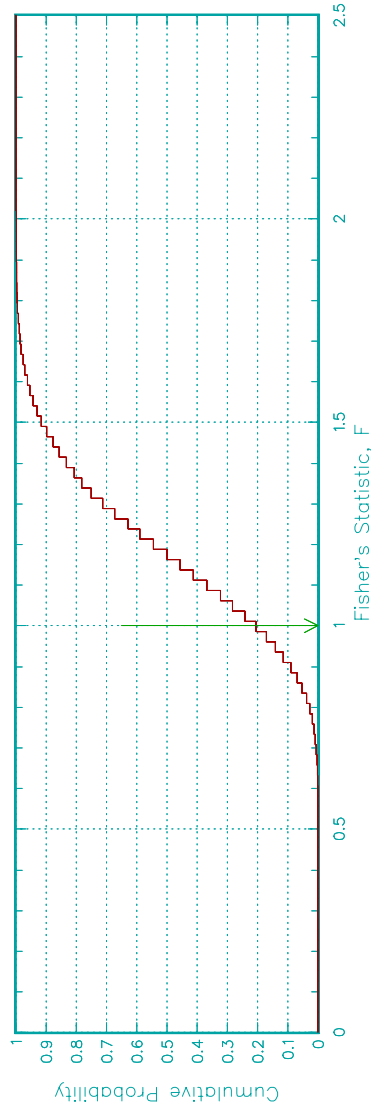
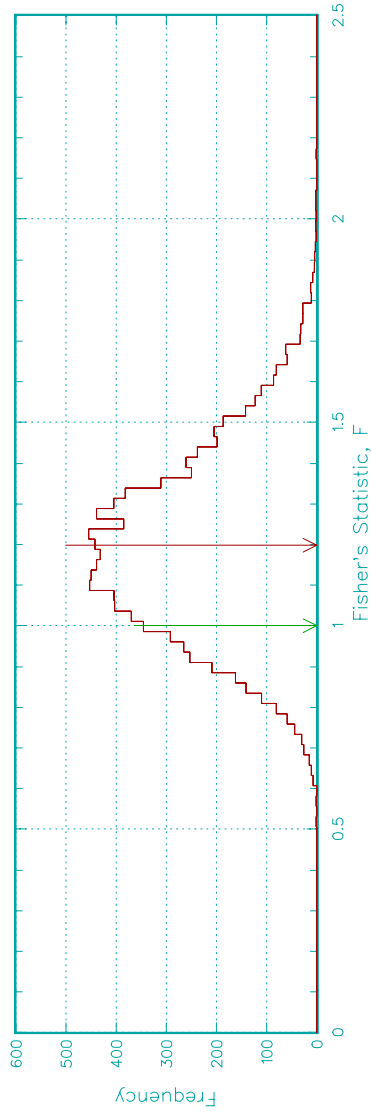


ENCOMPASSING MEASUREMENT ERRORS



Derived Tail Probabilities (p -values)

BOOTSTRAPPED FISHER STATISTICS



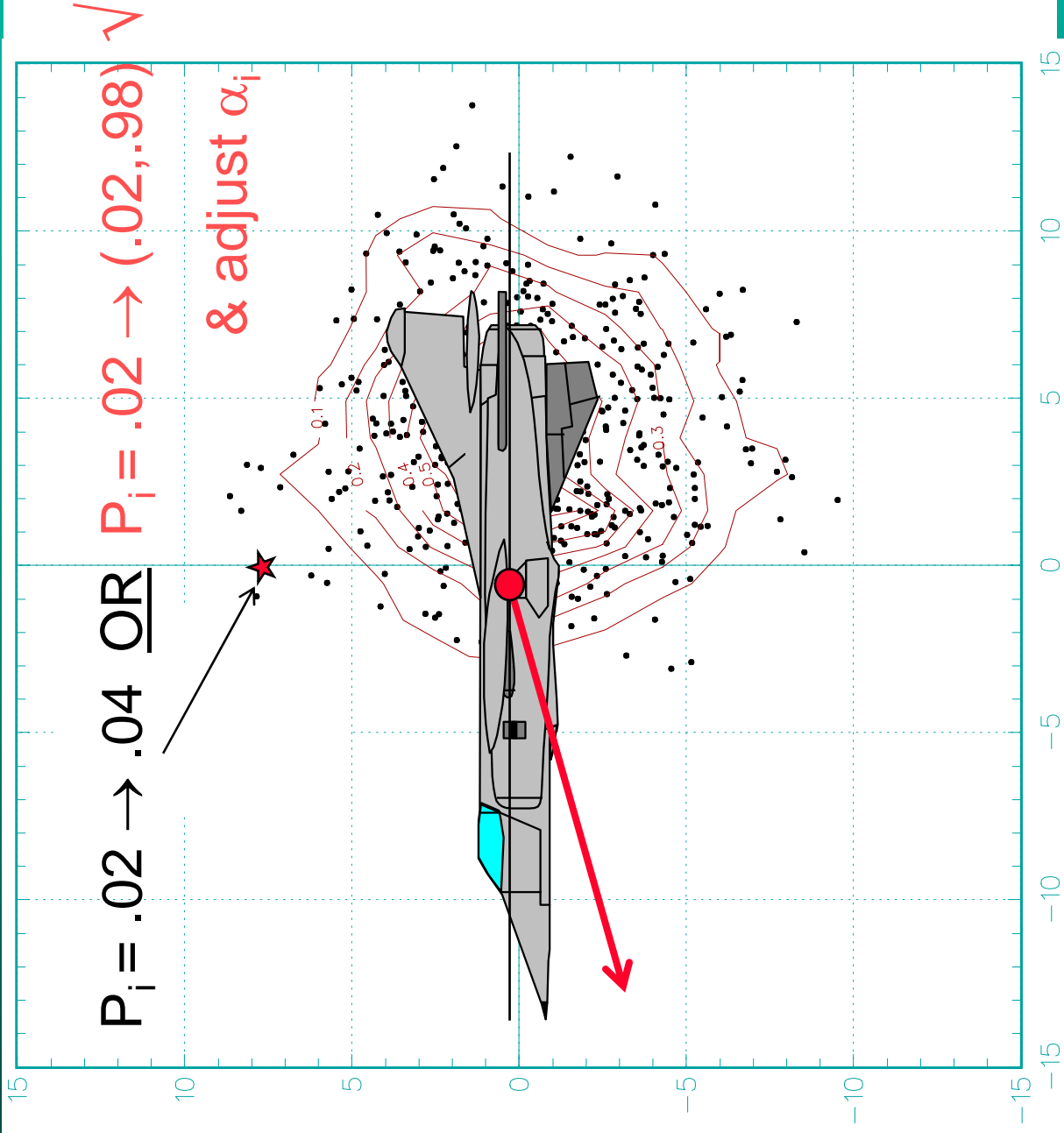
Distribution of 10,000 F -Values and Cumulative Probability



WHICH TAIL(S)?

- **Far-out** ↔ **Optimistic model**
(big misses are “bad”)
- **Close-in** ↔ **Pessimistic model**
(small misses are “bad”)
- **Both tails** ↔ **Both concerns**

CONVERSIONS TO 2-SIDED P_i



$P_i = .02 \rightarrow .04$ OR $P_i = .02 \rightarrow (.02, .98)$ ✓

& adjust α_i



EXAMPLE CALCULATIONS

- **Simulation Model Optimistic?**
 - $X = 2 [-\ln(.12) - \ln(.37) - \dots - \ln(.39)] = 14.4$
 - $X_O \ll \chi^2_{16}(0.05) = 26.3 \Rightarrow \text{NO!}$
- **Simulation Model Pessimistic?**
 - $X = 2 [-\ln(.88) - \ln(.63) - \dots - \ln(.61)] = 15.2$
 - $X_P \ll \chi^2_{16}(0.05) = 26.3 \Rightarrow \text{NO!}$
- **Simulation Model “Bad”?**
 - $X = \max(X_O, X_P) \ll \chi^2_{16}(0.025) = 28.8 \Rightarrow \text{NO!}$



META-ANALYSIS APPROACHES

- **One “Outlier” Test Result & Fisher ⇒ “Invalid”**
- **Other Tempered Meta-Analysis Procedures**
- **Endorse Classical Fisher Methodology**
 - Encourages investigation of “outliers”
 - » Test anomaly?
 - » Physical causes represented in model?
 - Promotes discussion
 - » Concept of “invalid”
 - » Statistical significance” vice
“practical significance”

GOODNESS-OF-FIT PROCEDURES



- **P_i's Consistent With A Uniform Distribution?**
 - Fisher procedure weights small values heavily
 - » $-\ln(p) \uparrow \infty$ as $p \downarrow 0$
 - Other standard goodness-of-fit tests
 - » Maximum difference in cdf's
 - » Averaged difference

STATISTICAL POWER



- **When the Simulation Model Is “Good”**
 - $S \equiv T$, for all factor combinations
 - With high confidence, **should not reject** H_0
 - Should attain nominally prescribed α 's
- **When the Simulation Model Is “Bad”**
 - $S \neq T$, for some/all factor combinations
 - With high confidence, **should reject** H_0
 - More “powerful” procedures reject more often



MONTE CARLO ANALYSES

- **Assume “w.l.o.g.”** $S_i \sim n(0,1)$
 - Each S_i is “known”
- **Consider** $T_i \sim [n(\mu_i, \sigma_i^2)]^{**} \rho_i$
 - Change mean, variance, shape
 - t -test (Z-test) optimal for $\mu_i = \mu, \sigma_i = \rho_i = 1$
- **Statistical Characterizations**
 - t -test (2-sample paired)
 - Fisher
 - Goodness-of-fit



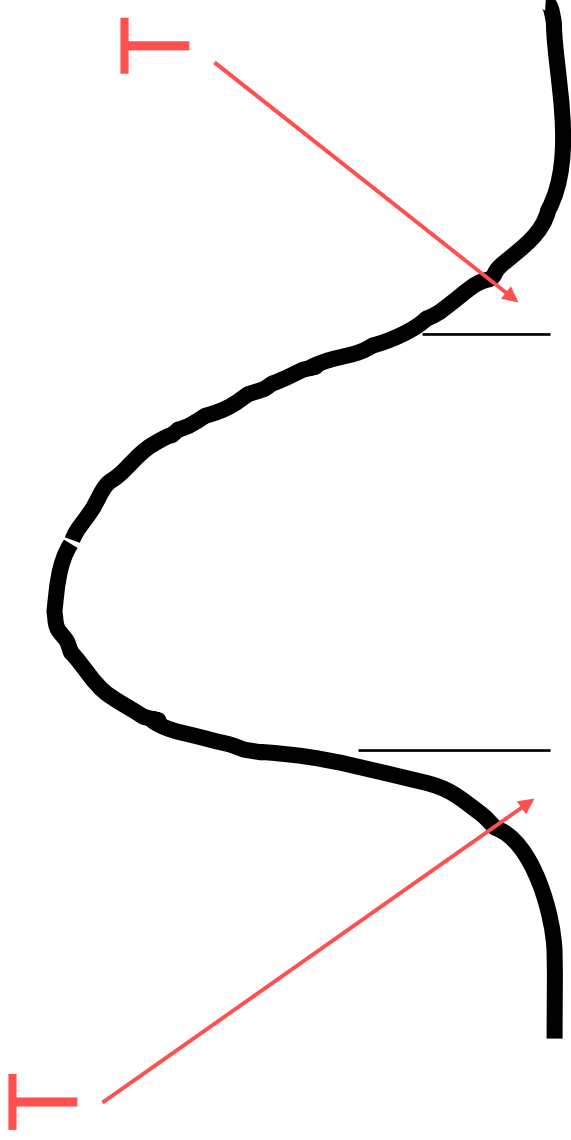
FALSE REJECTIONS

- **By construct**, all procedures attain pre-set α 's
- ***t*-test relies on “known” $S_i \sim n(0,1)$**
 - Could generalize to any $n(\mu_i, \sigma_i^2)$
 - But not to “irregular” prediction distributions
- **Other procedures rely merely on “known” S_i**
 - Appropriate for any prediction distribution

VALID REJECTIONS - CASE 1



S

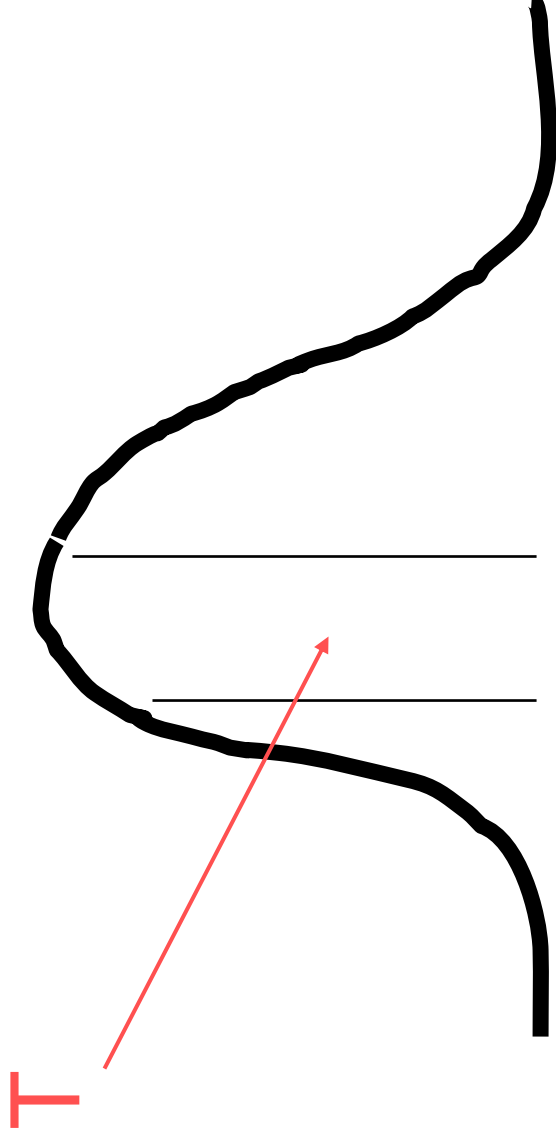


Fisher: Powerful for **T** in the tails of **S**
($\pm \mu, + \sigma^2, + \rho$)



VALID REJECTIONS - CASE 2

S



Fisher: Weak for **T** in “middle” of **S**
($-\sigma^2, -\rho$)



POWER - CASE 1

- **Fisher Dominates**
 - Fisher \gg t-test
 - Fisher $> \approx$ GOF's
- **Exception: Shift μ only**
 - t-test is “optimal”
 - Fisher is relatively efficient
 - » $S \sim n(0,1)$, $T \sim n(0.5, 1)$, $N = 25$
 - t-test **0.88** **0.78** **0.67** **0.40**
 - Fisher **0.85** **0.77** **0.65** **0.39**



POWER - CASE 2

- $S \sim n(0,1)$, $N = 25$, $\alpha = 0.05$
 - » $T \sim n(0, 0.6^2)$, Fisher = 0
 - GOF 0.55-0.83
 - t-test 0.05
 - » $T \sim n(0.4, 0.6^2)$, Fisher = 0.04
 - GOF 0.21-0.38
 - t-test 0.41

• **GOF >> Fisher**

• **GOF > \approx t-test**



“VALIDATION” CONCERNS

- **Case 1 >> Case 2**
 - **Optimistic Model!!**
 - **Pessimistic Model??**



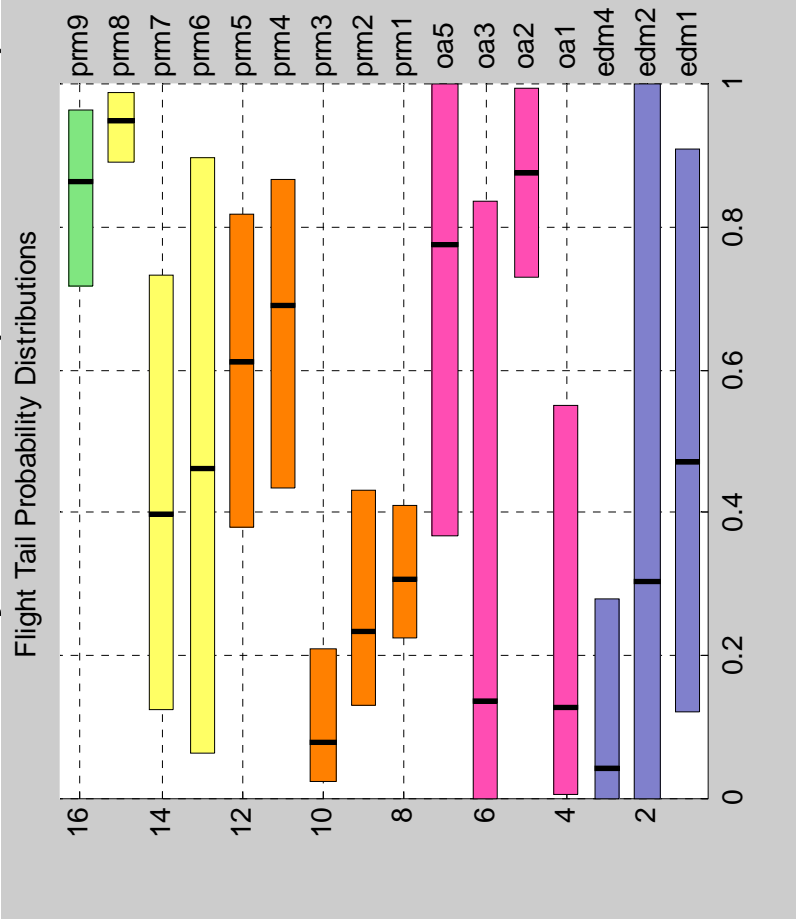
LIKELIHOOD OF CASE 2?

- **1-Dimensional Model**
 - Accurate μ 's
 - Overestimation of σ 's
- **Multi-Dimensional Model**
 - Annulli



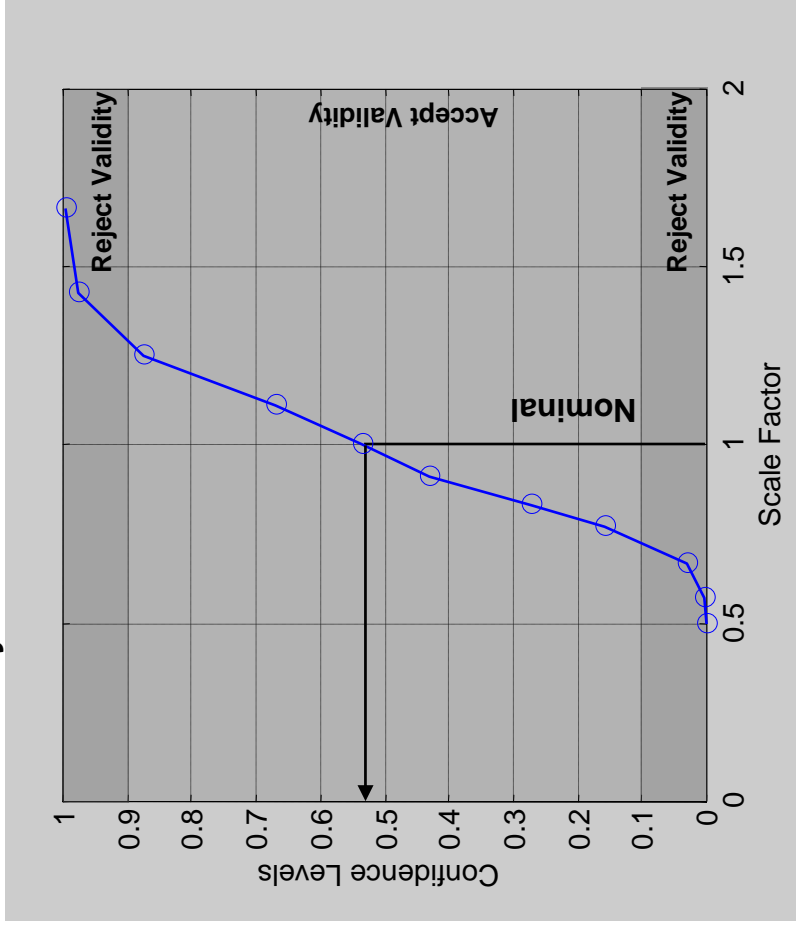
APPLICATION TO AIM-9X

Tail Probability Distributions (w/ uncertainties)



Tail probabilities distributed uniformly from 0 to 1 indicate simulation miss distance is valid

Summary Fisher Miss Distance Results



A Confidence level between 0.1 and 0.9 for nominal miss distance scaling factor (SF = 1) indicates sim miss distance prediction is valid



AIM-9X M&S VALIDATION HISTORY

- Don't Need To Do It
- Can't Do It
- Can Only Test At One Design Point
- **Will Do It!**
 - Implement rigorous managerial procedures & controls
 - Continually inform all T&E organizations
 - Utilize Fisher methodology
- **Awards From DoD Modeling & Simulation Office**
 - Government M&S award for 2001
 - Contractor M&S award for 2001
- **Recent Flight Test Anomaly**

OTHER CHALLENGES



- **Non-Point Predictions**
 - Curves Generated
 - No Single Primary Point Of Emphasis
- **Deterministic Models**
 - Single Prediction For Any Set Of Inputs

**Research within the Department of Homeland
Security: Key Issues and Priorities**

Statistical Research Opportunities

U.S. Army Conference on Applied Statistics

Napa Valley, CA

October 2003

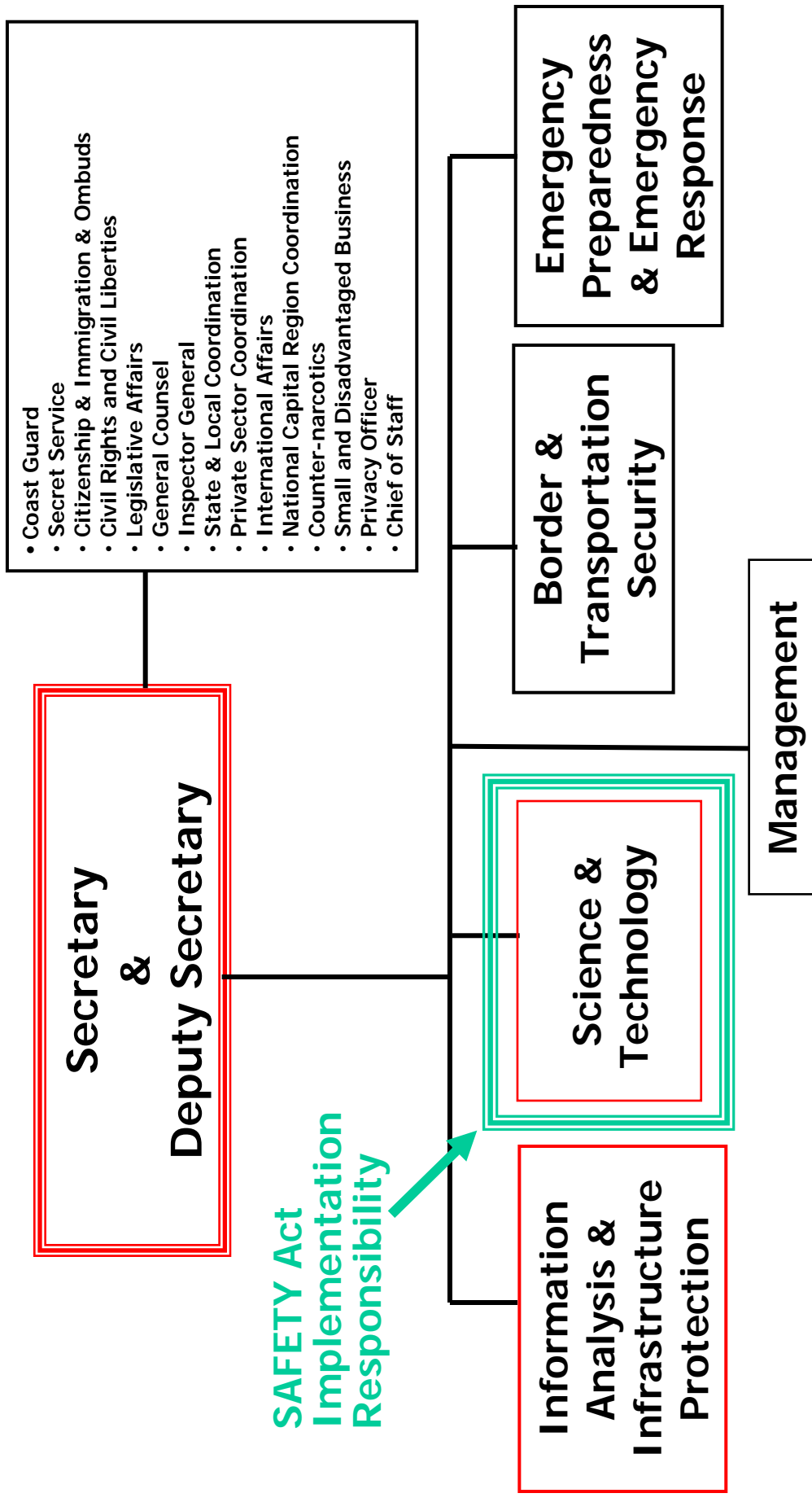
Parney Albright

Department of Homeland Security

Outline

- Disclaimer
- DHS Overview
 - Science & Technology Division
- Specific Research Opportunities
- SAFETY Act

Department of Homeland Security



S&T Division Components

- Office of Plans, Programs & Budget
- Homeland Security Advanced Research Projects Agency (HSARPA)
- Office of Research & Development
- Office of Systems Engineering & Development

Office of Plans, Programs & Budget

- Identify Requirements
- Create Strategic Initiatives
- **Prioritize Investments**
- Match Goals To National Policies
- Ensure Goals Are Met

HSARPA

- Engages industry, **academia**, government, etc.
- Innovative R&D
- Rapid prototyping
- Technology transfer

Office of Research & Development

- Cadre Of Scientists & Engineers
 - Office of National Laboratories
 - Homeland Security Laboratories
 - Components of DOE Laboratories
 - Federal Laboratories
- **Supports University & Fellowship Programs**

Office of Systems Engineering & Development

- Executes Transition Of Large-Scale Or Pilot Systems To The Field

Specific Research Opportunities

- Postgraduate & Postdoctoral Fellowships
- Scholarships
- Grants & Contracts
- Centers Of Excellence
- Homeland Security Institute

Homeland Security Scholars & Fellows Program

- 2003 Results
 - 2,500 applicants
 - 101 recipients of ~ \$2M
 - Required internships
- 2004 Cycle Unopened
 - Oak Ridge Institute for Science and Education
 - <http://www.ornl.gov/dhsed/>
 - Increased funding

Broad Agency Announcement (BAA) On S&T Requirements

- \$30M For R&D Funding
 - Mostly technologies
 - Asset protection prioritization schemes
- 2003 Completed
 - Announced May 14
 - June 13 deadline
- Technical Support Working Group
 - www.tswg.gov
 - Initial one-page proposal process

Centers Of Excellence

- **Scope**
 - One Center by November 2003
 - Up to 9 others by end of 2004
 - \$10M for FY04
- **2003 Completed**
 - Announced July 23
 - Deadline August 11
 - Risk-based economic modeling on the impact and consequences of terrorism
- **Administration**
 - Oak Ridge Associated Universities
 - <http://www.ora.gov/dhsuce/>

Homeland Security Institute

- Dedicated FFRDC
- RFP out November?

SAFETY ACT

- Purpose
- Implementation
 - Reviewers
 - Technical
 - Business
 - Insurance
- Relation To Studies

A Taxonomy of Terrorisms

James R. Thompson, Rice University

*Supported in part by W911NF-04-1-0354,
Army Research Office: Durham

What is Terrorism?

Violent Acts Designed to Modify the Behavior of an Opponent.

Some Famous Terrorists and Counter Terrorists

- **Francis Marion versus Banistaire Tarleton**
- **John Mosby versus George Armstrong Custer**
- **Michael Collins versus Wilson's Black and Tans**
- **Home Army (Poland) versus the Gestapo and the KGB**
- **Dzhokhar Dudayev versus *Spetsnaz* General Rokhlin**

Matricizing Barely Quantitative Data



Herman Kahn 1922-1983

Categories of Terrorism

- Psychological=1
- Chemical=2
- Biological=3
- Nuclear=4
- Agricultural=5
- Conventional weapons=6

Sources of Terror

- State=1
- Separatist/Independence Group=2
- Economic Cartel=3
- Cultural and Linguistic Group=4
- Religious Group=5
- Distraught Individual (aka nutcase) = 6

Levels of Lethality

1. Nonlethal
2. Low
3. Medium
4. High
5. Massive

Complexity (Cost) Level

1. Slight
2. Low
3. Medium
4. High
5. Advanced

Threat Level

1. Slight
2. Low
3. Medium
- 4. High**
- 5. Imminent**

Example: Sarin Attack by al-Quaida in NY Subway

- Category = Chemical
- Source = Religious Group
- Lethality Level = 4
- Complexity Level = 3
- Threat Level = 3

Example: Aryan Nations Time Delayed Incendiary Bombs in National Forests

- Category = Agricultural
- Source = Separatist
- Level of Lethality = 1
- Complexity Level = 1
- Threat Level = 1

Example: Palestinian Rocket Attack on Commercial Airliner

- Category= Conventional Weapons
- Source =Religious/Separatist/Independence Group
- **Lethality Level = 4**
- Complexity Level = 3
- Threat Level = 1

Problems of Aggregation

- The threat level of a rocket attack from a Palestinian group may be small. However, the aggregate threat level from all potential groups is the important matter
- Consequently, simple scoring might be replaced by probabilities
- This is not easily done

Lethality Level Must Play An Important Role

- Even if it is deemed unlikely that invasion and capture of a nuclear power facility occur, the consequences are so dire that security measures must be taken
- Trying to multiply threat level by probability of occurrence is probably not a good way to decide whether measures should be taken for the threat

Allocation of Security Resources Requires

Aggregation Across Sources

- **Guarding against rocket attacks from al-Quaida, Hamas, Hezbollah, al-Aqsa Martyrs Brigade, Russian Mafia, Mexican Mafia, Aryan Nations, Nuts, etc., must be pooled when considering whether to put countermeasures on aircraft**
- **Guarding against pollution of reservoir must pool across sources when deciding whether to put security system around reservoir**

Importance of Psychological Terrorism

- Acts do not produce terror unless the target population knows about them.
- Psy-ops can give the target population a feeling of insecurity
- Propaganda targeted toward potentially sympathetic subpopulations can produce sympathizers, even fifth columns

Active Psychological Terrorism Against Opponents

- Generally very ineffective in WW II.
- Ineffective in Korean Conflict
- **Effective in Indochina War Only Due to Agitprop Facilitation**

Propaganda Works Best If It Is Based On Reality

- British Class System
- Exploiter (William Joyce)



Lord Haw Haw

DVD

DIGITALLY RESTORED IN 35MM
STARRING BASIL RATHBONE AND NIGEL BRUCE

Sherlock Holmes
AND
THE
**VOICE OF
TERROR**



WITH HILARY BLOOD AS MRS. LINDA BASS AND IN "THE LAST RITE" BY SIR ARTHUR COSSAN BRYLE
PRODUCED BY HOWARD BENEDICT DIRECTED BY JOHN BARLON



Vera Lynn



Osama bin Laden appears to have no active propaganda activity targeted to the West.

Missed Opportunities for bin Laden to Take Credit

- Beltway Snipers
- LAX Gunman
- Army Fragging
- Two Major Forest Fires in 2002
- San Bernadino Fires in 2003
- West Virginia Sniper(s)
- Anthrax in Florida and in DC Area

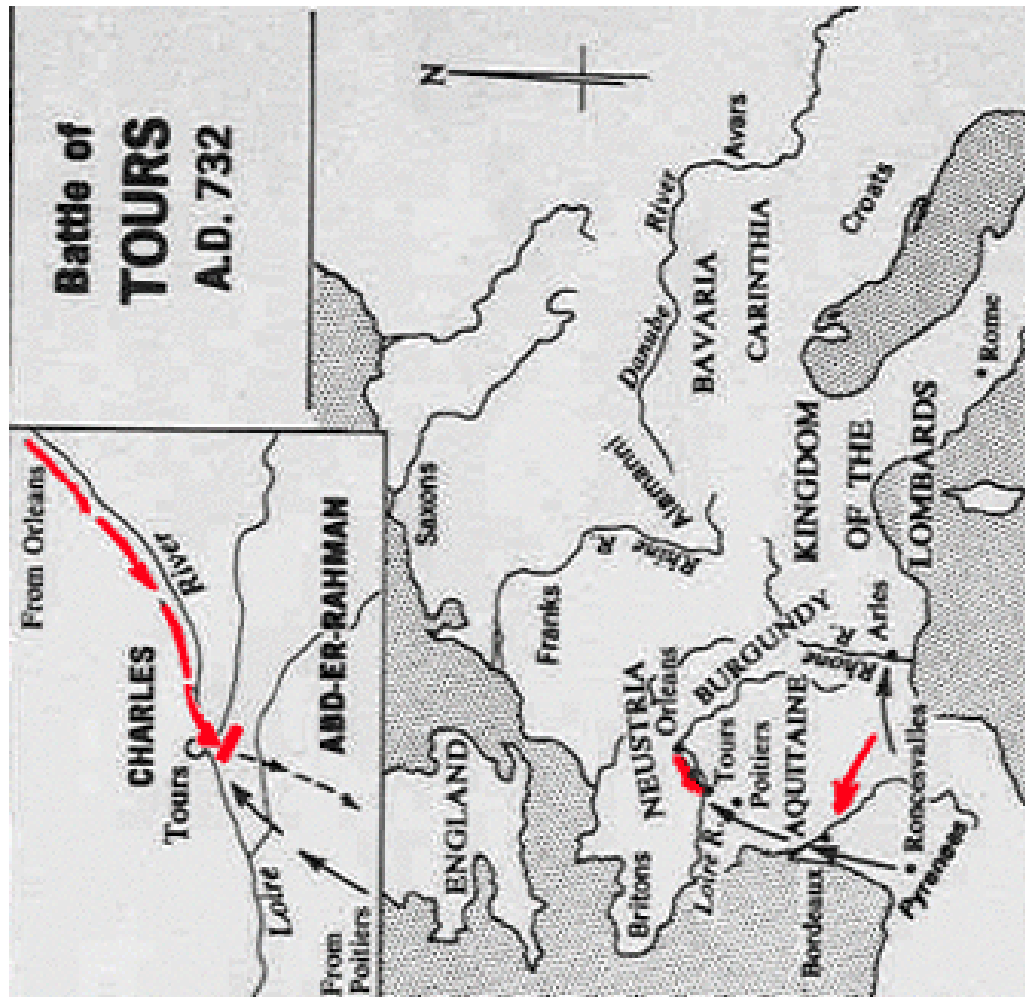
Why No English Language Propaganda Against the U.S.A.?

- Propaganda Against Enemies Not Effective
- “Victims” Strategy Works Best With
Natural Sympathizers
- Logistical Reasons
- Cultural Reasons

Some Dates in Early Islam

- 632 AD Death of the Prophet
- 636 AD Jerusalem Captured
- 640 AD Alexandria Captured
- 711 Spain falls (the Reconquest took 781 years)
- 732 Martel drives back Muslims at Tours

Battle of TOURS A.D. 732





Battle of Lepanto, Oct 7, 1571

The First 9-11



Battle of Vienna Sept 11 1683



Battle of Vienna Sept. 12, 1683

Advantages of a Non Active Victims Strategy

- **al-Jazeera provides gateway to Arabic-speaking populations**
- **Paranoia is induced in the target population**
- **Copycats are encouraged**
- **Other terrorist organizations may be encouraged to hide under the implied al-Quaida label**

Examples of Victims Strategy

Disadvantages of *Schrecklichkeit*

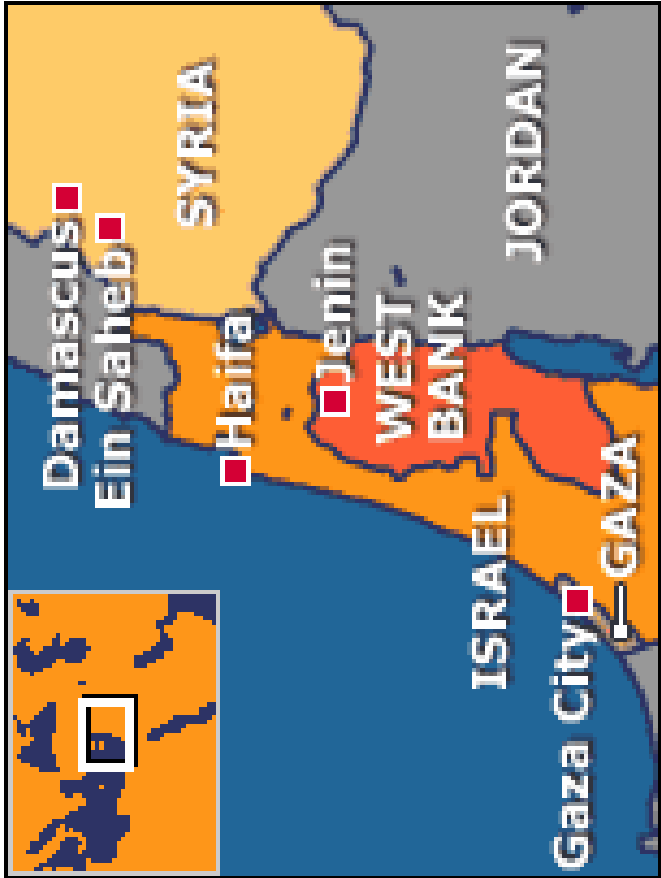
- The Alamo and Goliad
- Destruction of Warsaw 1939 and 1944
- Attack on Pearl Harbor
- Bombing of Coventry
- Assorted Japanese Atrocities





Jenin after Pacification





Avoiding Fifth Columns

1. The Germans were completely ineffective in building a fifth column in the USA
2. Eisenhower and Nimitz were CICs in Europe and Pacific respectively
3. The Japanese built only an ineffective fifth column amongst the Japanese-American community
4. This was suppressed by massive overkill
5. We may be on a trajectory which could produce a fifth column in the USA

A Scenario Analogy to Understand Situation of Muslim Americans

1. Ian Paisley Becomes UK Prime Minister on Platform of Ending IRA Terrorism
2. The Irish Republic Is Invaded by UK
3. The United States President Supports the British “Crackdown on Terrorism”
4. The U.S. Provides F-16s, Blackhawks To UK
5. The US President Describes Catholicism As A “Religion of Peace”
6. Attorney General Cracks Down on Irish-American and Catholic Organizations

7. US President describes Rev. Paisley as a “man of peace”
8. Cork and Galway are bombed by UK “in retaliation for IRA attacks”
9. US President urges Irish to replace their Taoiseach (President) with a leader more acceptable to the UK

UK Builds Security Fence



Conclusions

- **Terrorism is a multivariate problem**
- **Risk should be dealt with by taking an EDA approach rather than one based on expectations**
- **Aggregation of threat levels must be addressed**
- **Counter-strategies must be developed in the light of a realistic representation of the problem**
- **Al-Quaida may be using psy-ops much better than supposed. The establishment of a fifth column in the USA is a real possibility.**

DoD's Role in Homeland Security – Experimental Opportunity and Experimental Results

The DoD role in homeland security can be summarized as follows: (1) **homeland defense**, the protection of United States territory, domestic population, and critical defense infrastructure against external threats and aggression; and (2) **civil support**, providing military support to civil authorities at the federal, state, and local levels across a range of conditions. The Secretary of Defense described the three circumstances under which DoD assets would be involved in homeland defense and civil support missions:

- In **extraordinary circumstances**, DoD would conduct military missions such as combat air patrols or maritime defense operations. DoD would take the lead in defending the people and the territory of our country, supported by other agencies. Included in this category are cases in which the President, exercising his constitutional authority as Commander in Chief and Chief Executive, authorizes military action to counter threats within the United States.
- In **emergency circumstances**, such as managing the consequences of a terrorist attack, natural disaster, or other catastrophe in support of civil authorities, DoD could be asked to act quickly to provide capabilities that other agencies do not possess or that have been exhausted or overwhelmed.
- In **non-emergency circumstances of limited scope or planned duration**, DoD would support civil authorities where other agencies have the lead—for example, providing security at a special event such as the 2002 Winter Olympics, or assisting other federal agencies to develop capabilities to detect chemical, biological, nuclear, and radiological threats.

The DoD cannot provide for all aspects of homeland security. The homeland security mission requires the use of the full range of political, economic, diplomatic, and military tools, including in particular enhanced intelligence to improve potential detection of future attacks. The purpose of this presentation is to present areas whereby the Army experimental design and research community can assist DoD in establishing, defining and preparing training and acquisition programs for its Homeland Security role, to include means and methods of assisting in the preparation of concepts, plans, and contingencies. This presentation will also present a discussion of the HLS results from a recent Army/JFCOM experiment

2. Author Information:

Paul J. Deason, Ph.D. TRADOC Analysis Center (TRAC) – White Sands Missile Range

3. Type of Paper: Technical

4. Equipment Needed: Projector with cable; computer w/CD, Zip or floppy drive (we can bring a laptop if no computer is provided)

5. Telephone numbers: Dr. Paul Deason (505) 678-1610 (DSN 258)

6. Email: paul.deason@us.army.mil

Research Opportunity - DoD's Role in Homeland Security



Paul J. Deason, Ph.D.

US Army TRADOC Analysis Center (TRAC) – White Sands Missile Range

October 2003

Contents

- **Definition -- Federal Government's role in Homeland Security.**
- **DoD's role in Homeland Security.**
- **Strategy for Homeland Security/Homeland Defense.**
- **Example of Homeland Security/Homeland Defense in a recent Joint war game experiment.**
- **Challenge for statisticians and researchers in conducting Homeland Security oriented experiments and analyses.**

Role of Federal Government in Homeland Security

- **Homeland Security is best accomplished by building on state and local capabilities.**
- **The role of the Federal Government in Homeland Security is to enhance the capabilities at the lowest level of government.**
- **The Federal Government's Department of Homeland Security (DHS) is to:**
 - **Consolidate Federal activity.**
 - **Integrate national preparedness and response.**
- **The DoD established US Northern Command to consolidate under a single unified command homeland defense and civil support missions that were previously executed by other military organizations, specifically: to conduct operations to deter, prevent, and defeat threats and aggression aimed at the United States, its territories, and interests within the assigned area of responsibility; and as directed by the President or Secretary of Defense, provide military assistance to civil authorities including consequence management operations.**
- **US Northern Command plans, organizes, and executes homeland defense and civil support missions. The command will be assigned forces whenever necessary to execute missions as ordered by the President.**

Decision Flow - Homeland Security

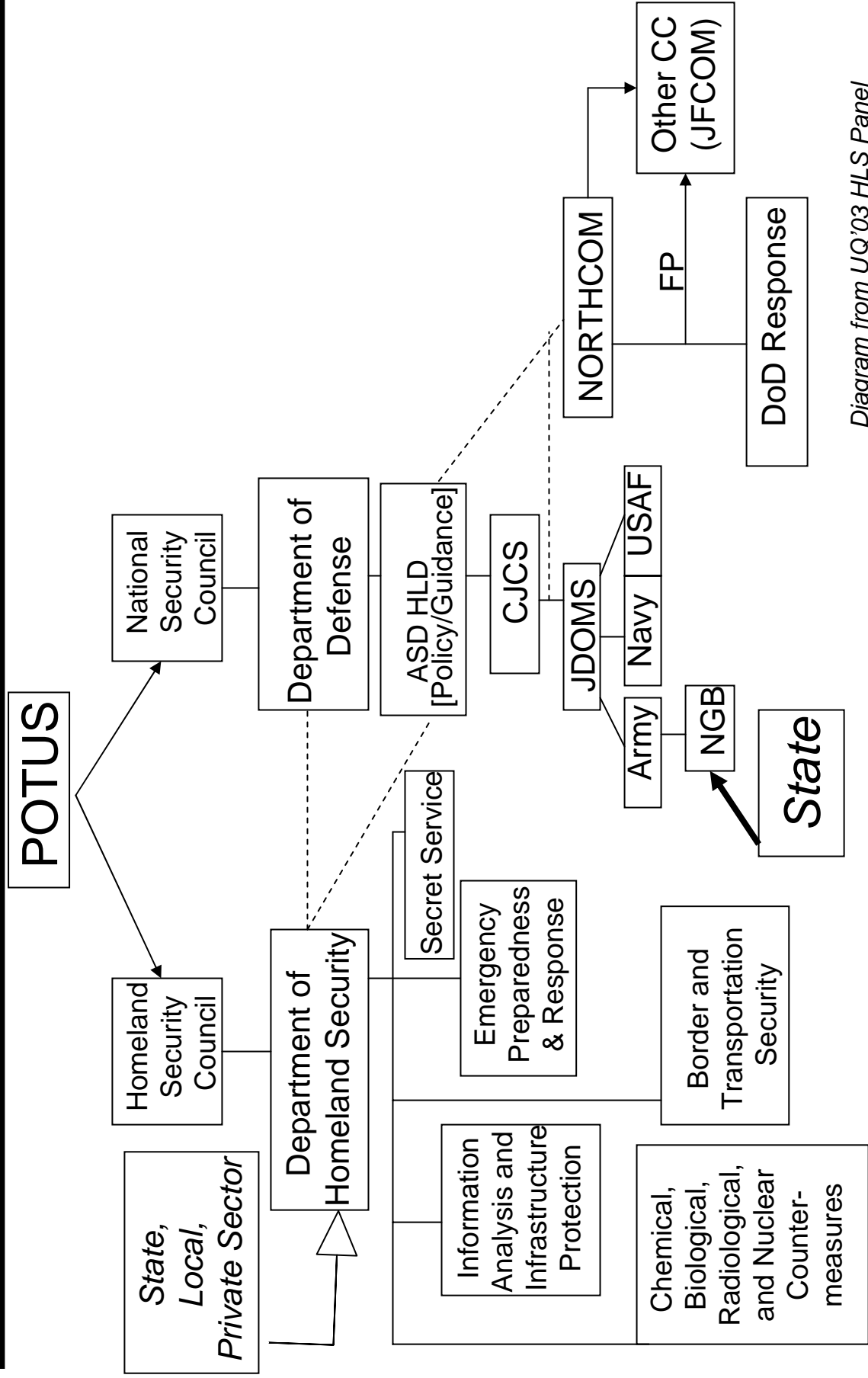


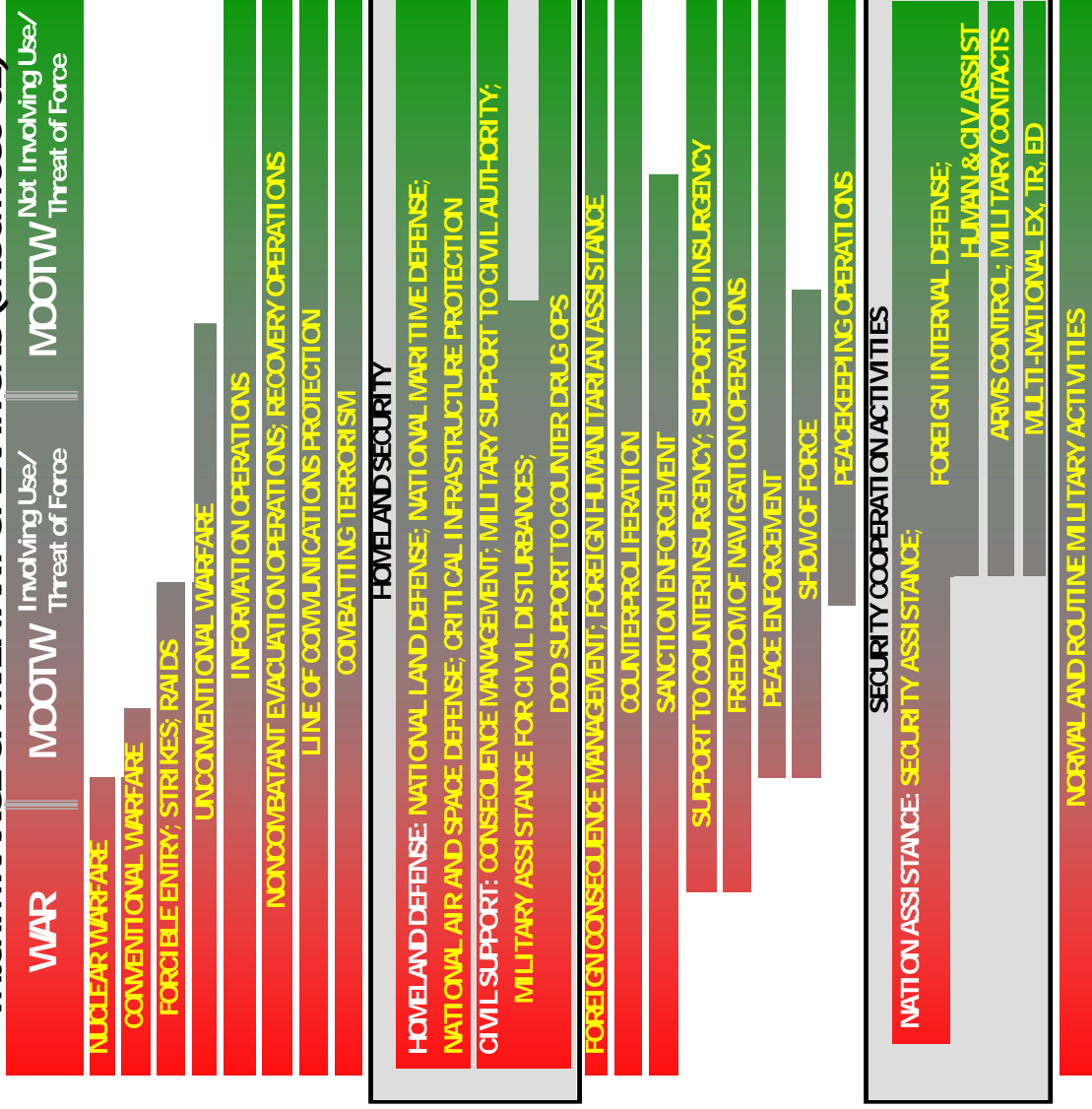
Diagram from UQ'03 HLS Panel
Senior Mentor

Responsibilities of DoD in HLS

- **DoD is responsible for:**
 - Defense of the US.
 - DoD capability support to the civilian authorities.
- **Area of responsibility is US, Canada, Mexico and the land, sea, and aerospace approaches. NORTHCOM and PACOM are to oversee:**
 - Military support to Domestic Consequence Management.
 - Support and protection of the Critical Defense Infrastructure in the Homeland.
 - Support Civilian Support operations (an Army core competency) (Army Modernization Plan 2003).
- **Homeland Security is best accomplished by building on state and local capability.**
 - Training and equipping state National Guards and militia.
 - Developing, coordinating, and cooperating with the First Responder community.

Defining the Range of Military Operations

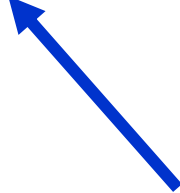
Interim RANGE OF MILITARY OPERATIONS (JROCM080-02)



Homeland
Defense



Civil Support



Jopsc
Version 4.0
Page 24

Responsibilities of DoD in HLS (2)

- **In the event of national need, DoD will have a role in three broad circumstances:**
 - **Limited in scope/temporary in time, requiring DoD to assist/train state and local actors (DoD supports):**
 - Special events.
 - Training First Responders.
 - Support to law enforcement.
 - **Emergency, requiring DoD to augment capabilities of civil authorities (DoD Supports):**
 - Post-event management.
 - Logistics, supply, mobility.
 - **Extraordinary, requiring DoD unique capabilities (DoD leads):**
 - Combat air patrols (CAP) around high-value sites (NYC, DC, etc.).
 - Explosive ordnance detonation (EOD).
 - Interagency augmentation using combat forces.

Analysis of DoD Support to HLS in Temporary and Limited Circumstances

- **Temporary in time/limited in scope, requiring DoD to assist/train state and local actors:**
 - **Special events**
 - Represent venues such as the Boy Scout National Jamboree or the Super Bowl and use simulation to develop plans, and contingencies to enhance control and mitigate manmade or natural disaster.
 - **Training First Responders**
 - Use simulation tools to develop training scenarios based on operational plans and contingencies.
 - Simulations like the Emergency Preparedness Incident Command Simulation (EPICS)** to train the command element of First Responders and their interface to other agencies.
 - **Support to law enforcement**
 - Assist the Border Patrol and area law enforcement to plan sensor arrays.
 - Assist law enforcement in developing contingency plans and asset use plans through the use of models and simulation. (As an example, a version of SimCity suitable for realistic plans and contingency development sufficient in detail to represent applicable C4ISR.)

**EPICS is an HLS command and control rehearsal model developed and maintained by TRAC-WSMR.
POC is Dr. Julie Seton (505) 678-4949

Analysis of DoD Support in HLS – Emergency Circumstances and Situations Analyses

- **Emergency situations, requiring DoD to augment capabilities of civil authorities, include:**
 - **Post-event management**
 - Effectiveness of means and methods to control and isolate the forward range boundaries of effect of WMD/WME.
 - Identify means to open and secure APOD/SPOD and protect the borders and littorals from insurgents.
 - **Logistics, supply, mobility**
 - What are the underlying assumptions in place regarding the availability of transportation and stores of emergency medical supplies like vaccines and antiagents?
 - What protection means are required to secure lines of communication within the affected site?
 - What effectiveness is provided by available complementary capabilities provided by private organizations, other agencies, or other governments?
 - **Support civilian support operations (an Army core competency (*Army Modernization Plan 2003*))**
 - Military support to Domestic Consequence Management.
 - Support and protection of the Critical Defense Infrastructure in the Homeland.

Analysis of DoD Support in HLS –

Extraordinary Circumstances Planning, Experimentation, and Analysis

- **Extraordinary, requiring DoD unique capabilities – use of combat forces within the US:**
 - **Combat Air Patrols**
 - Evaluation of materiel, sensor, and communication systems for target detection, control, and reduction.
 - **Explosive ordnance detonation (EOD)**
 - Projection of areas of influence, areas of effect.
 - Evaluation of contingency plans based on risk assessment of potential target areas, and means and methods of risk mitigation.
 - **Interagency augmentation and protection using combat troops**
 - Augmentation of US Customs Service, Border Patrol, and Immigration and Naturalization Service.
 - Evaluation of security vulnerabilities for sites of interest such as military bases, ports, airports, lines of communication, embassies, and chemical and nuclear storage facilities.
 - Evaluation of the effectiveness of sensors, lethal and less-than-lethal weapons.
 - Evaluation of joint force contribution for desired effect.
 - DoD and interagency force projection to detect, identify, preclude, preempt, or prevent a terrorist attack on the US, its territories or allies.

Introduction to HLS Strategy

- **The discussion above was the roles and expectations of DoD in Homeland Security/Homeland Defense or more generally, Homeland Operations.**
- **Next, the three phases for which a strategy must be planned:**
 - **Prevention**
 - **Immediate Response**
 - **Reaction and Recovery**

Strategy Planning Phase -- Prevent Terrorism

Pre-attack, anticipatory response

- **Preemption**
 - Strike at terrorist operations, infrastructure, motivations.
 - “Follow the money.”
- **Protection**
 - Physical security to harden targets, limit damage.
- **Preparation**
 - Blunt psychological impact.
- **Prevention**
 - Even one success will save lives and property.

Strategy Planning Phase -- Prevent Terrorism

(Continued)

America's Critical Infrastructure Sectors

- **Critical Sectors**
 - Agriculture
 - Food
 - Water
 - Public Health
 - Emergency Services
 - Government
 - Defense Industrial Base
 - Information and Telecommunication
 - Energy
 - Transportation
 - Banking and Finance
 - Chemical and Shipping
- **Department of Homeland Security is working with federal departments and agencies and regional, state, and local agencies to develop national infrastructure protection plans. Modeling and simulation tools will be developed to understand how the complex and connected infrastructure behaves.**

Strategic Response – Attack Executed

Trans-Attack – Immediate, preplanned response

- **Crisis Management:**
 - End terrorist attack. End further damage.
- **Consequence Management:**
 - Rescue survivors, limit damage.
 - Coordinate First Responders.
 - Initiate recovery and reconstitution.
 - Defend domestic and defense industrial base.
 - Defend National Security infrastructure.
 - Defend, support, and restore civil operational structure.
- **Mitigation – turn from defense to offense via planned response.**

Strategic Response – Attack Executed

(Continued)

Strategic Response -- Post Incident

Post Attack – Reactions and Recovery -- Tempered by laws and society

- Threat interdiction:
 - Defeat terrorist infrastructure.
 - Again, “Follow the money.”
- Attribution.
- Recovery from natural and man-made disasters.
- Free-flow of Information.
- Respond with civil support:
 - Protect
 - Recover
 - Calm the public and restore confidence
 - Re-energize
- Form the basis of long-term strategy:
 - Thwart terrorists
 - Prepare for natural disasters

Strategic Response -- Post Incident

(Continued)

2003 Army/JFCOM HLS War Game Experiment Guidance

- **Collaborative Effort**
 - Improve learning potential.
 - Create venue for innovation.
 - Examine Joint operational concept.
 - Stress strategic capabilities.
- **Coordinated Objectives and Research Issues using 2003 Army Objectives and 2003 JFCOM Objectives**
- **Common Global Scenario**
 - Scenario linked via economic and social-political interests, energy concerns, international terrorism, and international crime.
- **Game Inputs**
 - Thinking intelligent Threat in future strategic and operational environment.
 - Joint, interagency, coalition players.
 - National Security and National Military Strategies.
 - Coordinated Joint campaign plans.
- **Outcomes**
 - Strategic messages disseminated through Army/JFCOM produced insight pamphlets:
 - Wargame Insight Briefings
 - Integrated Analysis Reports
 - Issue Focused Articles

HLS War Game Purpose and Scenario

- **2003 Experimental Purpose:** In a Seminar War Game setting, explore the capabilities required for an integrated and comprehensive Homeland Security effort using local, state, regional, national, joint, and service competencies while conducting a global campaign.
- **It is 2015. Assume:**
 1. **Transnational crime funding and terrorist interests continue to support attacks on US interests.**
 2. **No major change in the nature of the nation state system has taken place. However, non-state actors are increasing in power.**
 3. **Alliances and coalitions are expected to be more complex and dynamic.**
 4. **Despite major increases in economic interdependency and the impact of information technology, there has been no basic change in the nature of the global economic system.**
 5. **No major changes have taken place in preexisting international treaties, agreements, and organizations, except as were specifically stated in the game scenario and materials provide during the game's conduct.**
 6. **Local, state, regional, and US Government agencies and organizations continue to exist in the game's timeframe, and are structured in accordance to Presidential directives.**

War Game Events

- **Prepare/Defend**
 - Intelligence confirms specific threats to the US requiring active protection.
 - Special purpose forces intend to attack seaports and airports to disrupt US force deployment.
 - Terrorist attacks against shopping malls, amusement parks, and special events are scheduled to coincide.
 - Attacks against public utilities are also planned.
- **Response to Hostile Event**
 - An attack has been conducted against a port:
 - Explosion aboard a merchant vessel creates immediate casualties and potential for many more.
 - Port as well as airport operations are suspended.
 - State, regional, and DoD assets are required to provide coordinated and synchronized measures to restore the port commercial and operational area.
- **Respond to Non-hostile Event**
 - A major earthquake occurs along the New Madrid fault (Mississippi River) causing extensive damage from St. Louis to Baton Rouge.
 - State, regional, and DoD assets are required to take immediate steps to provide coordinated and synchronized measures to save lives, prevent further property damage, reduce suffering, and restore the viability throughout the Mississippi River basin.

Results: Concern Areas

- **Changing Threat Perspectives**
 - Concern for the Homeland because of the Threat preemptive strikes and reach.
 - Countering area denial (Peacetime Theater Engagement) may prompt Threat action.
 - Concern for the Homeland because of the Threat cyber attacks and other activities.
 - Concern for the Homeland because of the Threat affecting airports and seaports as well as the physical safety of the US population.
 - Secure APOE/SPOE operations are required for force deployment.
 - Population safety concerns may cause limitations in plans for high demand/low density assets.

Results: Operational Concerns

- **Integrated Global Operations – Regional issues are global:**
 - A global view to the challenges of deployment, employment, and sustainment (DES) is required, of communicating coherent messages through the media, and of public perceptions.
- **Concern for the state of the security of the homeland communication assets.**
 - Concern for the state of the security of the homeland because that is where the APOE/SPOE.
 - Concern for the sustainment base in the defense industrial operations.
- **Concern because there may be a change in public and political intent in support of a war effort.**
 - The Threat may use actual or pseudo attacks in an attempt to change public support, or impart heightened concern due to the use of ports as DES gateways.
 - Public concern could slow port transiting, or even require ports to be operated with federal support.

Challenge for Statisticians and Analysts

- **The role required of DoD in HLS will serve as a catalyst for the transformation of the Army and DoD.**
 - **How do the requirements for Homeland Defense parallel those projected for the Future Force of the Army?**
 - **How to determine the importance and utility of expanded situational awareness.**
 - **How to model and evaluate efficient, networked systems that can immediately and accurately direct forces to perform critical missions.**
- **The community of applied statisticians, which offer capabilities in experimental design, data collection, analysis and interpretation, and modeling and simulation, can be a critical partner in accomplishing these and the following tasks:**

Statisticians Role in DoD's HLS Experimentation

- **Plan, coordinate, and execute experiments:**
 - Support procurement decisions.
 - Evaluate tactics, techniques, procedures, and concepts of operations.
 - Evaluate new training processes and/or devices.
- **Participate in experiments, “wargames,” and exercises to enhance the value of the experiment by providing techniques and procedures that result in valid experimental insights and conclusions.**
- **Assist in design, execution, and analysis of experiments to evaluate:**
 - Game theoretic risk assessment.
 - Screening criteria.
 - New equipment procedures.
 - Crisis management techniques, procedures, and approaches.
 - Potential value of specific security, counterterrorism, or crisis management activities.
- **Develop models and simulations:**
 - Develop the means to validly represent concepts, equipment, and behaviors in a model or simulation.
 - Implement those representations in high resolution computer simulations that can be used for quantitative analysis, as a training driver, or to support experimentation.
 - Develop both training and analytic scenarios.
- **Serve as the HLS experimental designers, war gamers, survey/test developers, data collectors, and analysts.**

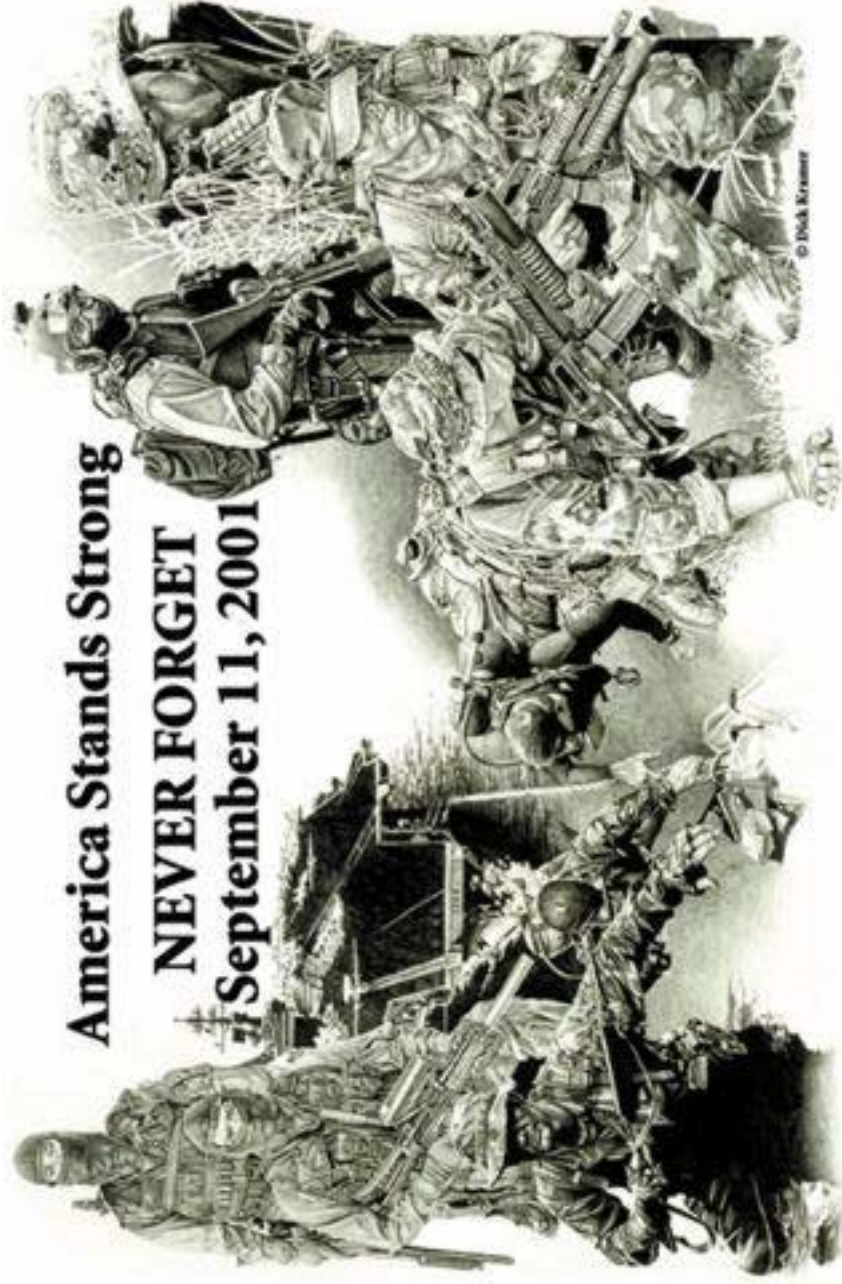
Studies & Analyses for DoD HLS Operations

- Determine the planning and operational requirements to support domestic disaster relief, antiterrorism, and consequence management operations.
 - In cooperation with HLD planners, conduct analysis of disaster relief, antiterrorism, and consequence management plans to flesh out plan details, validate the overall COA, and enhance the quality of those plans through operational plan (OPLAN) analysis.
 - Evaluate cost of alternate approaches to HLS.
 - Identify the need for materiel solutions necessary to accomplish the mission (analysis of requirements (AoR)).
 - Support the selection of optimum materiel solutions through the conduct of analysis of alternatives (AoA)).
 - Course of action and tactics, techniques, and procedures (COA and TTPs):
 - Assess the cost effectiveness of local or regional Homeland Security efforts against likely threats to those areas by conducting economic risk analysis.
 - Determine the cost impact resulting from heightened security measures (cost analysis).
 - Determine the importance and utility of expanded situational awareness.
 - Design efficient, networked systems that can immediately and accurately direct forces to perform critical missions.

Studies & Analyses for DoD HLS Training Development

- **Determine the training and exercise requirements to support domestic disaster relief, antiterrorism, and consequence management operations.**
 - **In cooperation with HLD planners, develop training and exercise scenarios of disaster relief, antiterrorism, and consequence management plan. Exercise the plan details, validate the overall COA and consequence anticipation.**
 - **Evaluate training system employed to train new HLS systems and/or process.**
 - **Evaluate cost of alternate approaches to HLS training and exercise. For example, the use of distributed simulation with a mix of human and virtual entities.**
 - **Develop the means to exercise developed COAs and TTPs:**
 - **Assess the effectiveness of training approaches, training events, and new equipment training through training impact analysis (training analysis).**
 - **Implement the representations of expanded situational awareness and networked systems.**

Questions?



Cyber Terrorism of Water Supply Infrastructure



MAJ John B. Willis
LTC Thomas M. Cioppa, Ph.D.
TRADOC Analysis Center (TRAC)
Monterey, CA

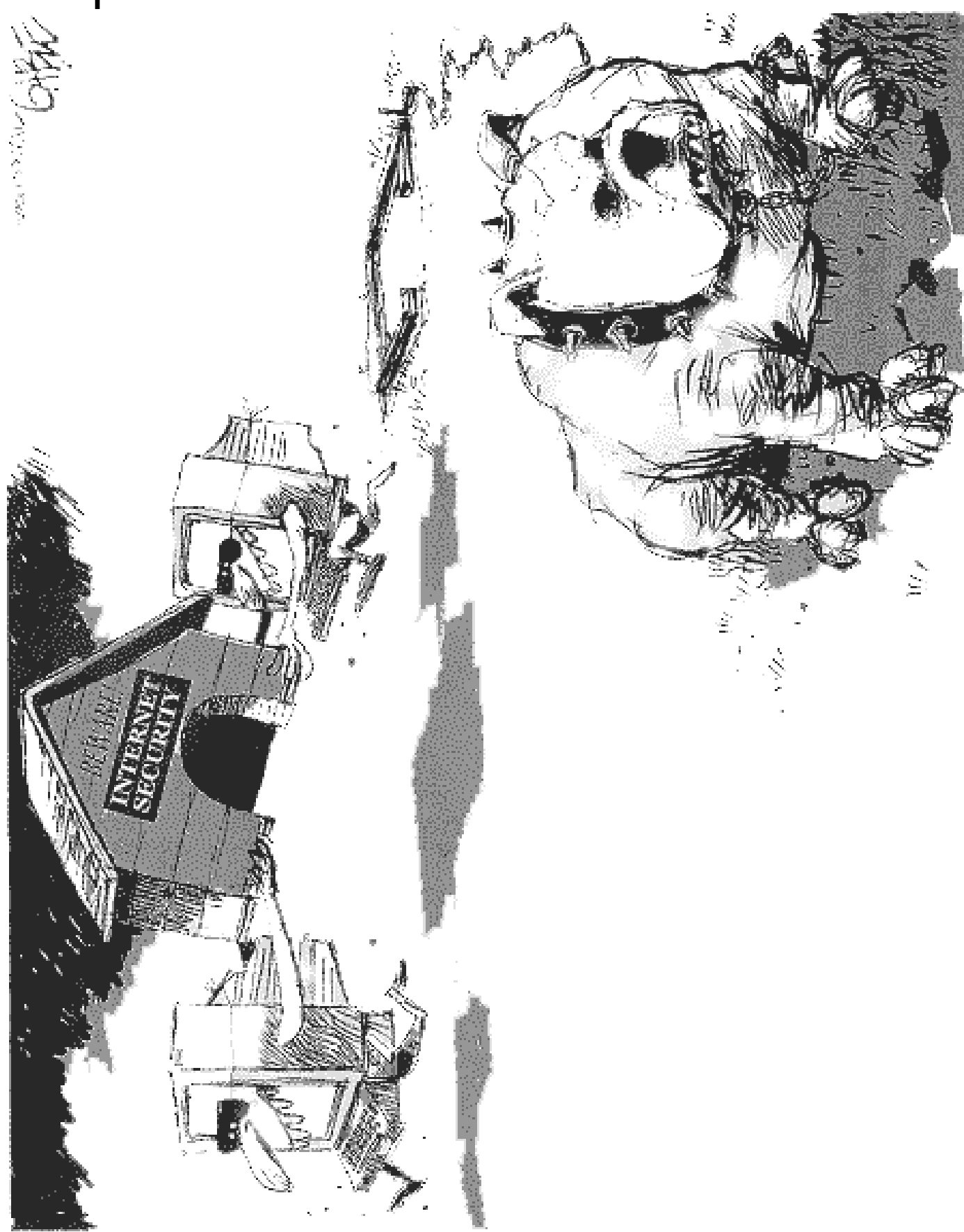


Cartoonist
TOMMY
STAR

OUTTA
THE WAY
...DUDE



MILO
www.milobooks.com



US water infrastructure faces large uncertainties in the character of threats and nature of system vulnerabilities



February 2001 Tabloid Headline

“URGENT! Last night, the FBI received a signed threat from a very credible, well-funded, North Africa-based terrorist group indicating that they intend to disrupt water operations in 28 US cities.

Because the threat comes from a credible, well-known source, with an organizational structure capable of carrying out such a threat, the FBI has asked utilities, particularly large drinking water systems, to take precautions and to be on the lookout for anyone or anything out of the ordinary.”

January 2001 AMWA E-mail to US Water Utilities

Time/CNN October 12, 2001 Poll

Do you think the following types of terrorist attacks are likely to occur in the US in the next 12 months?

Attack

At major public event
On some part of nation's water supply
Against the internet
On a nuclear power plant
Against another sky scraper

<u>Likely</u>	<u>Not Likely</u>
67%	29%
64%	32%
59%	32%
58%	37%
47%	50%

U.S. Forces discover files related to computerized water systems on al Qaeda computers in Afghan camps.

President's Critical Infrastructure Protection Board, Feb 2002

Cyber Terrorism of Water Supply Infrastructure

<p><u>Problem</u></p> <ul style="list-style-type: none">• President’s Commission on Critical Infrastructure Protection (PCCIP) study concluded that cyber threats are a clear danger (risk) to all infrastructures.• Among these critical infrastructures are the nation’s water supply systems.• Civilian water utilities support military installations and force projection	<p><u>Objective</u></p> <ul style="list-style-type: none">• Review/compare risk and vulnerability assessment methodologies• Conduct survey of water providers• Demonstrate value of these methodologies for military M&S
<p><u>Client</u></p> <ul style="list-style-type: none">• TRADOC HLS Directorate• TRADOC DCSINT HISTO• JFHQ-HLS NORTHCOM• ERDC Fort Future (Army STO)	<p><u>Deliverables</u></p> <ul style="list-style-type: none">• Vulnerability Assessment Analysis<ul style="list-style-type: none">• Initiated Oct 02 (complete)• Survey Analysis<ul style="list-style-type: none">• Web-based survey posted July 03• Technical Report<ul style="list-style-type: none">• Sep 03

Executive Summary ^(1/2)

1. Military/civilian leaders are responsible for **protecting our Nation's critical infrastructure**, communities, and symbols of national power from terrorists, home and abroad, as well as from natural disasters.
2. Public utilities support **military force projection** and DoD has a role in protecting critical infrastructure.
3. **Cyber risk awareness to SCADA systems** has increased significantly since 1996.
4. Internet-based attack trends indicate that the level of **sophistication in attacks is increasing**.
5. Consensus is forming that the **trusted insider/disgruntled employee** is more dangerous than other culprits of cyber attack.

Executive Summary ^(2/2)

6. Very little has been published in the way of rigorous **vulnerability assessment methodologies**.
7. Risk assessments are difficult to acquire because assessments are **proprietary or classified**.
8. Most risk assessments listed in the public domain are **soft system studies** relying almost exclusively on **qualitative measures and SME**.
9. A **quantitative systems-based risk assessment and management methodology** appears to provide the best approach.
10. **Military M&S** can support and benefit from risk assessment and management methodologies.

Motivation for Work

- **DOJ-NPS Interagency Agreement**
 - HLS Research and Technology Initiative
 - Focus on opportunities to strengthen U.S. capacity to deter, defeat and respond to threats to Homeland Security
- **TRADOC HLS ICT Charter (signed May 02)**
 - Requirement for Future Operation Capabilities that “...clearly address the Army’s Homeland Security requirements in the preparation, prevention, deterrence, preemption, defense, and response to threats and aggressions directed towards U.S. territories, sovereignty, domestic population and infrastructure...”
- **Army Homeland Operations (HLO) Concept**
 - Defines Army’s role in infrastructure protection and in defense against cyber attacks

Fort Future – Force Projection/Protection

- **Construction Engineering Research Lab (CERL), Engineer Research and Development Center (ERDC)**
- **Reliable utility systems are key to Force Projection and Protection**
- **Developing methods, tools, models to plan, assess, optimize, and monitor the ability of utility systems to support Army Force Projection**
- **Water utility applications**
 - **Pilot testing water dynamic system models**
 - **System vulnerability assessments**
 - **CBR contaminant scenarios**
 - **Analysis of system operation modifications based on real-time modeling data**

Definitions

- **Homeland Security (HLS)** is the prevention, preemption, and deterrence of, and defense against, aggression targeted at U.S. territory, sovereignty, domestic population, and infrastructure as well as the management of the consequences of such aggression and other domestic emergencies. Homeland security is a national team effort that begins with local, state and federal organizations.
- DoD and NORTHCOM's HLS roles include Homeland Defense and Civil Support.
- **Homeland Defense (HLD)** is the protection of U.S. territory, domestic population and critical infrastructure against military attacks emanating from outside the United States.

Critical Infrastructures

Government Operations



Water Supply Systems



Gas/Oil Systems



Banking & Finance



Electrical Energy



Transportation



Emergency Services



Telecommunications



Current National Situation

- **US is faced with a significant force projection challenges**
 - Southwest Asia
 - Eastern Europe
 - North Korea
 - Africa
- **Hostile groups continue to threaten all aspects of the critical infrastructure with a hybrid (asymmetric) attack (Physical and Cyber)**
- **Private industry understands the risks, but lacks a unified methodology and funding to harden their assets**

Major Agencies/Programs

- Critical Infrastructure Assurance Office (CIAO)
 - DHS
- Cybersecurity Tracking, Analysis and Response Center (CSTARC)
 - DHS
- National Infrastructure Protection Center (NIPC)
 - FBI
- Water Information Sharing and Analysis Center (ISAC)
 - EPA, AMWA, AWWA
- Computer Crime and Intellectual Property Section (CCIPS)
 - DoJ
- Critical Infrastructure Protection Program (CIPP)
 - GMU/JMU
- TRADOC HLS Directorate
 - US Army
- JFHQ-HLS NORTHCOM
 - DoD
- National Infrastructure Simulation and Analysis Center (NISAC)
 - LANL, SNL
- Homeland Infrastructure Security Threats Office (HISTO)
 - Critical Infrastructure Assurance Program (CIAP) and Critical Infrastructure Vulnerability Assessment Program (CIVAP)
 - LANL, NORTHCOM
- JTF-Computer Network Operations (“Cyber Army”)
 - DoD
- Terrorist Threat Integration Center (TTIC)
 - DHS, FBI, CIA, DoD

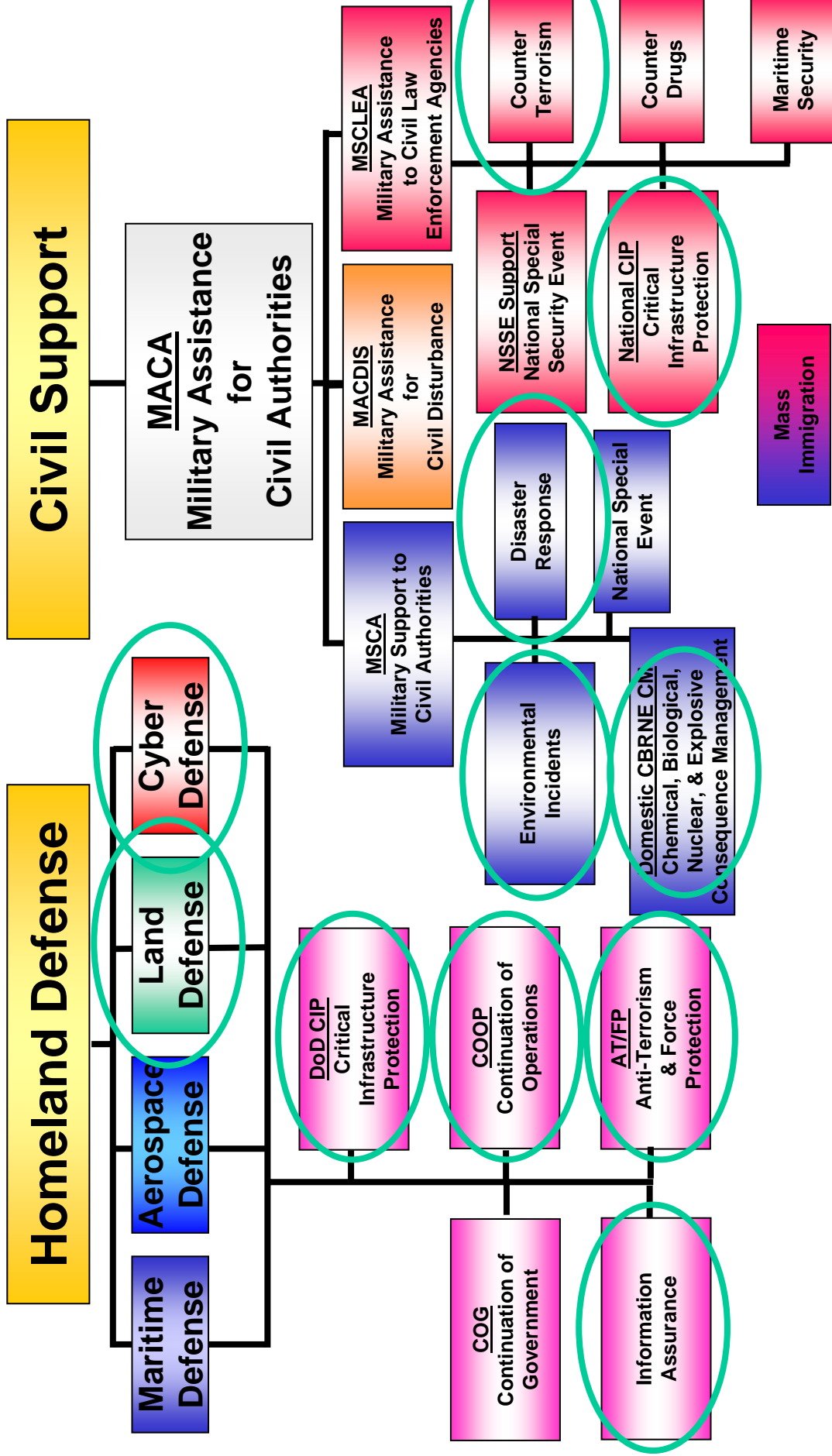
Stakeholders

- **General Public**
 - Expect uninterrupted flow of water service
- **Water Utility Companies**
 - Provide services that are potential targets of attack
 - ~160,000 public drinking water systems in US
 - ~370 systems serve over 100,000 customers each
- **Industry**
 - Designs utility systems and software
- **Government/Military**
 - Responsible for public safety/defense

DoD Focus

- **Secretary of Defense 2002 Annual Report and 2001 QDR**
 - Military forces need to be sized for defending the US
 - Homeland Security is DoD's primary mission
 - Reserve Component focus (e.g. National Guard WMD-Civil Support Teams)
 - Required capabilities/specific units – undefined
- **NORTHCOM**
 - Responsible for defending the US including:
 - Ocean approaches
 - Coastline
 - Seaports
 - Airspace
 - Assist civil authorities during emergencies within the US
- **Army**
 - Top Priority: Protecting military forces and their installations, embarkation ports/airfields, and information systems.

DoD HLS Key Functions



 = Link to Cyber Attacks vs. Water Infrastructure

Army Homeland Security Capabilities

- **Detection and Decontamination**
 - WMD-Civil Support Teams
 - Chemical/Chemical Recon and Decon
 - Biological Integrated Detection System
 - Technical Escort
 - Chem/Bio Rapid Response Team
- **Medical Services**
 - Medical Groups
 - Preventive Maintenance
 - Field Hospitals
 - Aviation-Evacuation
- **Perimeter Security**
 - Military Police
 - Infantry
- **Emergency Services**
 - Corps of Engineers
 - Quartermaster

Army Tasks for Critical Infrastructure Protection

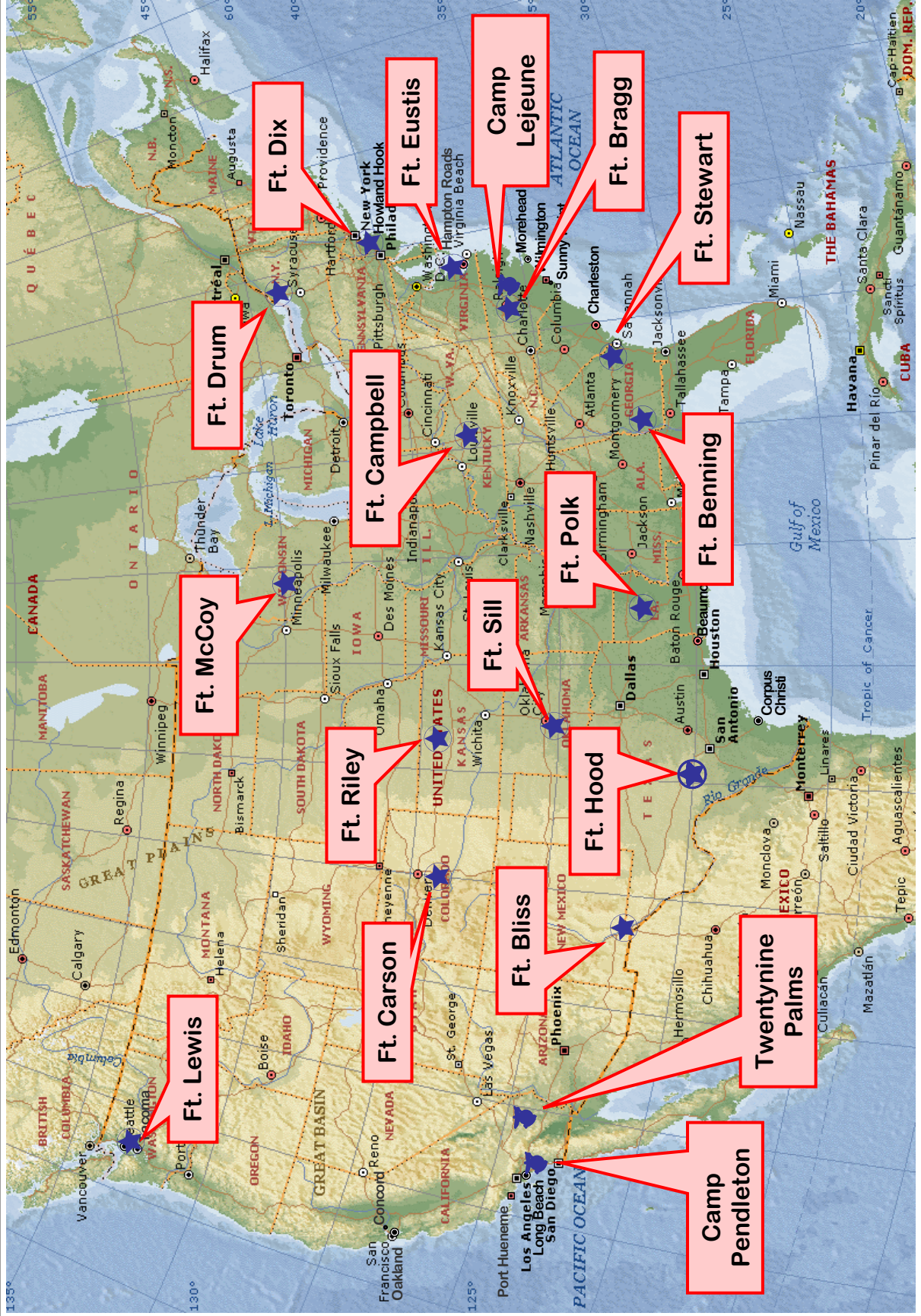
- **Ensure continuity of government operations by protecting facilities and personnel**
- **Reduce vulnerabilities of civilian physical infrastructure and information systems**
- **Provide area defense of critical infrastructure assets**
- **Consequence management**
- **Military presence to provide reassurance to American people**

Source: *The US Army and the New National Security Strategy*, RAND, 2003

Force Projection

- **Power Projection Platforms**
 - 15 Army Installations
 - 3 Marine Corps Installations
 - Navy and Air Force project power directly from home bases
- **Strategic Sea Ports**
 - 17 ports support Army and Marine deployments
- **Strategic Aerial Ports**
 - 17 airports support Army and Marine Corps deployments

Force Projection Platforms (18)



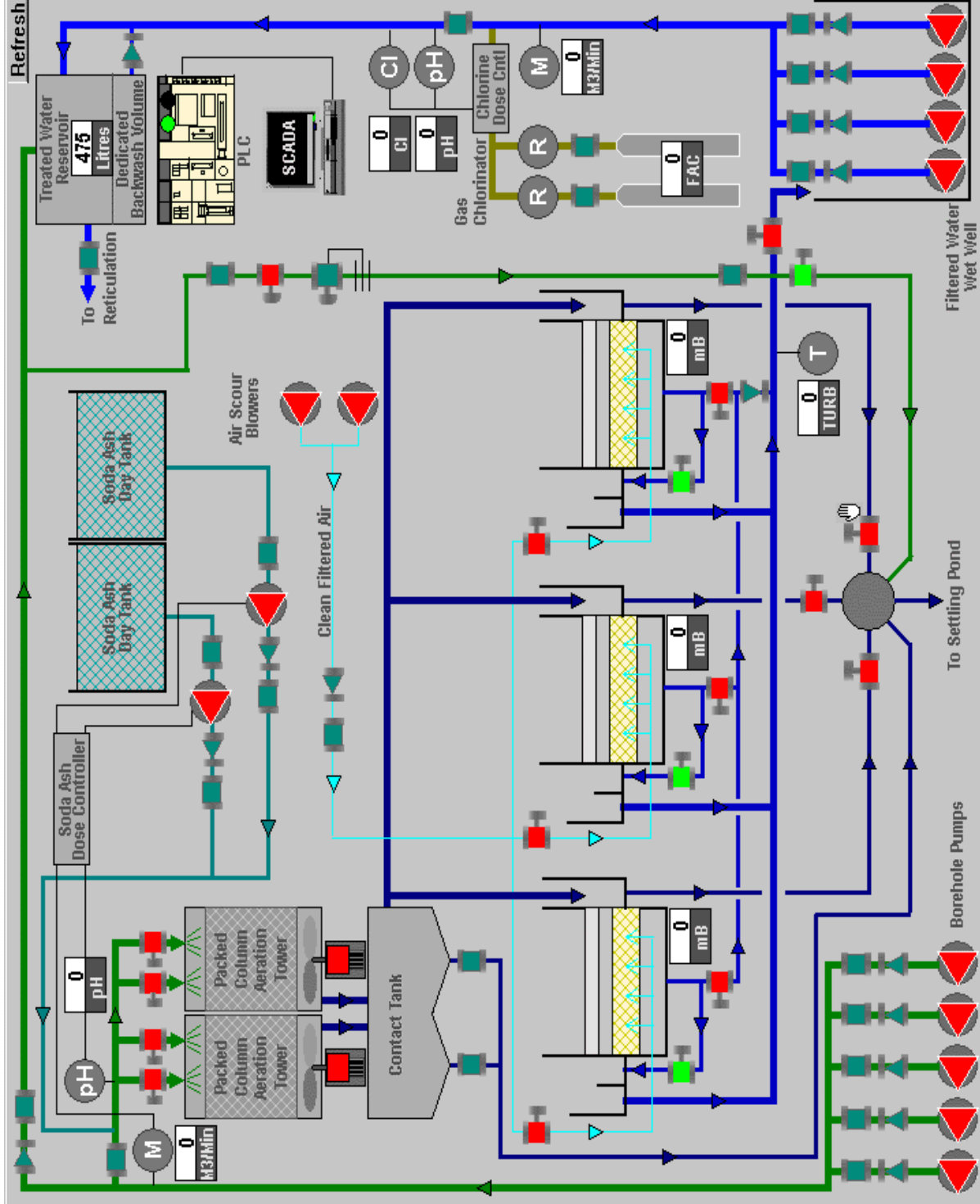
Water Supply Systems

- **Reservoirs; holding facilities**
- **Aqueducts; transport systems**
- **Filtration; cleaning systems**
- **Pipelines**
- **Cooling systems**
- **Waste water systems**
- **Firefighting systems**

Water Supply Control Systems

- **SCADA** – Supervisory Control and Data Acquisition
 - Uses MTU (Master Terminal Unit) and RTUs (Remote Terminal Units)
 - Open-loop; long distance
- **DCS** – Distributed Control System
 - Uses PLC (Programmed Language Controllers)
 - Closed-loop; local area network (LAN)

SCADA Display



Refresh

Nayagi Water Treatment Plant & Borehole Intake

System: **HEALTHY**

Status: **RUNNING**

Date/Time 08/01/1999 / 08:25:56

Alarm Information

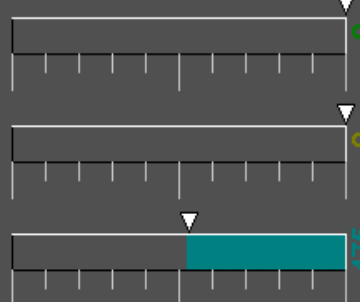
Plant Tools/Control

Reporting Information

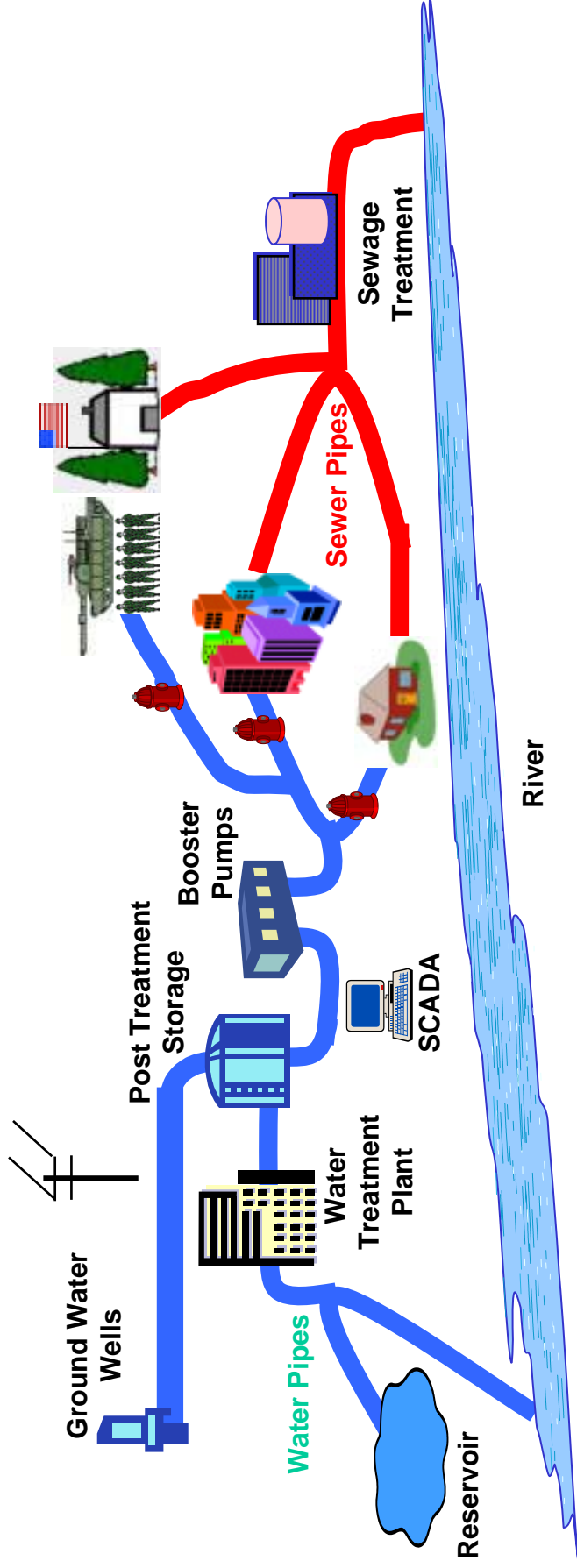
Security

Help

Engineering



Water systems are vulnerable to attack at multiple points: Risks and consequences vary by threat and location



Source → **Treatment** → **Distribution** → **Sewer/Treatment** → **Discharge**

Types of Attacks: Physical, Contamination, Cyber

Water System Vulnerabilities

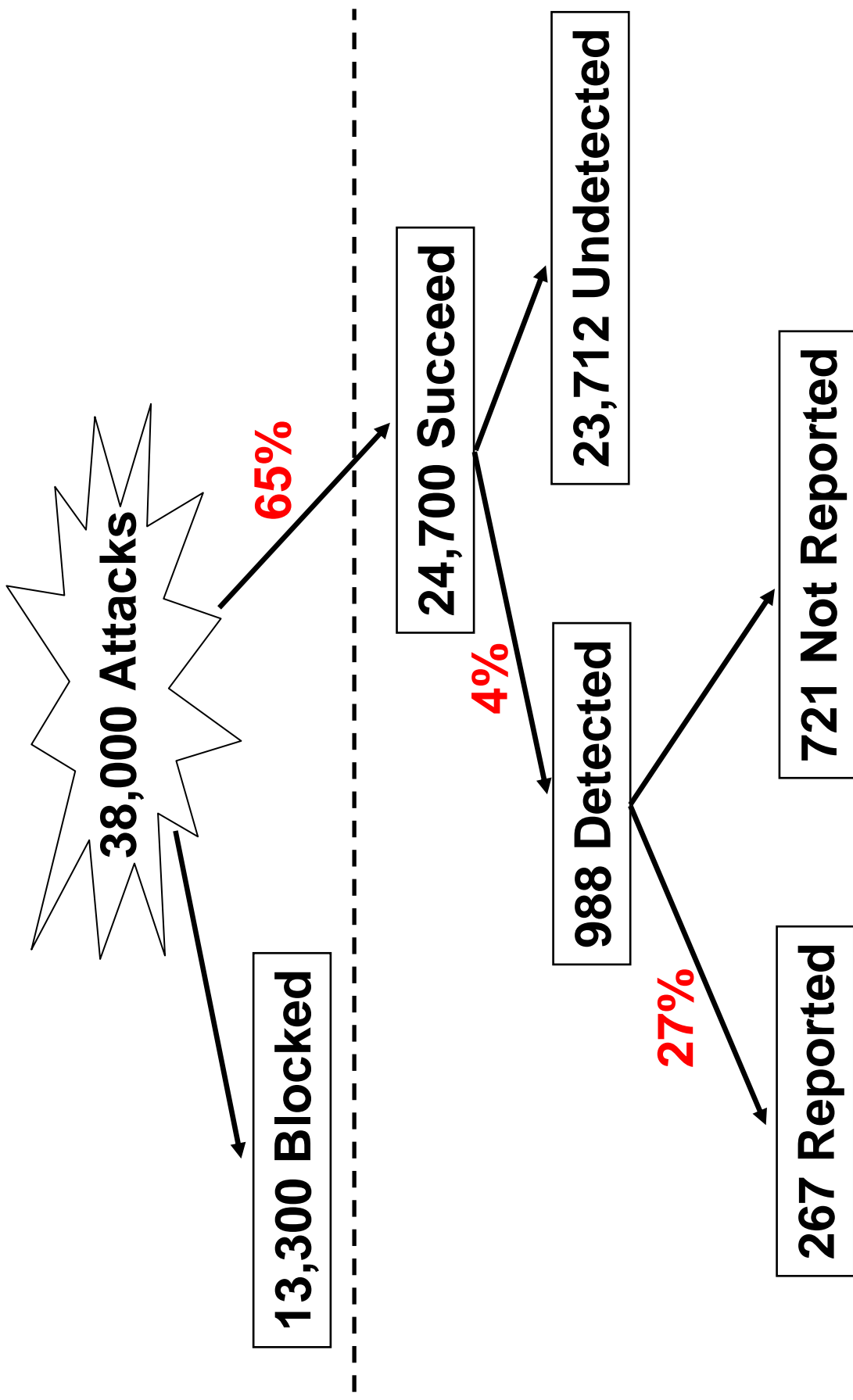
- **Source Water**
 - Drought, Flood, Contamination, Degradation
- **Treatment Plant**
 - Facility Attack, CBRNE
- **Pump Station**
 - Lack of Redundancy
- **Tankage**
 - Isolation
- **Distribution**
 - Hydrants, Valve Pressure Transients
- **SCADA**
 - Interdependency with Power, Interception, Lack of Encryption, Physical or Cyber Attack, Data Corruption

Vulnerabilities

- **Software and hardware weaknesses**
 - 40% of water facilities allow operators direct access via internet
 - 60% of water SCADA systems accessible by modem
- **Human weaknesses (e.g. training)**
- **Lack of a security culture**
 - Productivity vs. Security
 - Example: Power plant with all control systems set for access using the same password

Each of these vulnerabilities can be exploited to allow intruders unauthorized access to information systems, leaving the information or those systems subject to manipulation or other forms of attack.

1995 GAO Computer System Vulnerability Study
DISA (Defense Information Systems Agency)



Computer Security Breaches

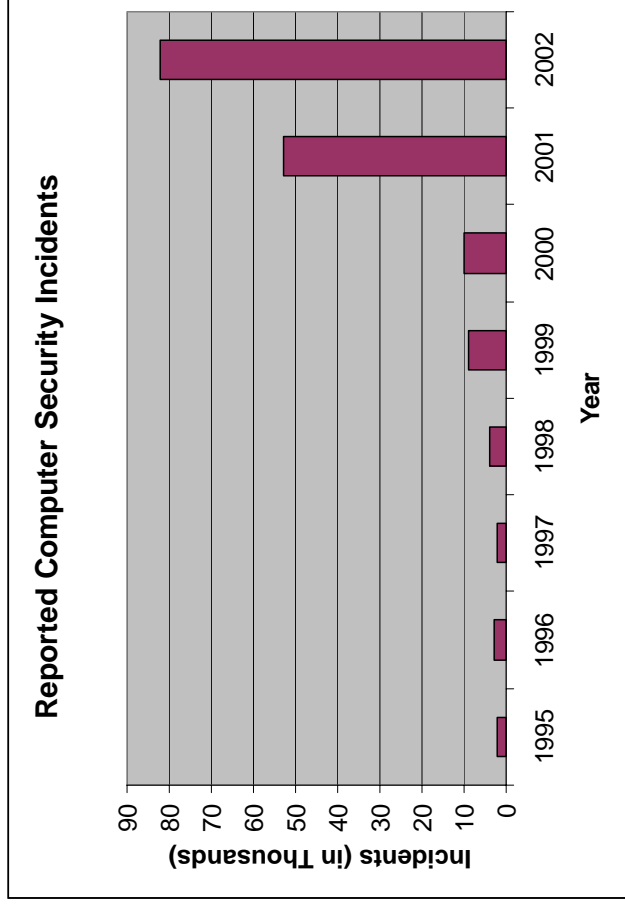
- **2002 Report: Computer Crime and Security Survey**

(Computer Security Institute and FBI's Computer Intrusion Squad)

- 90% of respondents had detected breaches

- **Reported incidents to CERT Coordination Center**

- 1999 – 9,859 incidents
- 2001 – 52,658 incidents
- 2002 – 82,094 incidents



- **80% of incidents likely go unreported**

- Not recognized as attack
- Reluctance to report

Two Forms of Cyber Attack

- **Against Data**
 - Stealing/corrupting data and denying services
 - Credit card number theft, web site vandalism, occasional denial-of-service assault
- **Against Control Systems**
 - Disable or take power over distributed control systems that regulate water, electricity, railroads, etc.
 - Systems are increasingly using the internet to transmit data or use LANs with “leaky” firewalls

Types of Cyber Attack

- **Denial of Service**
 - Preventing a computer from performing its function
 - Smurf, Spam, DNS redirect
- **Web Page Defacements**
 - Replace HTML files
 - Hacker war after US spy plane incident in China
- **Viruses/Worms**
 - Destructive code inserted and executed on victims' computer; worms are self-propagating
 - Nimda, Code Red
- **Trojan Horses**
 - Destructive code hidden within useful code
 - Badtrans – sends system and password info back to author
- **Network Intrusions**
 - Unauthorized network entry to gain “root” access
 - Gain network administrator functions

SCADA vulnerabilities

- Old software or standard vendor software with well known vulnerabilities
 - Hacker: “If you are running a Microsoft-based SCADA system, you have a target painted on your head. The volume of reading material on how to close the security holes in Windows NT is huge and the knowledge required to follow even the step-by-step instructions is very high and not the in the skill set of Joe Water Supply Company.”
- Can be penetrated without detection
- Can’t distinguish between attacks and system failure
- Can be manipulated remotely
- Could potentially inflict physical damage
- Courses on how to use control systems open to anyone with \$
- Psychological impact of attacks to critical infrastructures
 - Could affect the confidence of the population on a particular sector/activity related to SCADA systems.

Cyber Terrorists – Who are they?

- **Hactivists**
 - Politically motivated (e.g. Israeli vs. Palestinian hackers; website defacement)
- **Terrorist Organizations**
 - Web sites used for recruiting, fund-raising, target research
- **Foreign Governments**
 - Government trained hackers used to attack other nation's computer power
- **Individuals**
 - Hackers, Script Kiddies, Insiders

#1 Threat to SCADA

Are Cyber Attacks a Real Threat?

- **Skeptics:**
 - Cyber attacks do not have the shock effect sought by terrorists
 - Difficult to knock out infrastructure with cyber attacks; easier to bomb
 - Lots of money and sophisticated skills are needed for successful attacks
- **Believers:**
 - Millions of black boxes controlling infrastructure systems
 - Systems are now internetted (but were initially designed as stand alone systems)
 - Encryption programs and security culture needed now

Water Supply Survey

- **Control systems in use**
- **Access**
 - LAN, Dial-in, etc.
- **Threats**
 - Hackers, employees, terrorists, etc.
- **Attack tools**
 - User command, scripts, programs, etc.
- **Consequences of attacks**
 - Information corruption/disclosure
 - Denial of service



Critical Infrastructure Survey

Sponsor: The US Army Training and Doctrine Command Analysis Center (TRAC) conducts research on potential military operations worldwide to inform decisions about the most challenging issues facing the Army and the Department of Defense (DoD). TRAC serves our Nation's soldiers by helping to define and underpin the concepts, requirements and programs that enable our Army to be the best organized, equipped, trained and ready Army in the world. TRAC directly supports the mission of the Army's major command, the Training and Doctrine Command (TRADOC), to develop future concepts and requirements set in while also serving the decision needs of many military clients. TRAC develops scenarios of potential military operations set in the future for use in modeling and analysis and is a significant contributor to advanced modeling and simulation research and improved modeling methodologies in the military. TRAC's research in Homeland Security (HLS) supports the TRADOC HLS Directorate and the Integrated Concept Team (ICT) currently addressing HLS requirements including critical infrastructure protection. Additionally, TRAC is directing its efforts in support of the Army Homeland Operations Concept which outlines the Army's role in infrastructure protection including defense against cyber attack.

Survey Purpose: The purpose of the survey is to collect data in support of research regarding vulnerability assessment and quantification for critical infrastructure (telecommunication, water supply, electric power, natural gas, and hydroelectric power).

Background: A pilot survey was conducted in 1998 addressing the issues of infrastructure vulnerability. Much has changed since then. Results from this survey will facilitate the development of a systems-based vulnerability assessment methodology that allows critical infrastructure vulnerability to be assessed and quantified. For questions on this survey contact [Major Barry C. Ezell](#)

Section One (Administrative):

Name:	<input type="text"/>	Email:	<input type="text"/>
City:	<input type="text"/>	Phone number:	<input type="text"/>
State/Region:	<input type="text"/>	Infrastructure:	Other <input type="text"/>
Country:	Other <input type="text"/>	Job Description:	<input type="text"/>
Job Title:	<input type="text"/>	Do you want your administrative/demographic information to remain anonymous?	Yes <input type="text"/>

Survey of Water Utilities (1998)

- **47% believe the disgruntled employee is the number one concern followed by 13% for internet hackers**
- **41% spend less than 10% of time on system security. 51% spend no time**
- **Ten utilities (10 out of 50) reported attempted attacks, successful unauthorized access, or use of their system**
- **Corruption of information and denial of service were seen by respondents as the major concerns from a cyber intrusion**
- **17% did not know the number of valves and 11% were unclear regarding the number pumps they controlled**
- **39% felt their system was safe from unauthorized access and only 37% from unauthorized use**
- **55% agreed that the ultimate objective of an attacker is damage followed by challenge or status**

Survey of Water Utilities (2003)

Case Study – Monterey Peninsula

- **Water services provided by Cal-Am**
- **38,900 Customers**
 - **Monterey, Carmel, Pacific Grove, Pebble Beach, Sand City, Seaside**
- **Largest Customer Groups**
 1. **La Mesa Navy Housing Area**
 2. **Presidio of Monterey – Defense Language Institute**
 3. **Community Hospital of Monterey Peninsula**
 4. **Naval Postgraduate School**
- **Water Source: San Clemente Reservoir and 28 wells**
- **Components**
 - **59 tanks**
 - **49 pumping plants**
 - **Water treatment plant**

Case Study – Monterey Peninsula

- **Greatest perceived human threat:**
 - **Disgruntled employee or other insider**
- **Worst Consequence (physical attack):**
 - **Damage/contamination at treatment plant**
- **Worst Consequence (cyber attack):**
 - **Disrupting telephone line-based control system**
- **Scenarios studied in the past:**
 - **Earthquakes and floods**
- **Current Focus:**
 - **Installation of new SCADA system**

Cyber Security Efforts

- **DHS agencies addressing cyber security**
 - **Critical Infrastructure Assurance Office (CIAO)**
 - **National Infrastructure Protection Center (NIPC)**
 - **Federal Computer Incident Response Center (FCIRC)**
 - **National Communications System (NCS)**
- **In the summer of 2002, the U.S. EPA mandated that all community water utilities that serve more than 3,300 people complete vulnerability assessments by the end of June 2004.**
- **EPA provided \$51M in grants to assist water utilities in conducting vulnerability assessments and response plans**

Very little has been published in the way of rigorous vulnerability assessment methodologies.

- No agreed upon definition of vulnerability
- Vulnerability assessment guidance: ad hoc checklists
- SCADA should be viewed as system of a larger complex organization system.
- Implication: Classic risk assessment questions:
 - what can go wrong,
 - what is the likelihood, and
 - what are the consequences, should be preceded by the question:
 - **what is the system in focus**
- The system in focus and the context must be understood before meaningful risk and vulnerability assessment is undertaken.

Risk/Vulnerability Assessment Methodologies

- **Risk Assessment Methodology – Water (RAM-W)**
 - Sandia National Laboratory
- **CARVER + Shock**
 - Joint Pub 3-05.5
- **Infrastructure Risk Analysis Model (IRAM)**
 - ODU, Stevens Institute of Technology, Tek Soft

RAM-W Program

- **How to plan and prioritize for your assessment**
- **How to identify threats to your utility**
- **How to identify facilities and assets that need to be protected**
- **How to understand the consequences of the loss of an asset**
- **How to evaluate your system's effectiveness in preventing an attack**
- **How to assess your risks**
- **How to develop a plan to reduce risk through operational changes**

RAM-W Methodology

- **Risk Equation**

$R = P_a \times C \times (1 - P_e)$ where:

R = Risk

P_a = Probability that something will happen

C = Consequences if something happens

P_e = Effectiveness of the security and response mitigation system

R, P_a , C, and P_e all have values from 0 to 1

If P_a , C, or $(1 - P_e) = 0$, then there is no risk

RAM-W Methodology (cont.)

- **Identify System**
 - Population/customers served
 - Critical customers (hospitals, military/govt., high-usage commercial users)
 - High profile events attracting national attention (sports, conferences, etc.)
 - Location of water systems emergency command center
- **Assemble Risk Assessment Team**
 - Identify Lead Associate
 - Identify associates from management, operations, quality, engineering, loss control, and personnel.
- **Identify Facilities and Pressure Gradients**
 - Sources, pumps, plants, boosters, storage, mains, dams, distribution systems
- **Compile Information for each Facility**
 - e.g. for SCADA: hardware/software, users, access, redundancies

RAM-W Methodology (cont.)

- **Gather Facility Documentation**
 - Site plans, system maps, risk management plans
- **Identify System Demands**
 - Flow rates, gradients, critical levels
- **Weight Criteria Critical to Operations**
 - Water quality, service, critical customers
- **Rank/Weight all Facilities using Critical Criteria**

Threat Assessment

- Focus on “**Design Basis Threats**” (DBT) – Maximum credible threat against which a water system’s security and operational practices should be designed to defend
- **Intrusion**
 - Destroy/disable equipment, contamination, hazardous chemical release
- **Blast**
 - Destroy/disable facilities, hazardous chemical release
- **Cyber Threat**
 - Access SCADA/DCS and disrupt operations
- **Distribution System Contamination**
 - Toxic substances introduced through main or hydrant

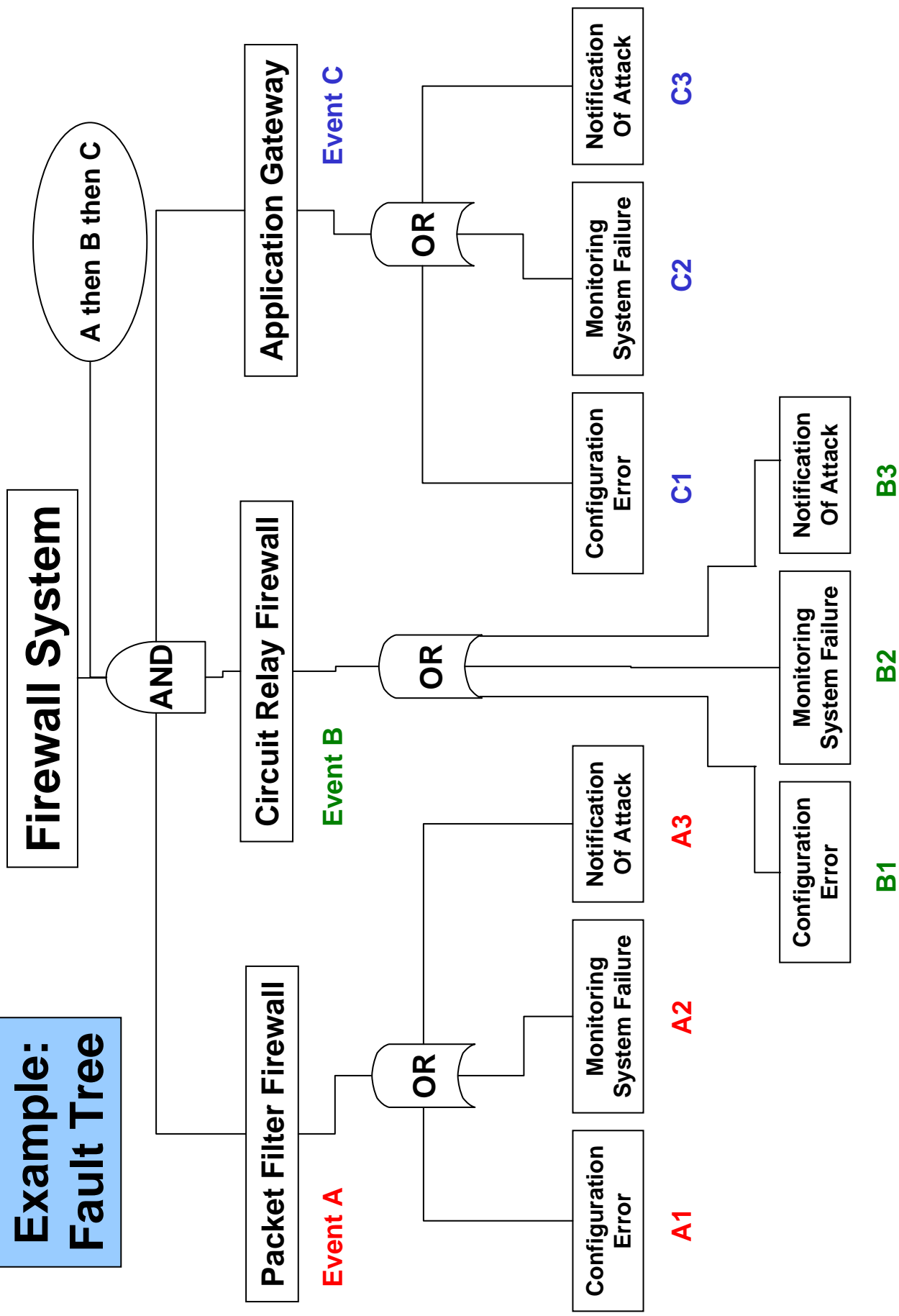
Facility Characterization (Calculating “C”)

- **Understand Redundancy and Reliability of Facility**
 - Multiple water sources, manual vs. DCS operation, emergency recovery means (generators, bottled water, etc.)
- **Create Facility Specific Consequence Table for each Threat Condition**
- **Prioritize and Weight Measures of Consequence**
- **Calculate Consequence (C) Value for Components and Overall System**
 - High: 1.0 – 0.7
 - Medium: 0.6 – 0.4
 - Low: 0.3 – 0.0

Security System Effectiveness (Calculating “ P_e ”)

- **Identify existing Physical Protection Systems (PPS)**
 - Fencing/gates, access control, lighting, alarms, sensors, guards, firewalls/passwords
- **Identify Operational Aspects of the System**
 - System storage/supplies, remote monitoring/control, emergency power
- **Create an Adversary Sequence Diagram (ASD) or Fault Tree**
 - Identify potential target assets, route/sequence of attack, centers of gravity
- **Perform Path Analysis**
 - Identify PPS(s) for the ASD that could detect/delay an adversary
 - Assign probabilities for each PPS to detect adversary and estimate delay times
- **Estimate PPS Effectiveness**
 - 0.0 < P_e < 0.2 Destruction and departure before response
 - 0.2 < P_e < 0.4 Destruction with arrival of response team
 - 0.4 < P_e < 0.6 Caught during destructive act
 - 0.6 < P_e < 0.8 Response before destructive act begins
 - 0.8 < P_e < 1.0 Response before intrusion

Example: Fault Tree



Risk Equation Analysis (Calculate “R”)

- Use values of C and P_e determined from analysis
- $R = P_a \times C \times (1 - P_e)$
- Determine if Risk is acceptable
 - $R > 0.75$ High Risk Facility
 - $0.75 > R > 0.50$ Med/High Risk Facility
 - $0.50 > R > 0.25$ Med/Low Risk Facility
 - $R < 0.25$ Low Risk Facility
- Identify Potential Improvements to Lower Risk
 - Prioritize (R x Critical Ranking)
 - Identify Improvements in Operations, Security, Response, etc.

CARVER + Shock

- **Criticality**
- **Accessibility**
- **Recuperability**
- **Vulnerability**
- **Effect**
- **Recognizability**
- **Shock**

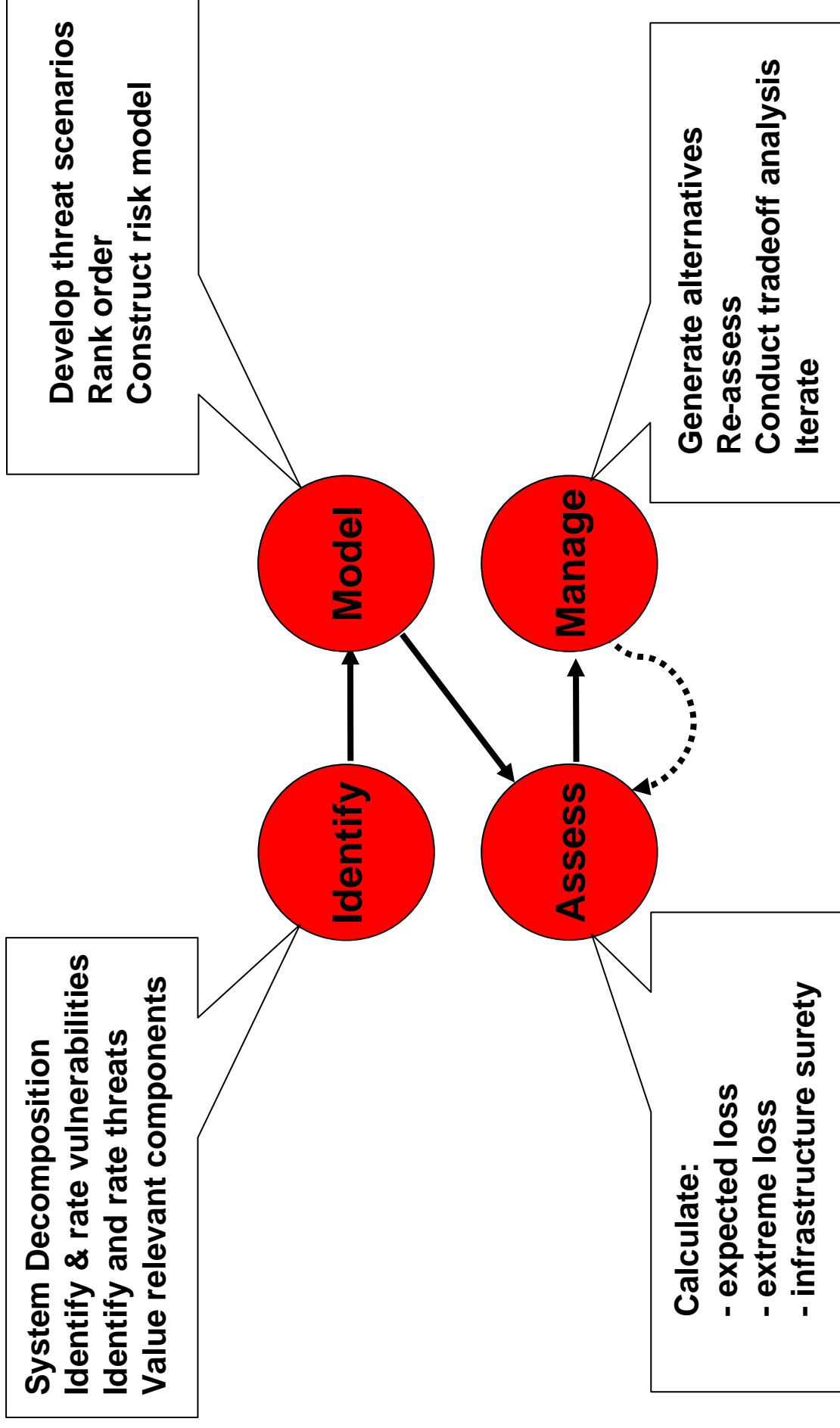
(Consider Asymmetric Threat, at a minimum cyber and physical)

CARVER + Shock Matrix

Target	C	A	R	V	E	R	Shock	Total
Intake Pump Station	7	3	6	2	3	9	6	35
Pump System	9	3	8	1	9	9	6	45
Source	5	4	4	3	3	9	6	34
Water Transfer	5	8	3	8	5	8	3	42
Mains	6	9	7	9	5	8	3	47
Backflow Valves	3	6	3	4	3	8	3	30
Water Treatment	9	9	9	10	9	10	8	64
Chemical Treatment	8	10	9	10	9	9	8	63
Control Center	10	8	10	9	9	10	8	64
Monitor/Control Center	9	7	9	9	10	5	9	58
SCADA/Switches	9	7	3	8	10	8	10	55
Computer Hardware	9	8	10	9	9	10	9	64

Ratings are 1 to 10 with 10 highest. Highest total scores annotate preferred targets.

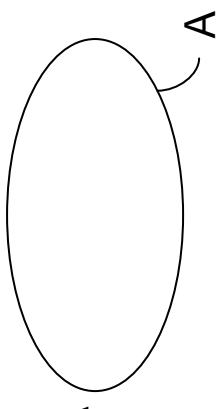
Infrastructure Risk Analysis Model (IRAM)



Risk Assessment and Management

- **What can go wrong?**
- **How likely is it?**
- **What are the consequences?**
- **What can be done?**
- **What are the tradeoffs in terms of all costs, risks, and benefits?**
- **What are the impacts of current decisions to future options?**

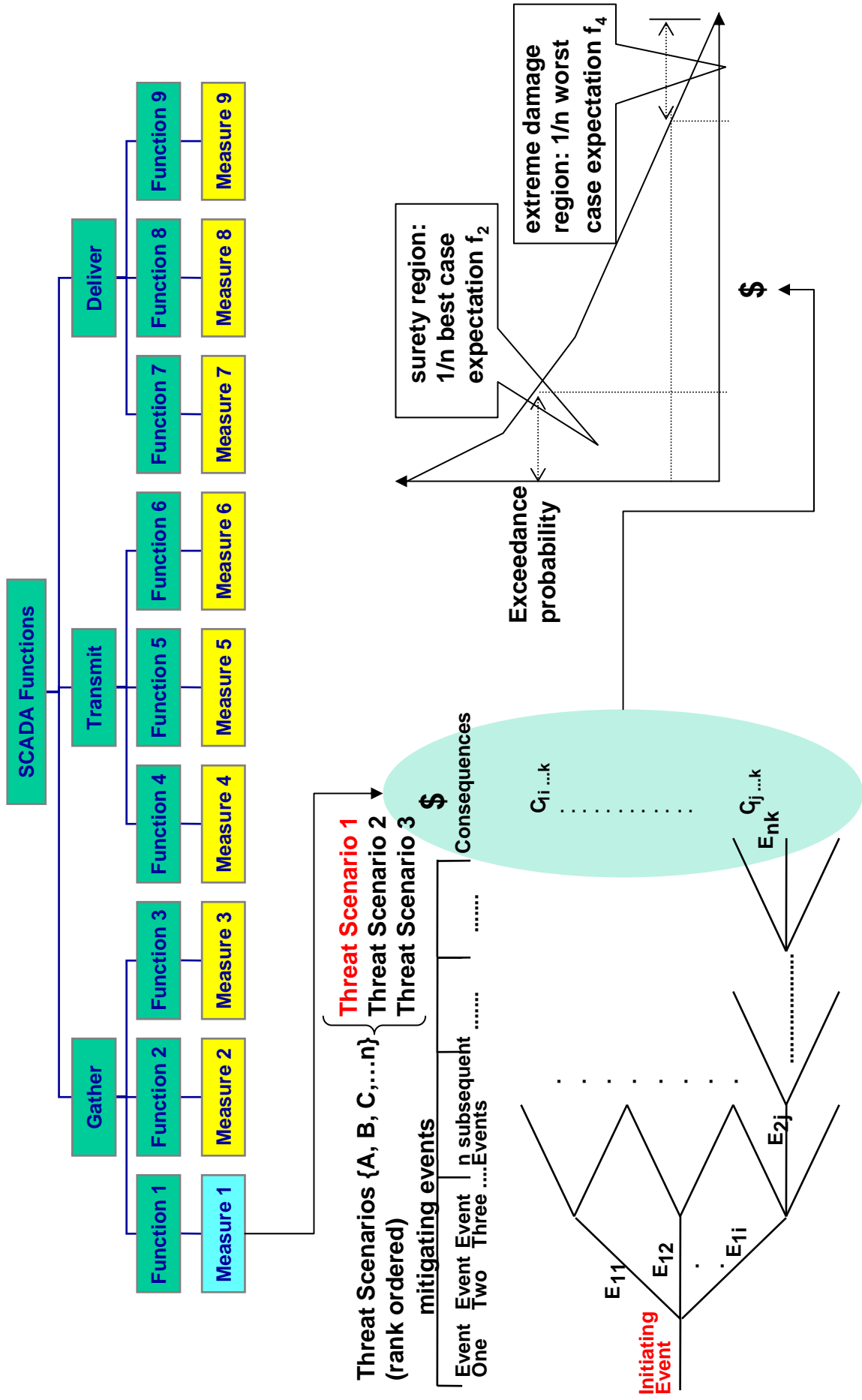
$$\text{Risk} = \{S_a, L_a, X_a\}_A$$



Kaplan and Garrick 1981. "On the quantitative definition of risk", Risk Analysis 1(1): 11-27

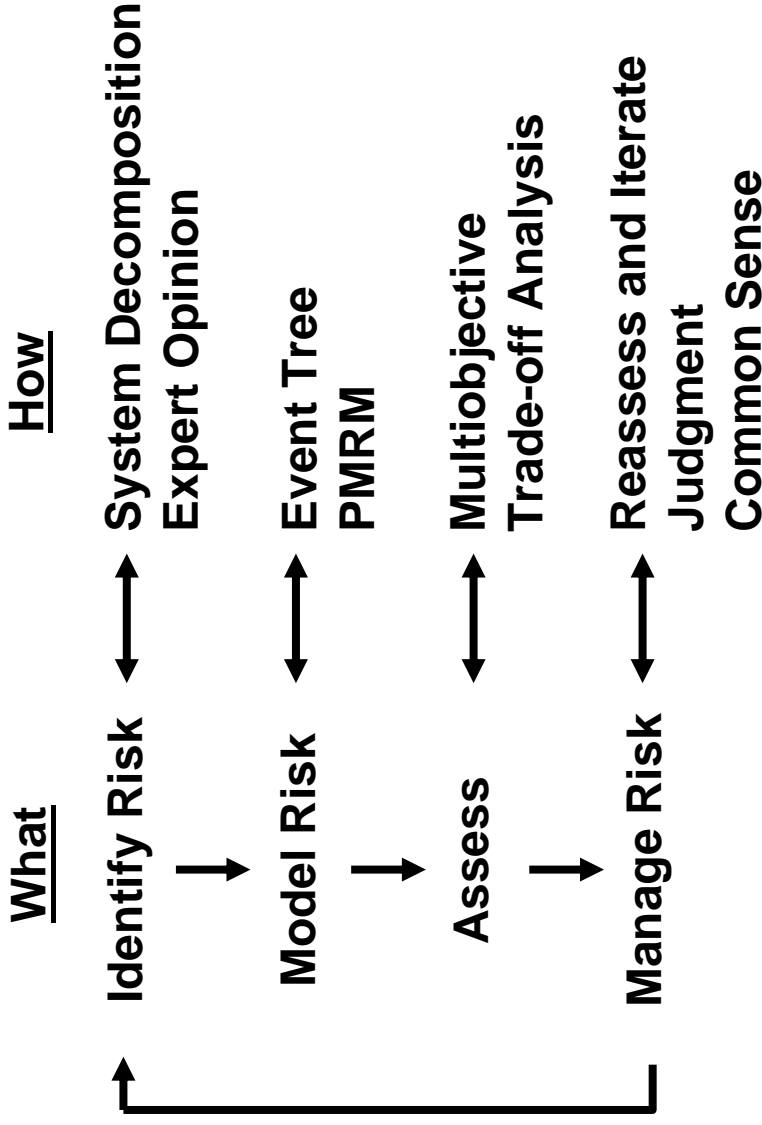
Haimes 1991. "Total risk management", Risk Analysis 11(2):169-171.

Dynamic multiple-objective decision analysis models incorporate risk and uncertainty by placing distributions on extreme event probabilities or on the weights assigned to value model evaluation measures. Component vulnerability is used to decide on threat scenarios.



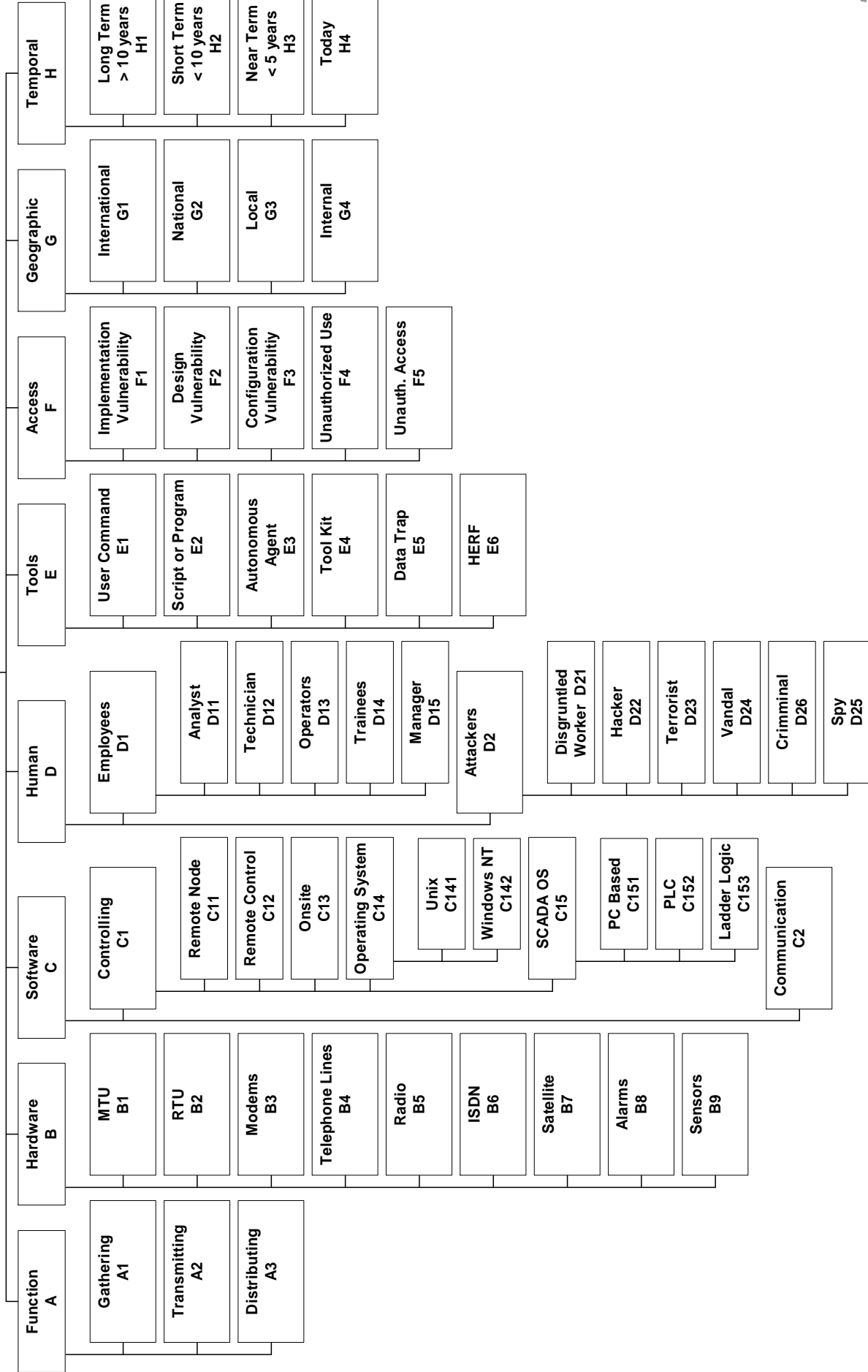
SCADA System Risk Management Framework

- Builds on existing probabilistic risk assessment (PRA) methodology
- May assist decision-makers in understanding cyber intrusion risk, consequences, trade-offs



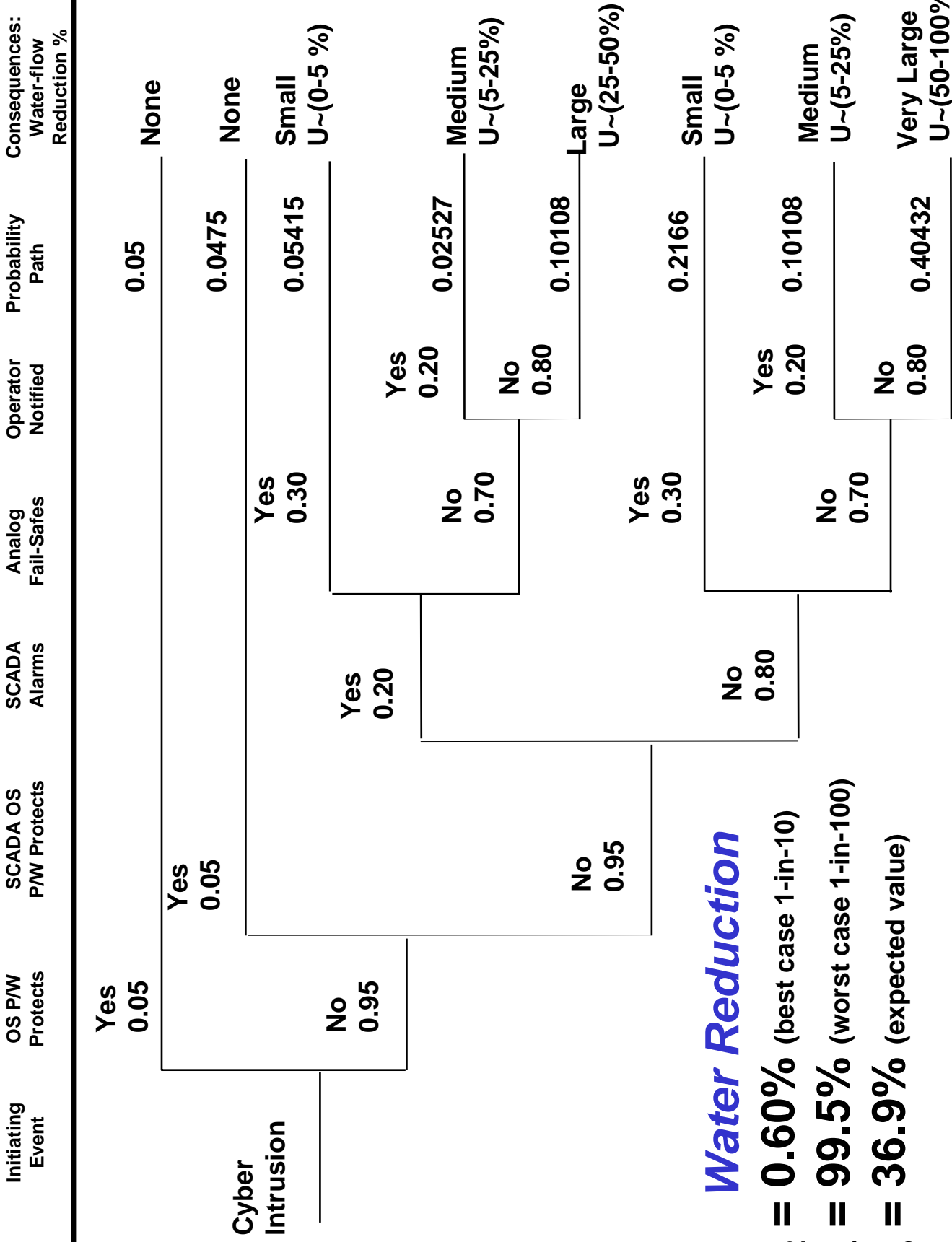
SCADA System Decomposition

SCADA
(Master-Slave)



Partitioned Multi-objective Risk Method (PMRM)

- Tool for quantifying/reducing risks of extreme events
 - Low probability and High consequence events
 - Event Trees used to develop probability functions
- Decision maker desires:
 - Expected % water-flow reduction in best 1-in-10 outcomes (f_2)
 - Expected % water-flow reduction in worst 1-in-100 outcomes (f_4)
 - Expected value of water-flow reduction (f_5)



Water Reduction

$f_2 = 0.60\%$ (best case 1-in-10)

$f_4 = 99.5\%$ (worst case 1-in-100)

$f_5 = 36.9\%$ (expected value)

Alternative Generation

- **Outsource web hosting**
- **Password sharing policy**
- **Filter firewall to isolate internal SCADA system from web server**
- **Configure call-back/logging features of dial-up modem**
- **Cancel access upon employee termination**
- **Token-based authentication**
- **Alarm suppression detection**
- **Alarms for unusual pump/tank usage**
- **Separate admin and operations servers**

Multi-objective Trade-off Analysis

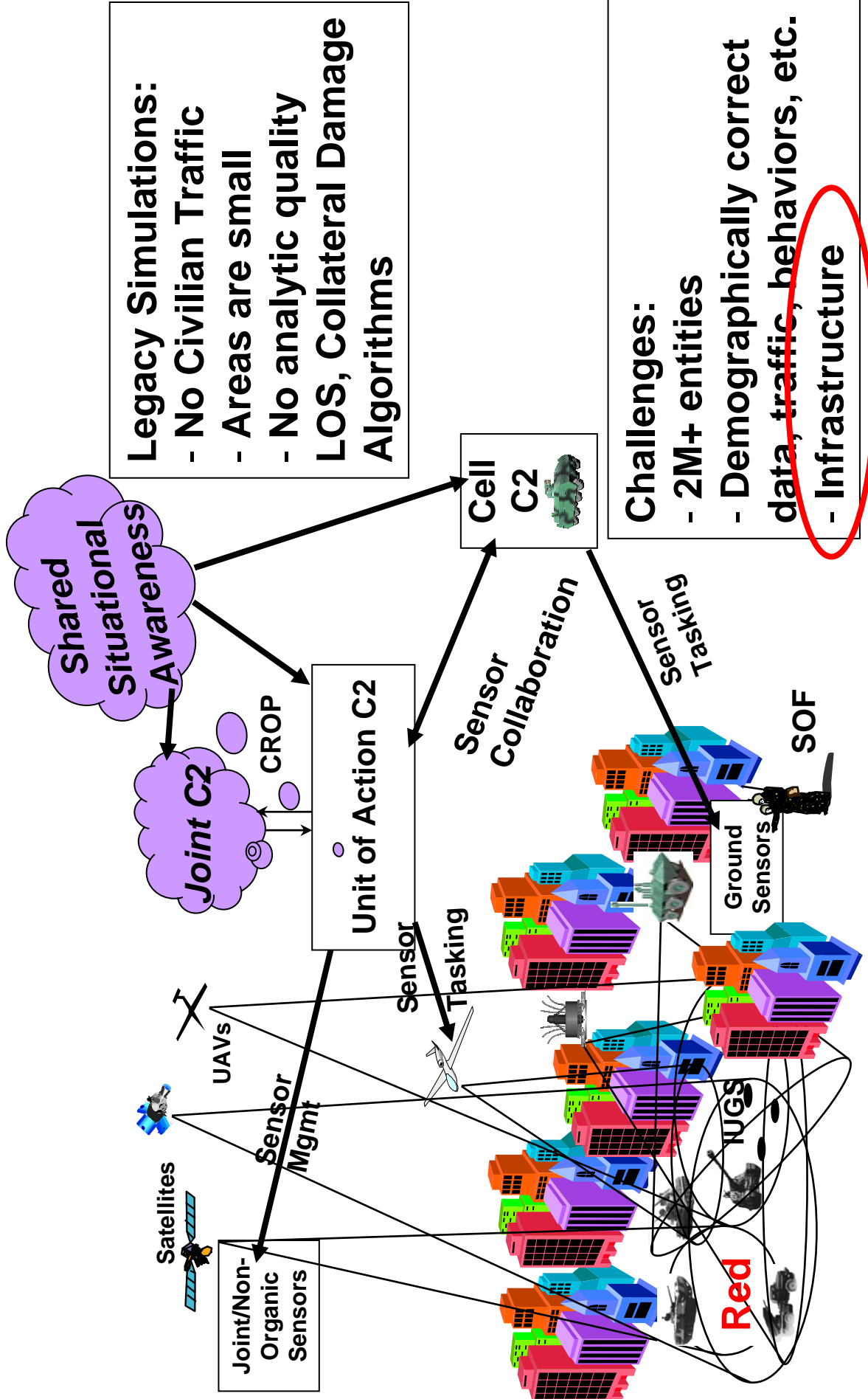
	OS Penetrated	SCADA Accessed	Alarms Defeated	Fail-safes (averted)	Operator Notified (paged)	Discount rate of 7% for 10 years (\$)
Alt 1	0.70	0.80	0.80	0.70	0.80	51,504
Alt 2	0.30	0.40	0.80	0.70	0.80	84,721
Alt 3	0.01	0.01	0.80	0.70	0.80	72,515
Alt 4	0.50	0.50	0.30	0.10	0.10	88,870
Alt 5	0.75	0.75	0.20	0.10	0.10	15,904
Alt 6	0.95	0.95	0.80	0.70	0.80	0.00

Applicability to Military Simulation

- **Reliability, Maintainability**
- **Survivability, Vulnerability**
 - **Human and Machine Failures**
- **C4ISR/Battle Command System Assurance**
- **Intelligence Analysis**
- **Risk Assessment**
- **Installation Infrastructure (e.g. Fort Future)**
- **HLS-Sim, AVERT, EPiCs**

Legacy Simulation Issues – Multi-Resolution

U.S. Joint Forces Command



Summary

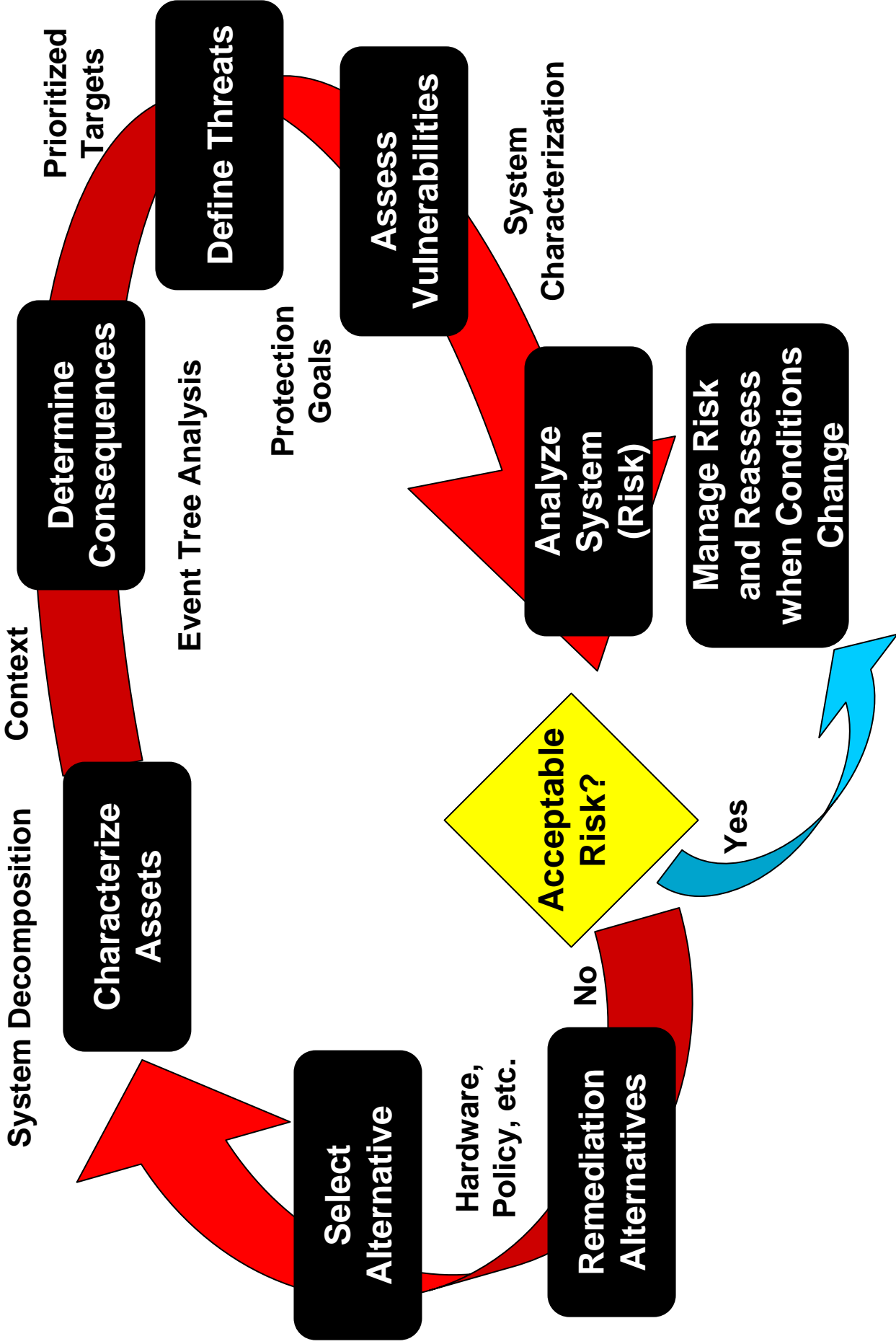
- **DoD has a role in critical infrastructure protection**
- **Public utilities support military force projection**
- **Cyber terrorism is an existing threat to utility SCADA**
- **Systems-based risk/vulnerability assessments will mitigate threats**
- **Military M&S can support and benefit from risk assessment and management methodologies**

QUESTIONS?



MAJ John B. Willis
TRADOC Analysis Center (TRAC)
Monterey, CA
john.willis@trac.nps.navy.mil
(831) 656-7580 DSN 756

Back-Up Slides



Consequences
Water Supply
Reduction

Probability
Path

Operator
Notification

Fail-safes

SCADA
Alarms

SCADA
Protection

OSP
Protection

Firewall
Protection

Initiating
Event

YES P1

P1

None

**Cyber
Intrusion**

YES P3

P2P3

None

NO P2

YES P5

P2P4P5

**Small
1-5%**

NO P4

YES P9

P2P4P6P7P9

**Small
1-5%**

YES P7

YES P13 P2P4P6P7P10P13

**Medium
5-15%**

NO P10

NO P14

P2P4P6P7P10P14

**Large
15-50%**

**Example:
Cyber Intrusion
Event Tree**

YES P11

P2P4P6P8P11

**Small
1-5%**

NO P8

YES P15 P2P4P6P8P12P15

**Medium
5-15%**

NO P12

NO P16 P2P4P6P8P12P16

**Large
15-50%**

- **Can't tell if the enemy has good weapons until he uses them (unlike counting bombers/tanks)**
- **“Swarming Attack” Scenario: Terrorist organization uses physical/cyber attacks on U.S. infrastructure combined with cyber attacks which disrupt the abilities of first responders (e.g. 911 system)**
- **Feb 02 Letter to Pres. Bush:**
 - ***“The critical infrastructure of the United States, including electrical power, finance, telecommunications, health care, transportation, water, defense and the Internet, is highly vulnerable to cyber attack. Fast and resolute mitigating action is needed to avoid national disaster.”***
 - ***Signed by 54 scientists, former national leaders, intelligence community recommending a Cyber Warfare Defense Project modeled on the Manhattan Project***

- **Recent cyber attacks**
 - **Legion of Doom – 89 – seized control of much of Southwestern Bell’s telephone network infrastructure. Could have tapped lines and shut down 911 service**
 - **Queensland, Australia – Mar 00 – man used the internet, a wireless radio, and stolen control software to release up to 1M liters of sewage into a river and coastal waters**
 - **Slammer – Jan 02 – worm that took down internet in South Korea and affected 911, airline, and banking systems in U.S. Exploited vulnerability in MS SQL Server 2000. 90% of damage in first 10 minutes**
 - **Nimda – 18 Sep 01 – worm attacked Wall Street and millions of computers by infecting email programs and slowing internet access**
 - **Moonlight Maze – Mar 98 – 2-year probing of DoD, DoE, NASA, universities, research labs traced to mainframe computer in Russia**
 - **Code Red – Jul 01 – worm that affected 300k U.S. computers and targeted the White House web site with denial of service attacks**
 - **Mountain View, CA – Aug 01 – “multiple casings of sites” emanating from the Middle East and South Asia looking for information on utilities, government offices and emergency systems**

- **Captured Al Qaeda computers/documents reveal:**
 - **Reconnaissance plans of U.S. critical infrastructure**
 - **Models of catastrophic dam failure**
 - **Information about digital switches used by power and water company system infrastructures**
 - **Ties to Inter Services Intelligence (Pakistani) which has contacts in various hacker groups**
 - **Use of sophisticated cryptography equipment**
 - **Use of one-time use email addresses**
 - **“Franchise-model” of independent partners**
- **Electrical control system engineer: “Worst case could be loss of power for 6 months or more.”**

- **Cyber Security History**
 - **Pres. Reagan first addressed the problem (Computer Security Act of 1987)**
 - **Pres. Clinton established PCCIP**
 - **Pres. G.W. Bush established PCIAB**
 - **Feb 03, Pres. Bush released *National Strategy to Secure Cyberspace***
 - **Implement public-private partnerships**
 - **85-95% of cyberspace owned/managed by private sector**

-
- **Awareness of infrastructure vulnerabilities**
 - **Precipitating event: Oklahoma City bombing**
 - **DoD “Eligible Receiver” exercise**
 - **NSA Red Team (hackers) used to attack Pentagon systems**
 - **Could only use publicly available computer equipment/software**
 - **Results – Infiltration and control of:**
 - **PACOM computer network**
 - **Power grids and 911 systems in 9 major U.S. cities**

- **Amit Yoran**
 - **Years from now, cyber attack will be a primary method of war; cyberspace a primary theater of operations**
- **Richard Clarke**
 - **“Red Teams” (government employed hackers) have succeeded every time in hacking into sensitive government computers, and gained total control of the networks involved, without the owner/operators even knowing it happened.**
 - **Focus needs to be on the fact that attacks are possible and not on who is doing it.**
 - **Osama bin Laden is not going to come for you on the Internet.**
 - **Cyber disruptions could be similar to anthrax attacks or DC area sniper.**

- **O. Sami Saydjari**
 - **Some systems use “honeypots” which are fake systems that don’t have any critical content but have interesting keywords/content that might attract cyber intrusion.**
- **James Lewis**
 - **The people thinking about the seriousness of cyber warfare tend to be computer people. We need to broaden the debate and get the involvement of more national security people, military people, etc.**
 - **You’d be shocked to discover how infrequently we have assessed nuts and bolts vulnerabilities, for instance the links from a guy sitting in front of his keyboard all the way to the floodgate on the dam.**

-
- **John Arquilla**
 - **We have to worry about the possibility of a campaign approach being taken by the cyber attackers in which they mount several attacks over a period of hours or days. Think about the economic impact of deploying a Nimda virus once a week for three months.**
 - **Cyber attacks will transform 21st century warfare, as militaries which are highly dependent on secure information systems will be absolutely crippled.**

- **Michael Skroch**
 - **To secure SCADA systems, we need end-to-end authentication and encryption to help prevent attacks.**
 - **SCADA systems currently don't have the firewalls, routers, anti-viral software, etc. that are needed to secure them from attack.**
 - **Industry has not developed a business case for cyber security. It may take a cyber Pearl Harbor before we implement the security that is required.**
 - **We understand physical security and know how to achieve it. SCADA systems are a component of U.S. critical infrastructure that we don't understand well today.**
 - **Most of the U.S. infrastructures that use SCADA systems underestimate the vulnerabilities associated with those systems, particularly because they're not interested in security, they're interested in delivering a product, and security is not viewed as a part of that process.**
 - **Industry is using common internet technologies, IP-based communications and operating systems that are popular and prevalent in our economy. In so doing, they are adopting the broad base of vulnerabilities and adversaries that are able to take advantage of those vulnerabilities.**
 - **SCADA is not protected to the same degree as IT infrastructure.**

- **Hacker:**
 - Penetrating a water or electrical SCADA system running a Microsoft operating system takes less than 2 minutes.
 - Typical hacker is 12-16 year-old
 - “Exploits” (published sample bit of code or program that demonstrates a flaw in software or hardware) are shared among hacker community. Digital Millennium Copyright Act seeks to prevent the sharing of such things.
 - Can use “zombies” to make the attack appear like it’s coming from a computer in another country.
 - U.S. is the most vulnerable nation-state because IT is so essential. This makes vital infrastructures highly vulnerable.
 - Scenarios:
 - “Pearl Harbor” – multiple cyber attacks by terrorist/rogue groups against water, power, etc.
 - Insider – e.g. altering calibration of precision parts production of military equipment.
 - Information acquisition costs – the knowledge necessary to launch a sophisticated attack is provided by our own side for free. Groups like Al Qaeda are willing to spend the time, energy and money to learn what they need to know to carry out effective attacks.
 - In 90 days, a team of 6 to 10 people could acquire very cheaply the equipment and knowledge needed to take out huge sections of U.S. infrastructure.
 - If you bring in the FBI to lecture you on computer security and to tell you what you need to worry about, SCADA is at the very top of the list.
 - If somebody was surveying our infrastructures for a potential missile attack, we would be very excited about it. The fact that water and electrical supply SCADA systems are being probed and mapped for holes and choke points doesn’t get attention because there’s no flash, bang or blood.
 - In my world, firewalls are referred to as “speed bumps”.
 - Hackers have broken into SCADA systems of critical infrastructure but had no idea what they were looking at. To them, it was just another insecure Windows box.
 - Government regulation won’t solve the problem. What’s needed is to stop providing protection to the Microsofts and the business models of security providers.
 - DoD networks, including SIPRNET, have numerous poorly secured nodes.
 - FBI and CIA are fully competent and are working diligently to secure their systems. NSA and DoD are not.
 - Neither the government nor the private sector firmly grasp the dangers. Vulnerabilities will stay in place until there are bodies in the streets.

-
- **Diane VanDe Hei**
 - **Most of water system security funds go to physical protections, not cyber**
 - **Rep. Lamar Smith (R-TX)**
 - **A mouse can be as dangerous as a bullet or a bomb**

Vulnerability, Risk, and Systems Theory are emerging as three fundamental literature streams

What is vulnerability as it applies to critical infrastructure systems?

How does risk and systems theory apply to vulnerability assessment of critical infrastructure?

How can system vulnerability be quantified?

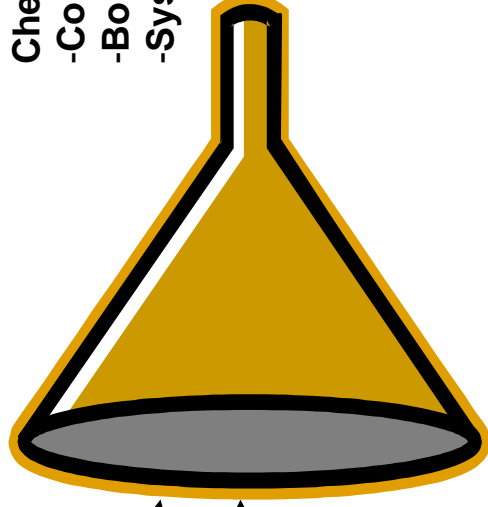
What results from deployment of the system vulnerability methodology?

Research Disciplines

Vulnerability

Risk

Systems



Gaps/Implications

Divergent views regarding vulnerability

Confusion on definitions of risk and vulnerability

Checklists without:

-Context

-Boundary

-System view:

-Emergence

-Equifinality

-Holography

-Variety

-Complexity

-Interdependencies

-Interactions

-Open, Natural, Rational, Cybernetic

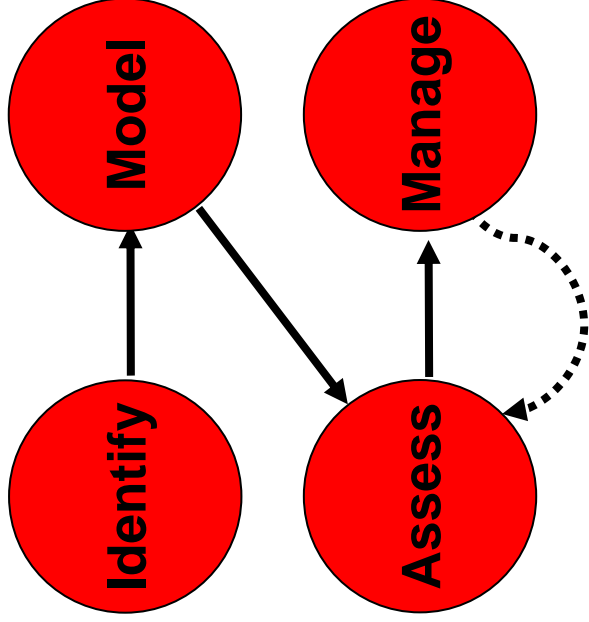
-Social-technical

-Negative Entropy

Relationship between risk and vulnerability is under study

- **Risk: a measure of the probability and severity of adverse effects**
- **Vulnerability: suggests *susceptibility* to risk**
- **Risk = $\{S_a, L_a, X_a\}_A$**
 - **The issue with extreme event probability L_a , is that it may cause a misleading rank ordering of threat scenarios**
 - **Perhaps threat scenarios and risk mitigation strategies should focus on system points of vulnerabilities rather than the expected value of damage**

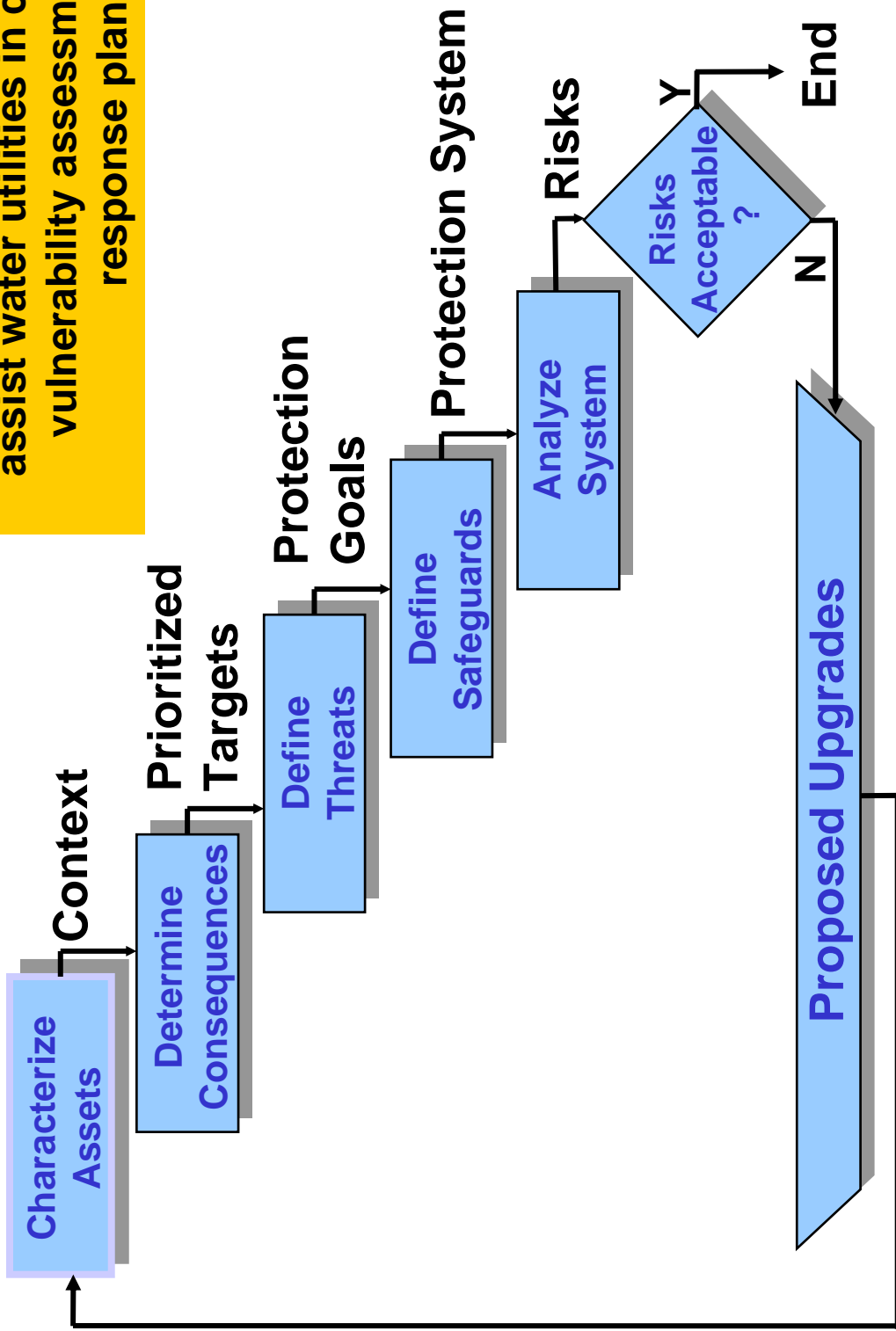
Improving IRAM would include systems perspective as a preceding step, making an systems description an explicit step in the assessment.



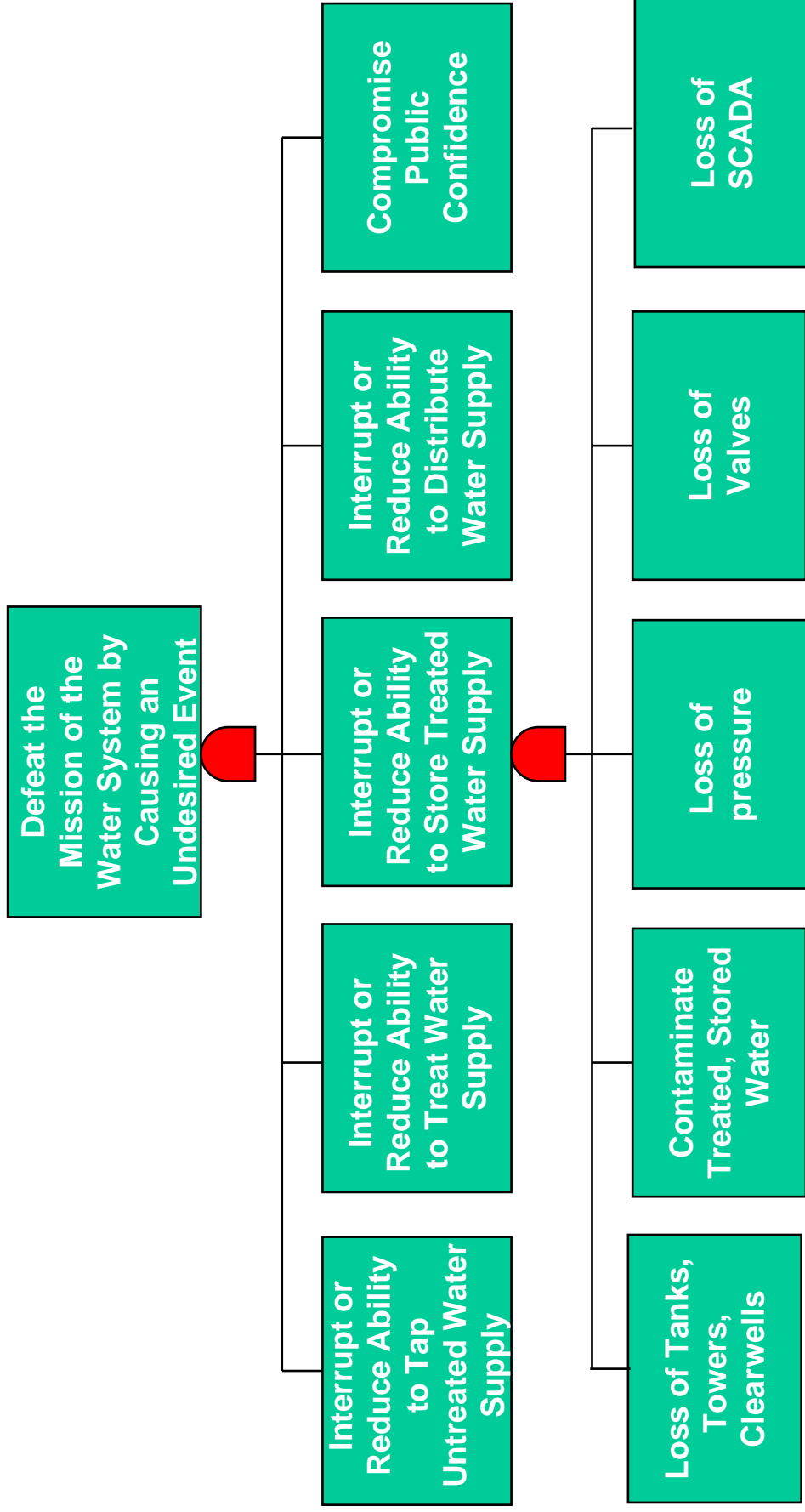
- Identify system and boundaries
- Decompose
 - Function
 - Component
 - State
- Ideate threat scenarios
- Assess mitigating aspects of the systems
- Assess Consequences
- Identify vulnerable points in the system
- Ideate risk mitigation strategies
- Model risk mitigation strategies
- Assess strategy performance
- Decide
- Implement

Risk Assessment Process

EPA has provided \$51M in grants to assist water utilities in conducting vulnerability assessments and response plans



Generic Modular Fault Tree



***Debunking the Threat to Water Utilities*¹ asserts the threat to water utilities is fear mongering.**

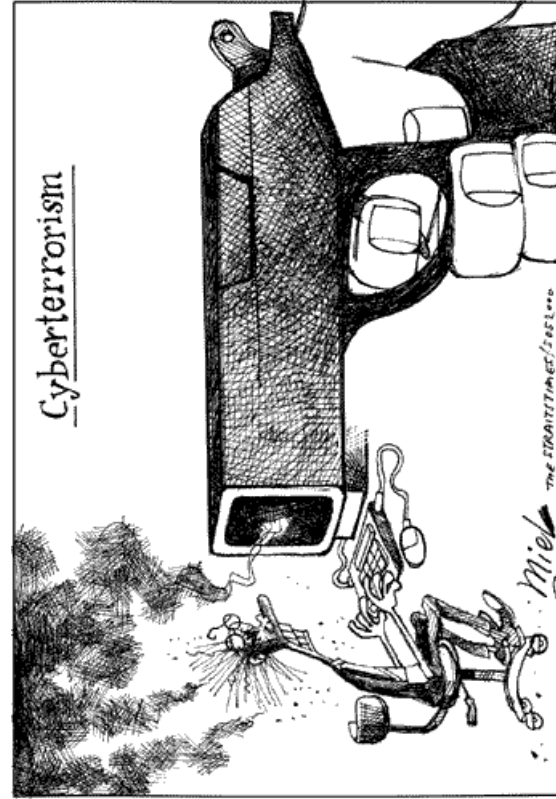
- Really?
 - “As you've probably heard, there was an interesting case of hacking at Maroochy, on the Queensland coast just north of Brisbane.” [...] “Over several months Maroochy Water Services experienced intermittent faults with its computerized sewerage SCADA system including numerous pump stations shut downs without any alarms, resulting in the first recorded conviction of a computer hacker causing serious environmental harm.”²
 - “As many of you know, the SQL Slammer worm struck last weekend (1/25/03) and caused overload conditions in the world wide Telco infrastructure.”³
- Wishful thinking or simply waiting for a disaster is dangerous and foolhardy.
- We know there are risks but how do we assess and mitigate?

1. CIO Magazine

2. http://www.courts.qld.gov.au/qjjudgment/ca02_151.htm

3. SCADA Mail List

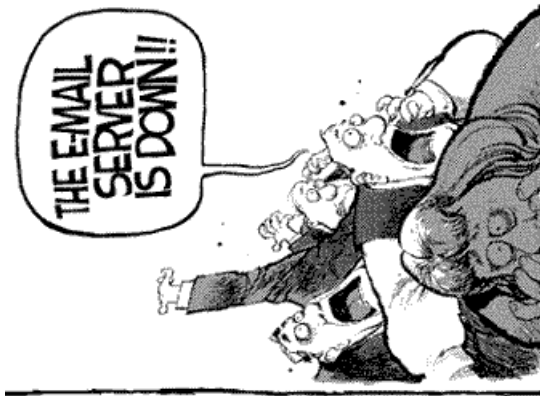
BRINGING CIVILIZATION TO ITS KNEES...



Cyberterrorism



TERROR IN THE 20TH CENTURY...



TERROR IN THE 21ST CENTURY...



© 2005 Cytizen
http://www.cytizen.com

Illustration by Bleeding

Abstract for

“Current and Future Challenges of Software Reliability Assessment”

By William H. Farr.

This presentation will present an overview of what software reliability is, why it is important to assess, how it is modeled, and some of the do's and don'ts in the modeling approach. This will provide a current state of practice for the methodology. The talk will then cover current challenges for the methodology and future challenges as the software engineering environment evolves.



Current and Future Challenges of Software Reliability Assessment

October 30, 2003

**US Army Conference on Applied Statistics
Napa Valley, CA**

William Farr, PhD

B35, Combat Systems Branch Head

Phone: (540) 653-8388

Fax: (540) 653-8673

Email: farrwh@nswc.navy.mil



Outline Of Presentation

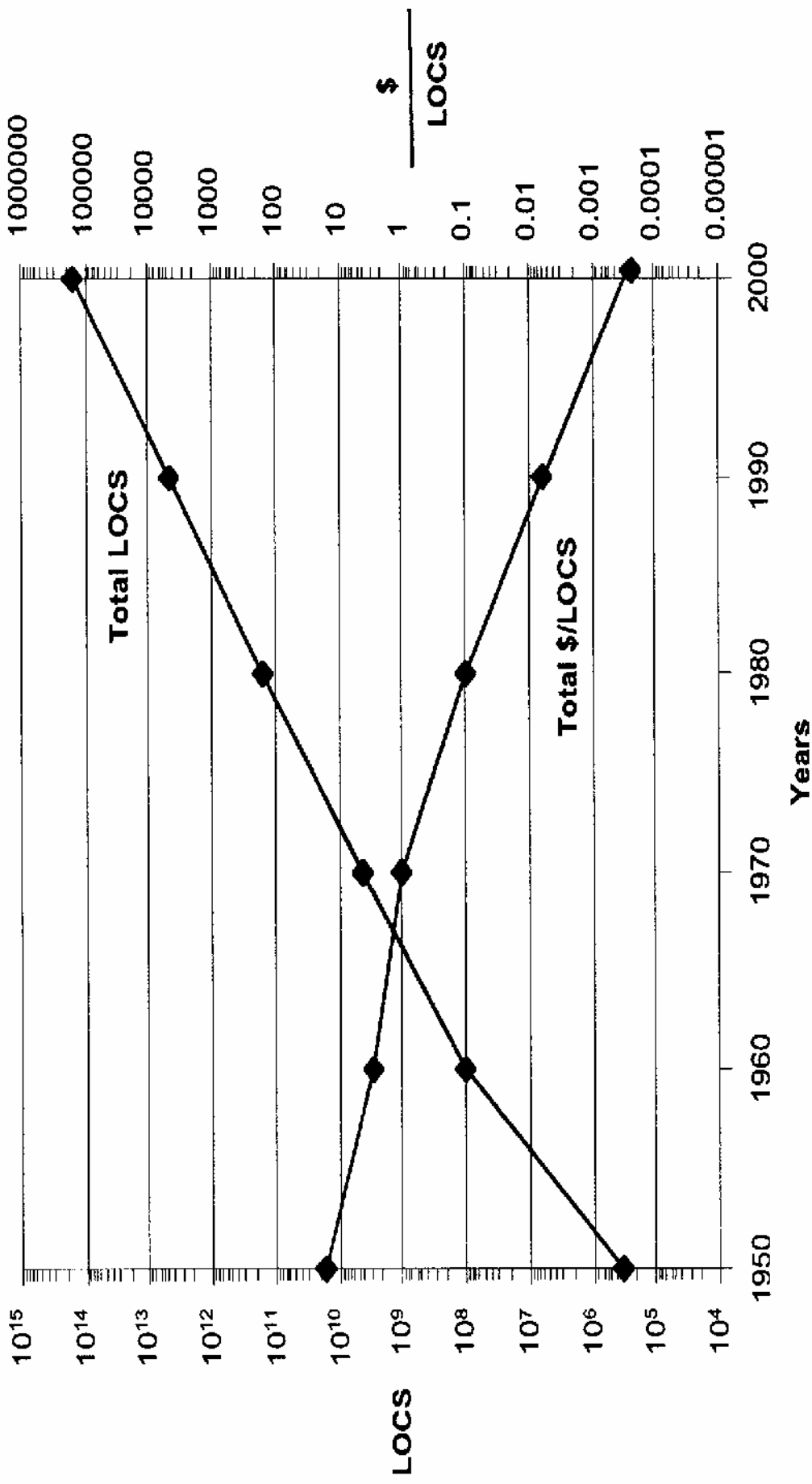
- **Introduction**
- **What is Software Reliability and how is it assessed**
- **Current Challenges**
- **Future Challenges**
- **Summary**

Introduction

- Software is playing an ever greater role in our systems.
 - More intrinsically difficult tasks are now undertaken by software.
 - Software has moved from an auxiliary role to a primary role in providing critical services.
 - Software based applications are becoming an accepted part of our life.
- Software Reliability assessment is therefore of major concern to developers and users.

Lines of Code in Service:

U.S. DoD





What is Software Reliability and How is it Assessed

DEFINITIONS

SOFTWARE RELIABILITY IS THE PROBABILITY THAT A GIVEN SOFTWARE PROGRAM WILL OPERATE WITHOUT FAILURE FOR A SPECIFIED TIME IN A SPECIFIED ENVIRONMENT.

SOFTWARE RELIABILITY ENGINEERING (SRE) IS THE APPLICATION OF STATISTICAL TECHNIQUES TO DATA COLLECTED DURING SYSTEM DEVELOPMENT AND OPERATION TO **SPECIFY, PREDICT, ESTIMATE, AND ASSESS** THE SOFTWARE RELIABILITY OF SOFTWARE-BASED SYSTEMS.

Errors, Faults, Failures

- **ERROR** – Human action that results in software containing a fault.
- **FAULT** – A defect in code that can be the cause of one or more failures (synonymous with “bug”).
- **FAILURE** – The inability of a system or system component to perform a required function within specified limits. A departure of program operation from program requirements.

ERROR \Rightarrow FAULT(S) \Rightarrow FAILURE(S)

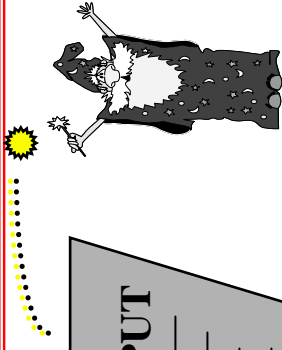
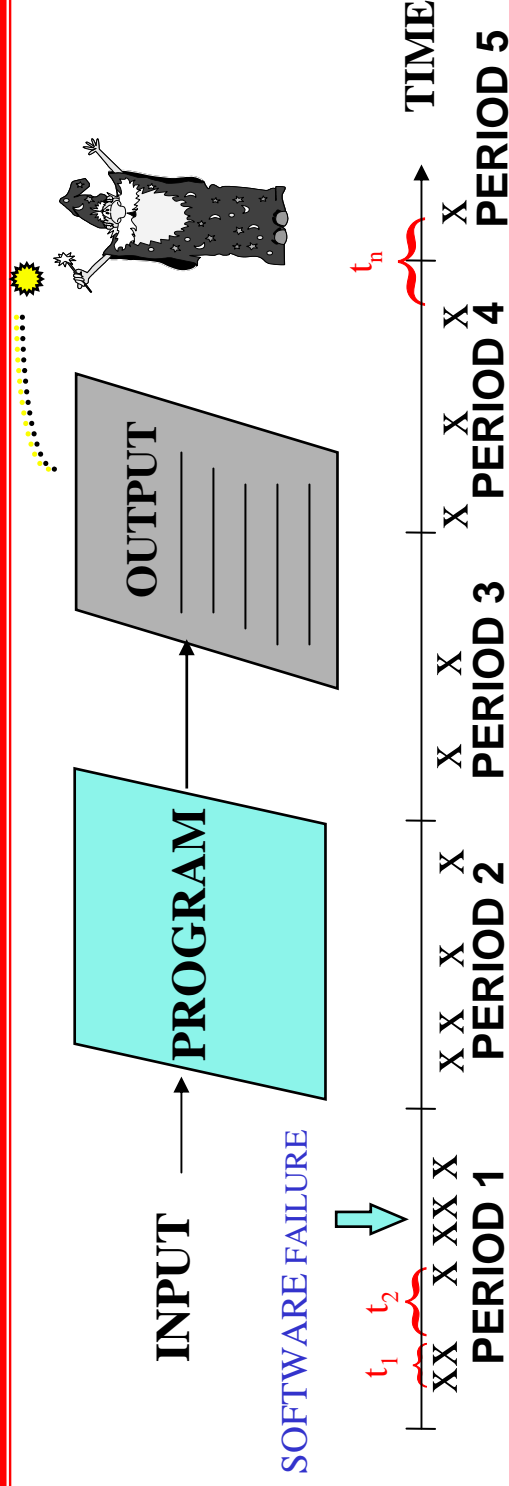
Comparisons Between Software & Hardware

- Both are stochastic processes and can be described by probability distributions.
- Software faults are design faults not physical faults.
- Software does not “wear-out, burn-out, or deteriorate” over time.
- If there is a fault in a particular version of the software it will be in every copy.
- Software reliability is much more difficult measure to obtain and analyze.
- Software is continuously modified throughout its life cycle.

APPROACHES OF S/W RELIABILITY ESTIMATION

- ERROR SEEDING/TAGGING
MODELS
- DATA DOMAIN
- TIME DOMAIN

APPROACH TO ESTIMATING S/W RELIABILITY IN THE TIME DOMAIN



SOFTWARE RELIABILITY DATA:
NUMBER OF FAILURES/PERIOD (e.g.,
6,4,2,3,1,...)

OR

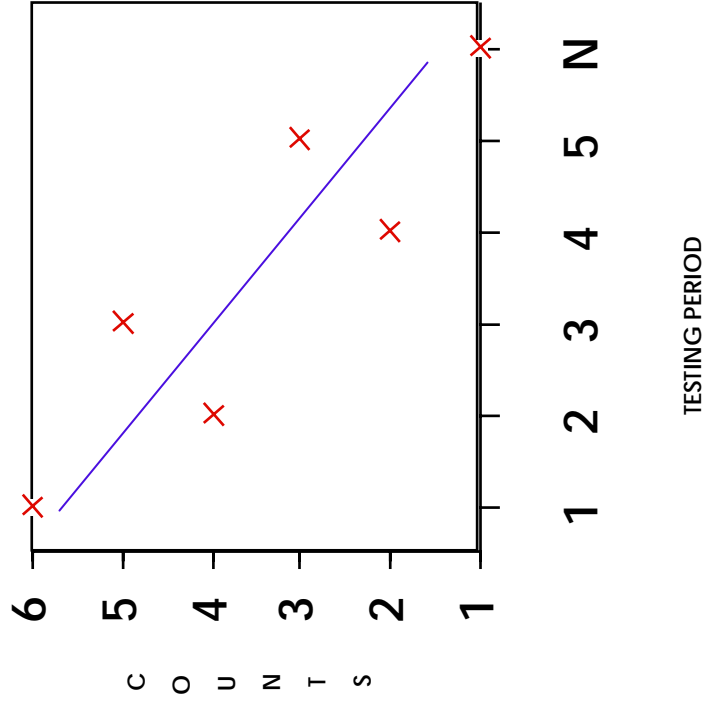
TIME BETWEEN FAILURES (e.g., t_1, t_2, \dots, t_n)

TIME UNITS

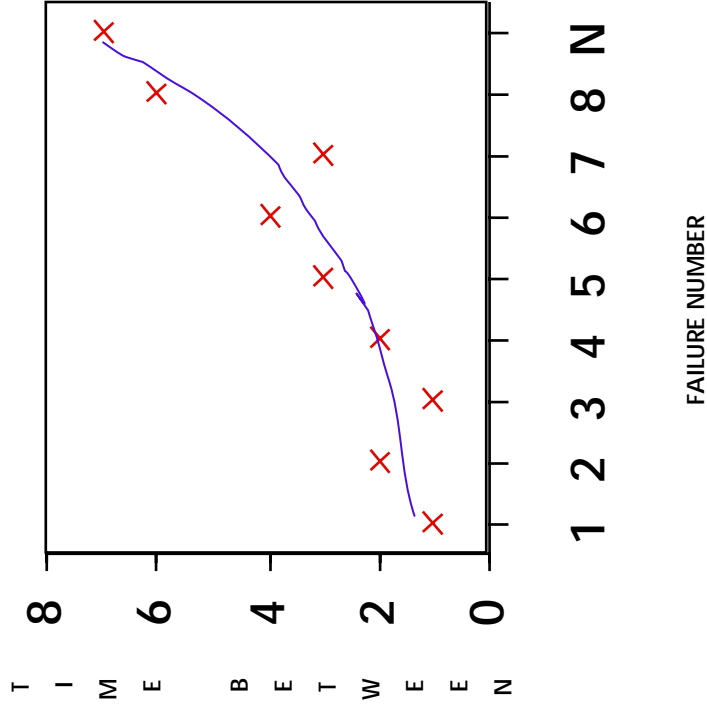
- **“Counts”** – by day, week, month, quarter, etc.
- **“Time Between”** – Wall Clock, CPU, Number of Test Cases, Natural unit (e.g. # transactions, # printed pages, etc.)

Personal Experience: Counts by month provided the easiest to implement without impacting prediction accuracy.

SOFTWARE RELIABILITY MODELS IN THE TIME DOMAIN



FAILURE COUNTS PLOT



FAILURE TIME PLOT

EXAMPLE SOFTWARE RELIABILITY MEASURES

ERROR COUNT MODELS

- TOTAL NUMBER OF ERRORS
- EXPECTED NUMBER OF ERRORS IN FUTURE TESTING PERIODS
- FAILURE RATE
- TESTING TIME REQUIRED TO ELIMINATE THE NEXT K ERRORS OR TO ACHIEVE A SPECIFIED FAILURE RATE
- CURRENT PROGRAM RELIABILITY

TIME BETWEEN MODELS

- TOTAL NUMBER OF ERRORS
- MEAN TIME TO NEXT FAILURE
- FAILURE RATE
- TESTING TIME REQUIRED TO ACHIEVE A SPECIFIED RATE
- CURRENT PROGRAM RELIABILITY

EXAMPLES OF SOFTWARE RELIABILITY MODELS

Time Models

- Littlewood & Verrall Bayesian Model
- Musa Basic Execution Model
- Geometric Model*
- Non-homogeneous Poisson Process Model
- Musa Logarithmic Poisson Model*
- Jelinski-Moranda Model

Error Count Models

- Generalized Poisson Process Model
- Non-homogeneous Poisson Process Model for Counts
- Brooks and Motley's Model
- Schneidewind Model*
- S-shaped Reliability Growth Model*

Note: The above models are all in SMERFS³ and CASRE.

The models with the "*" have shown themselves to be especially applicable to many data sets.

CURRENT STATUS OF RELIABILITY MODELING

- Currently over 100 models in the Literature
- SRE Newsletter, International Symposium (14 held – next in Denver), Numerous Publications on the subject including 4 key books, (John Musa – “Software Reliability Engineering”)
- Standards already exist, e.g. IEEE 982.1 and 982.2 and AIAA’s ANSI/AIAA “Recommended Practice for Software Reliability” (R-013-1992)
- Software Reliability modeling is applied in a wide variety of applications, e.g. Communications (AT&T Best Practice), Space Program (Shuttle, Galileo); and the DOD

Current Challenges

- **Diversity of Reliability Concerns**
- **Ultra reliability requirements ($\leq 10^{-8}$) of some software systems**
- **Demonstrated value of software reliability assessment**
- **Basic Assumptions**
 - Independence, fault correction, & the environment
- **Changing nature of software development**
 - Component-based multi-tiered software; COTS; quality standards; CMMI; web development; growth of programming languages; programmers without formal training
- **Emphasis on shorter development time; a push to market**

Future Challenges

Min Xie – “Software reliability is ready for a new phase of development where emphasis is on practical implementation.”

- **Specify and write high quality, modular, high performance software for a wide variety of applications**
- **Combining various views of reliability (testing, code inspections, audits, reliability modeling, etc.) into an overall assessment**
- **Develop techniques for early reliability assessment**
- **Develop reliability assessment for evolving technologies (Example: autonomous agents)**
- **Develop reliable software for web based systems**
 - Support 24/7 applications
- **Demonstrate ultra-high reliability for safety critical systems**
- **Develop reliability models for complex systems**
 - System of systems
 - Hardware & software
 - Factor the human element into the equation

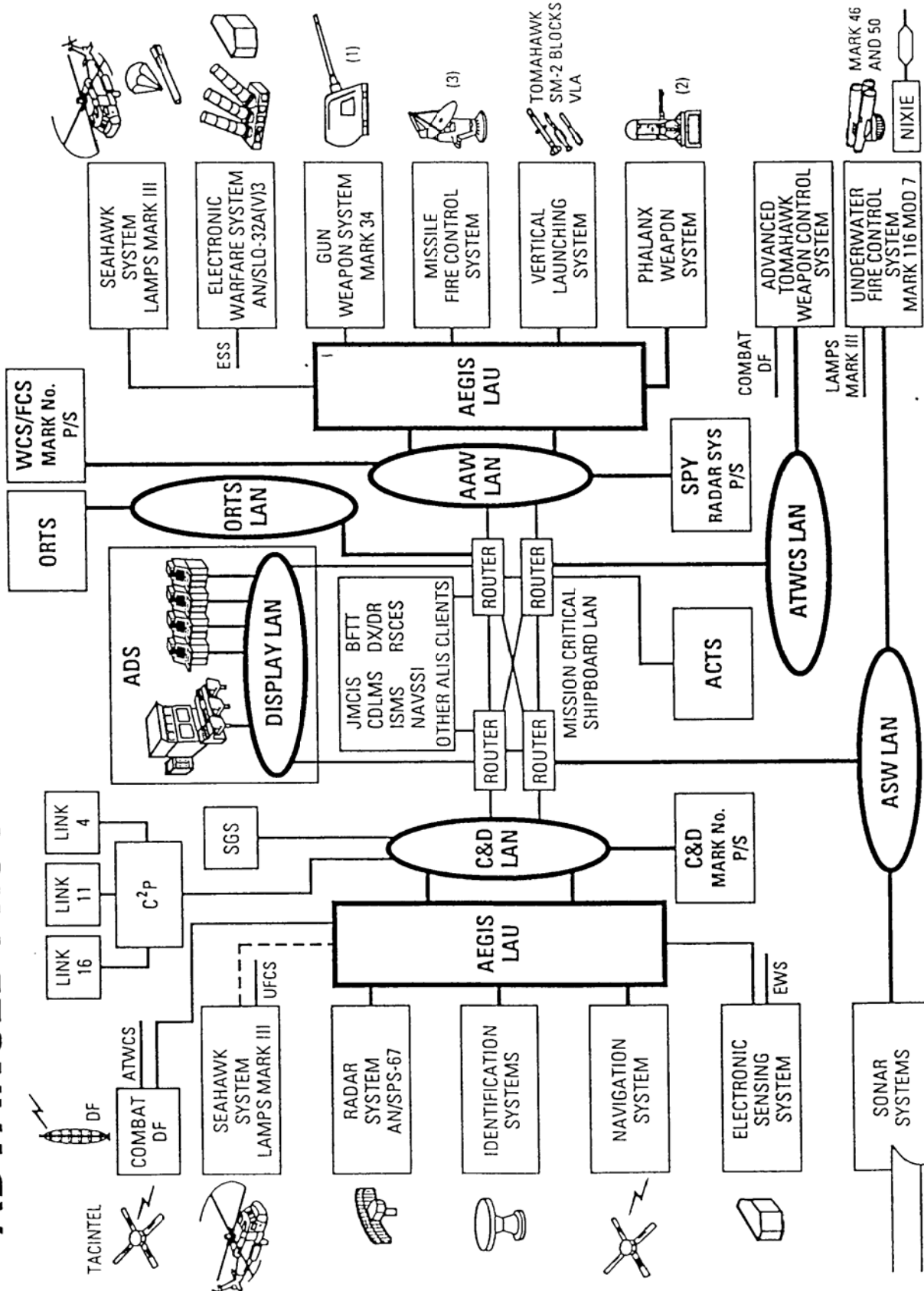
Examples of Complex Systems

- Air Traffic Control System
- Baggage Handling System at the Denver Airport
- Space Shuttle
- Medical Test Equipment
- Banking & E-commerce
- Navy warship

AEGIS DDG



WARSHIP ("SYSTEM of SYSTEMS")



NOTE: P/S = PRIMARY AND SECONDARY NODES / LAU = LAN ACQUISITION UNIT

Human Factors Considerations on Reliability



Do's and Don'ts of Software Reliability

Do's

1. Consider it as a quantitative measure to be used in determination of software quality.
2. Track reliability over time and/or software version to determine whether reengineering is needed.
3. Develop specific objectives for a software reliability effort.
4. Develop an automated data collection process and database.

Don't

1. Consider one software reliability model is applicable for all situations.
2. Extrapolate results beyond the environment in which the data is generated.
3. Consider it as the sole data point upon which to make judgments.

SUMMARY

- **RELIABILITY ANALYSIS CAN PROVIDE USEFUL INFORMATION TO ASSESS BOTH THE PRODUCT AND THE ORGANIZATION'S DEVELOPMENT PROCESS**
- **Challenge: WHEN MAKING AN OVERALL PRODUCT ASSESSMENT, CONSIDER RELIABILITY MODELING AS PROVIDING ONE PERSPECTIVE - OTHERS NEED TO BE CONSIDERED AS WELL**
- **Challenge: THERE IS A NEED TO QUANTIFY BOTH COMPONENT (HARDWARE & SOFTWARE) AND SYSTEM'S RELIABILITY**
- **Challenge: FOR TOTAL SYSTEM'S RELIABILITY THE HUMAN FACTOR MUST BE WORKED INTO THE FRAMEWORK ALSO**
- **Challenge: RELIABILITY NEEDS TO BE CONSIDERED OVER THE ENTIRE LIFECYCLE OF THE SYSTEM**
- **Challenge: NEED TO SHARE INFORMATION ON APPLYING THE METHODOLOGY TO DEMONSTRATE ROI**

KEY REFERENCES

- *Software Reliability Engineering*, by John Musa, McGraw-Hill, 1999.
- *Software Reliability Measurement, Prediction, Application*, by J. Musa, A. Iannino, and K. Okumoto, McGraw-Hill, 1996.
- *Software Reliability Modeling*, by M. Xie, World Scientific, 1991.
- *Handbook of Software Reliability Engineering*, Edited by Michael Lyu, McGraw-Hill, 1997.
- *Recommended Practice for Software Reliability*, ANSI/AIAA, R-013-1992, ISBN 1-56347-024-1
- Center For Software Reliability - <http://www.csr.city.ac.uk/>
- Links to Others - <http://members.aol.com/JohnDMusa/>
<http://www.dacs.dtic.mil/>

Useful Software Reliability Modeling Practices in Industry Environments

*Daniel R. Jeske
University of California
Riverside, CA*

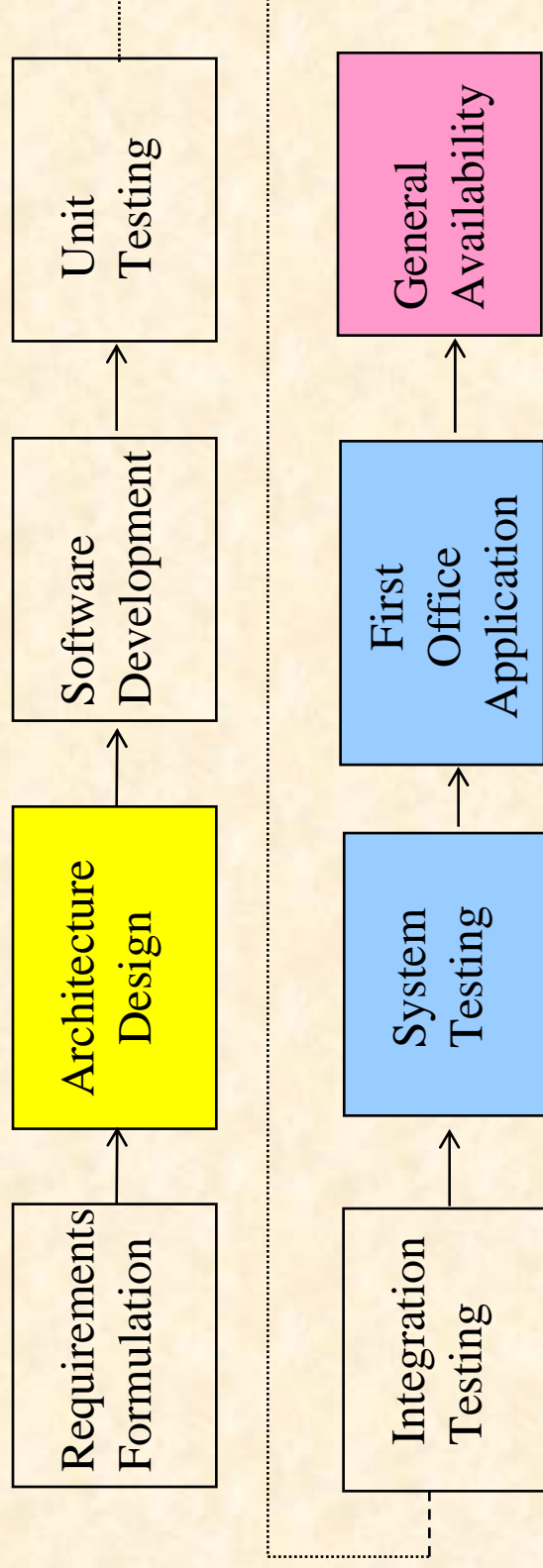
- 1. Introduction and Context*
- 2. Architecture-Based Software Reliability Models*
- 3. Software Reliability Growth Model*
- 4. Case Studies*
- 5. Summary*

U.S. Army Conference on Applied Statistics
October 29-31, 2003

When are Software Reliability Models Typically Applied?

Architecture-Based Reliability Models

.....decisions are being made as to how to design reliability into the system



Software Reliability Growth Models

.....reliability is quantified and influences the release decision

Software Reliability Growth Models

.....reliability predictions are verified

.....parameters useful for modeling reliability of next release are estimated

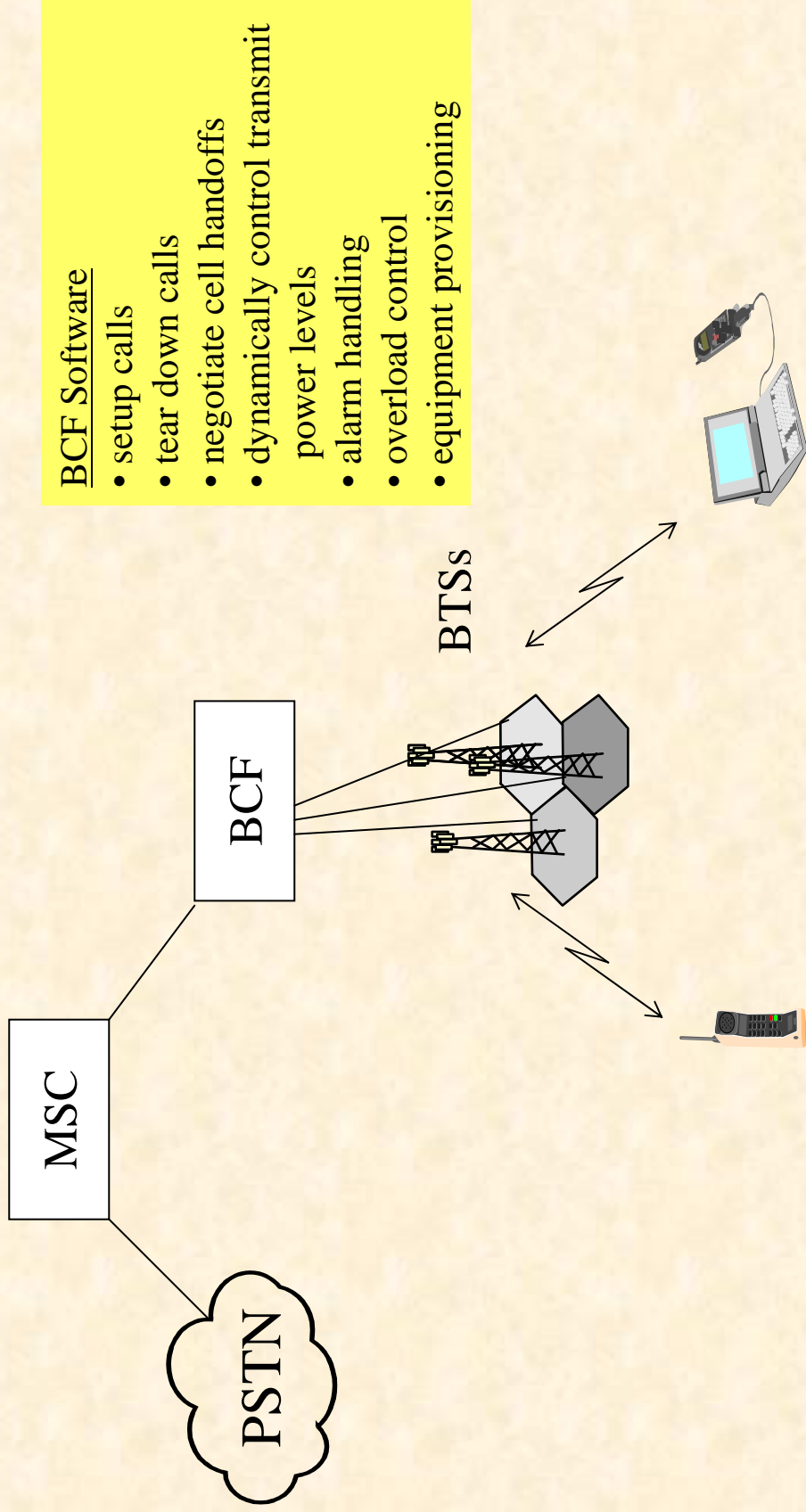


Architecture-Based Software Reliability Models

Questions to Answer

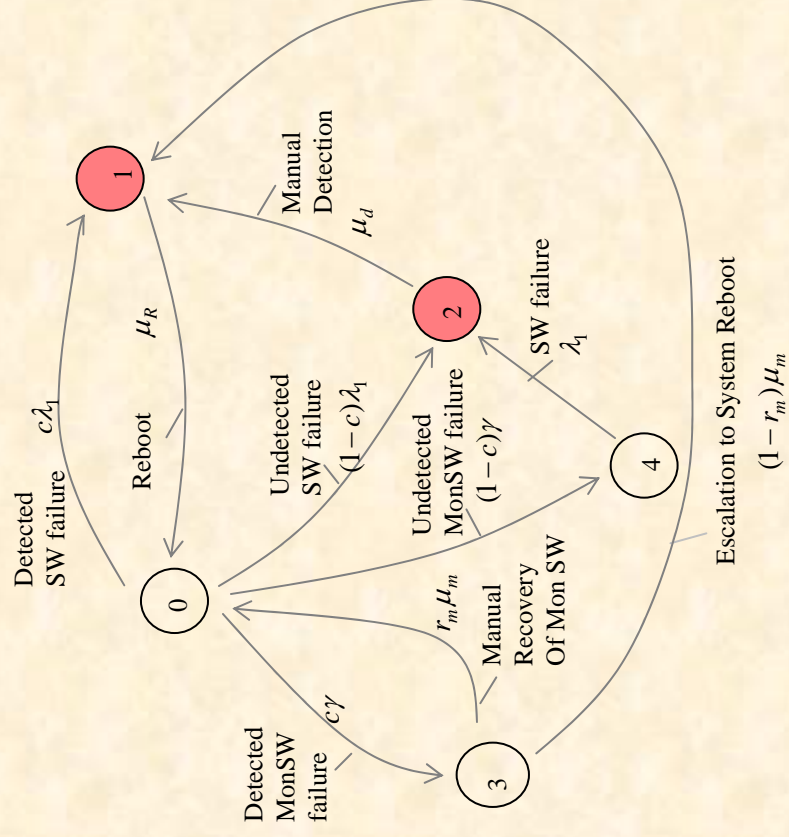
- What type of software redundancy, if any, is needed in order to achieve the reliability target?
- How fast do fault recovery times need to be?
- Should software processes try to restart before a failover is attempted?
- How thorough do system fault diagnostics need to be?
- What is the required software failure rate?

Case Study: GSM Wireless System

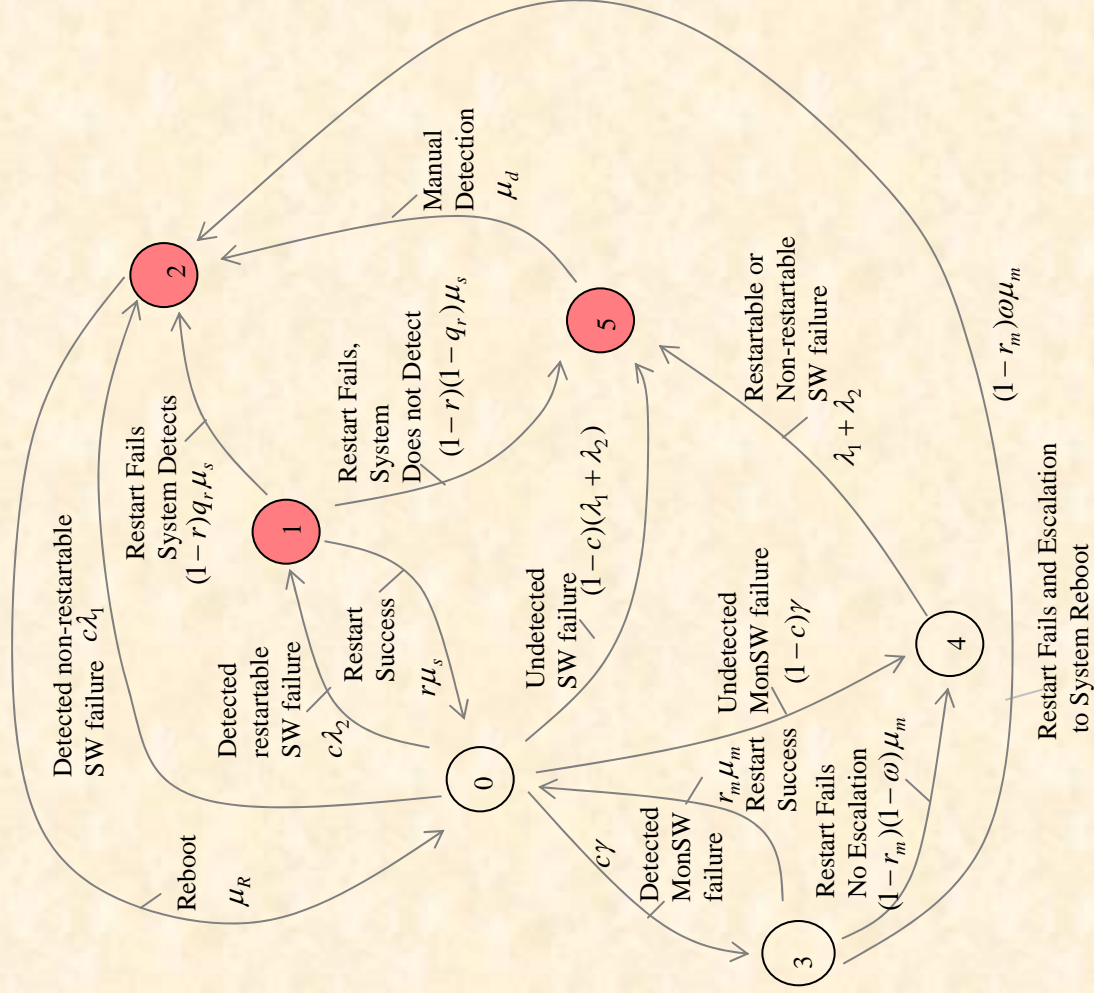


- BCF - Base Station Controller Frame
- BTS - Base Transceiver Station
- MSC - Mobile Switching Center
- PSTN - Public Switched Telephone Network

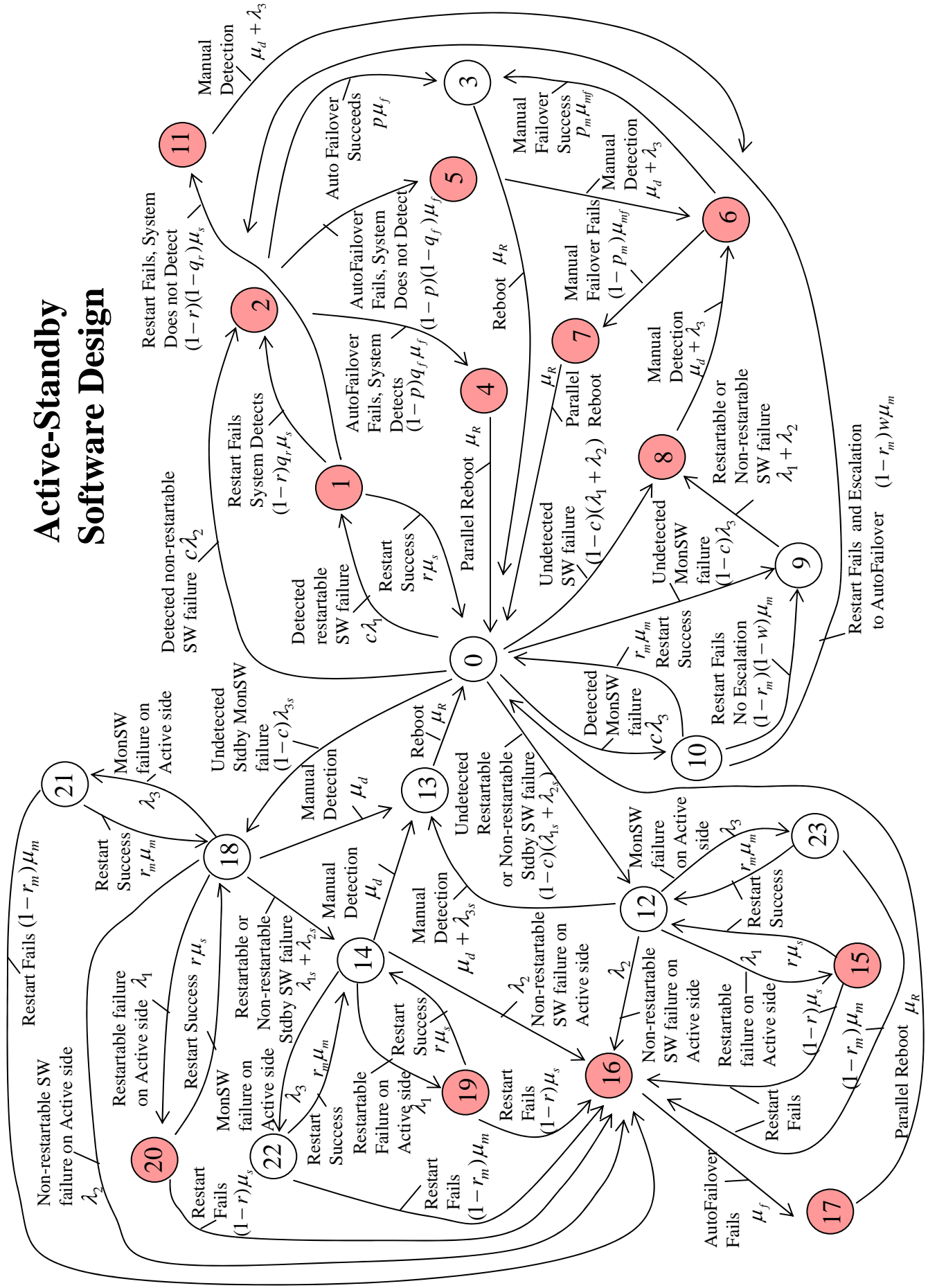
Simplex Design Without Restart Feature



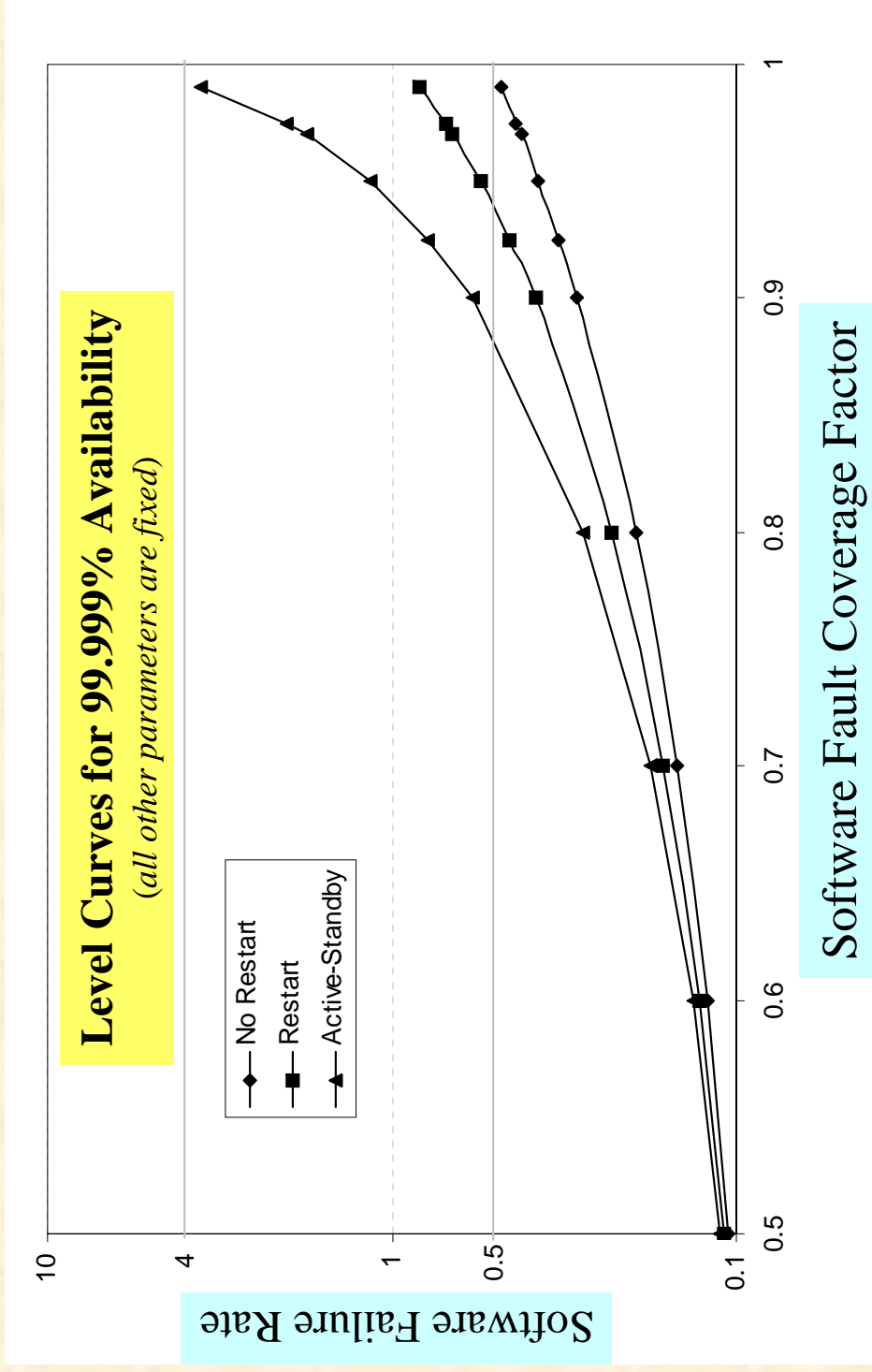
Simplex Design With Restart Feature



Active-Standby Software Design



Feasibility Regions for Critical Parameters



Identifying Influential Parameters

Challenges

- Availability formula is straight-forward to obtain, but is complicated in that it is highly non-linear and it depends on numerous (23) parameters
- Analysis of derivatives, while also straight-forward is not particularly insightful, since they themselves depend on the same set of parameters.
- “Tornado Graphs” depends too much on the fixed values for the other parameters.

Approach

- Identify likely ranges for each parameter to define the domain of the availability function
- Uniformly sample N times from the domain to obtain A_1, A_2, \dots, A_N
- Use an efficient second-order polynomial to approximate the availability formula
- Use the t-statistics as the associated measure of influence

Case Study – Continued

(Likely Ranges for Parameter Values in Active-Standby Architecture)

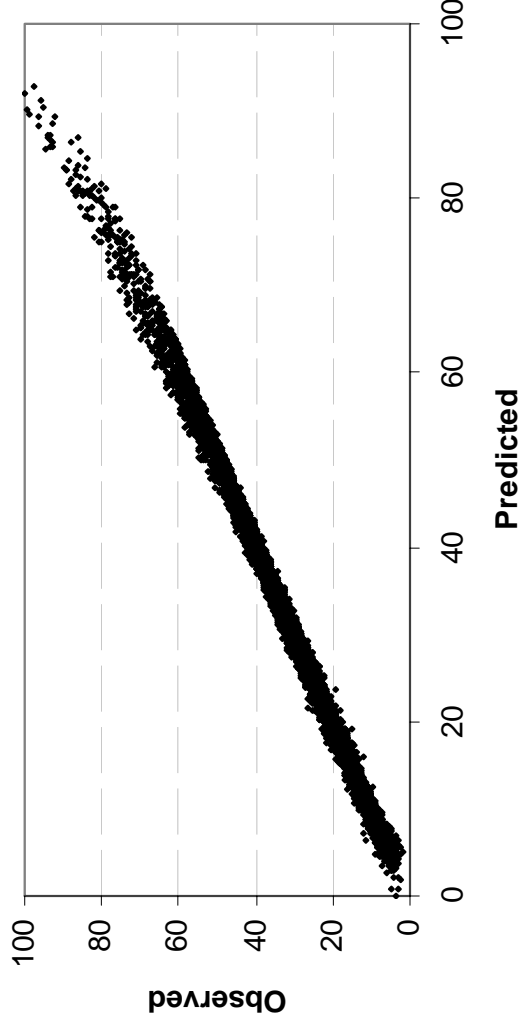
λ_1	Fails/yr	[0.1, 1.0]	μ_{mf}	Failovers/yr	26,280
λ_2	Fails/yr	[0.1, 1.0]	μ_m	Recoveries/yr	17,520
λ_3	Fails/yr	[0.1, 0.5]		r	[0.9, 0.99]
λ_{1s}	Fails/yr	0		r_m	[0.9, 0.99]
λ_{2s}	Fails/yr	[0.1, 0.5]		q_r	[0.9, 0.99]
λ_{3s}	Fails/yr	[0.1, 0.5]		q_f	[0.9, 0.99]
μ_s	Restarts/yr	31,536,000		p	[0.9, 0.99]
μ_f	Failovers/yr	3,153,600		p_m	[0.9, 0.99]
μ_r	Reboots/yr	26,280		c	[0.9, 0.99]
μ_d	Detections/yr	[1095, 2190]		w	[0.9, 0.99]

Case Study – Continued

(Efficient Second-Order Regression Model)

$$D(\text{min/yr}) = 497 + 379\lambda_1 + 378\lambda_2 - 358\lambda_1c - 356\lambda_2c + 284(\mu_d / 1000)c - 268(\mu_d / 1000)c - 461c \\ - 12.2\lambda_1(\mu_d / 1000) + 319\lambda_3 - 12\lambda_2(\mu_d / 1000) - 298\lambda_3c \\ - 22p - 11.6\lambda_3(\mu_d / 1000) - 11.7r - 9.6q_f - 9q_r - 5.8r_m - 5w - 2.1p_m$$

Adequacy of Fit



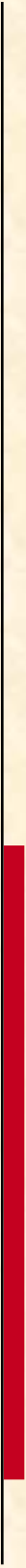
Interpretation

Influential Parameters

- Software failure rates on active side
- Silent failure detection time
- Coverage factor
- Cross-product of coverage factor with software failure rates and silent failure detection time

Weakly Influential Parameters

- All of the success probabilities



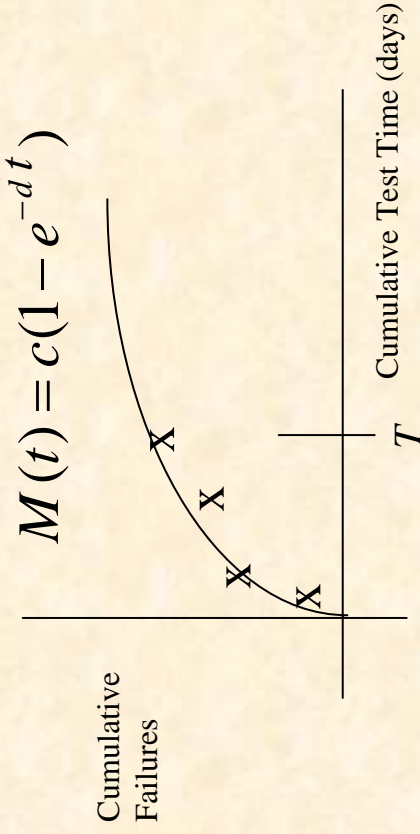
Software Reliability Growth Models (SRGMs)

SRGMs: Questions to Answer

- What would the user-perceived failure rate of the software be if the software was to be released now?
- How much more testing is needed to achieve the reliability targets?
- How many faults exist in the code at the end of system test?

A Well-Known SRGM

Goel-Okumoto Model



c = initial number of faults

d = average per-fault failure rate

t = cumulative test time, aggregated across all installations

T = current value of t

Failure Rate Estimate

$$\lambda(T) = M'(T) \\ = cde^{-dT}$$

Challenges With Applications

Traditional SRGMs are applied to system test data with the hope of obtaining an estimate of software failure rate that will be observed in the field. The following issues have to be taken into consideration:

- The usage profile in the field is typically very different than the testing profile
- Fault removal in field environments is not instantaneous
- Quality of software reliability data

Non-Instantaneous Fault Removal Times In Field Environments

G-O Model

$$M(t) = c[1 - e^{-dt}]$$

$$\lambda(t) = cde^{-dt}$$

where:

c initial number of faults at time of field deployment

d average per fault failure rate in a field environment

p probability that a detected fault is successfully removed

G-O Model with Imperfect Debugging

$$M(t) = \frac{c}{p}[1 - e^{-dpt}]$$

$$\lambda(t) = cde^{-dpt}$$

Non-Instantaneous Fault Removal Times In Field Environments


- Perfect debugging assumption is an acceptable assumption (based on thorough regression tests)
- We relate non-instantaneous fault removal to imperfect debugging

Under the imperfect debugging model:

$$E(\# \text{ of occurrences of each fault}) = 1/p$$

Under the situation of non-instantaneous fault removals, if μ denotes the mean time to remove a fault and there are n systems in the field

$$E(\# \text{ of occurrences of each fault}) = 1 + n\mu d$$


$$\frac{1}{p} = 1 + n\mu d$$

Non-Instantaneous Fault Removal Times In Field Environments

$$M(t) = c(1 + n\mu d) \left[1 - e^{-\frac{d}{1+n\mu d} \times t} \right]$$

$$\lambda(t) = cde^{-\frac{d}{1+n\mu d} \times t}$$

- where:
- c initial number of faults at time of field deployment
 - d average per fault failure rate in the field environment
 - μ average time to remove a fault ($\mu = 0$ gives back the G-O model)
 - n number of systems in the field
 - t cumulative exposure time, aggregated across all installations

Field Failure Rate Prediction **if** Testing Profile Matches the Field Usage Profile

The number of initial faults at time of field deployment time is the same as the number of residual faults after testing has completed

The average per fault failure rate in the field environment is identical to the average per fault failure rate in the test environment

For the field failure rate model, we can use the G-O model replacing c with ce^{-dt} and adjusting for non-instantaneous removal times

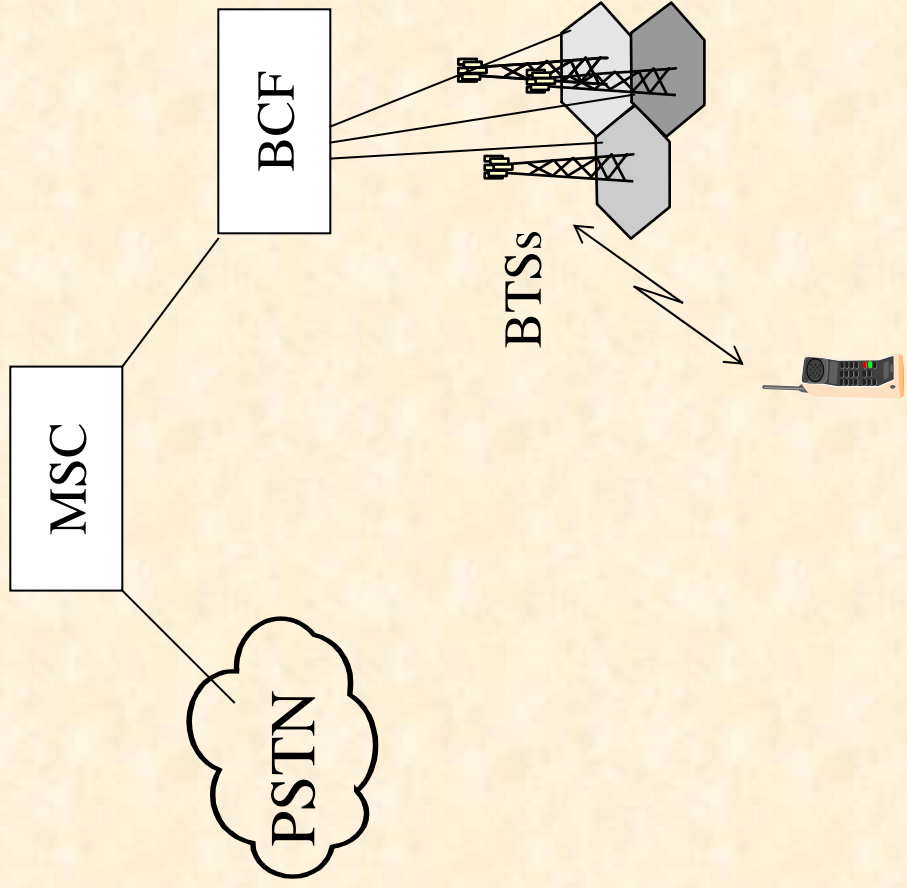
$$\lambda(t) = (ce^{-dt})de^{-\frac{d}{1+n\mu d} \times t}$$

Field Failure Rate Prediction When Testing Profile Does Not Match the Field Usage Profile

- 1) Testing and field usage profiles drive the average per fault failure rate. Usually the failure rate of faults is smaller in field environments than in test environments
- 2) Define $K = d / d^*$, where d^* is the average per fault failure rate in the field environment
- 3) K is the “per fault failure rate” calibration factor
- 4) Estimate K from previous releases of the software, or from related projects

$$\lambda_{adj}(t) = (ce^{-dT}) (d / K) e^{-\frac{d / K}{1+n\mu d / K} \times t}$$

Case Study - Continued



- BCF - Base Station Controller Frame
- BTSS - Base Transceiver Station
- MSC - Mobile Switching Center
- PSTN - Public Switched Telephone Network

BTS Controller Software

- allocate radio resources
- negotiate hand-offs
- manage connections between BTS and BCF
- provision and maintenance

Alternative Architectures

- Simplex, no process restart capability
- Simplex, with restart capability
- Active-Standby, with fail-over capability

Objective

Goals

Estimate the field failure rate of R3 BCF software which is currently in system test

Data Available

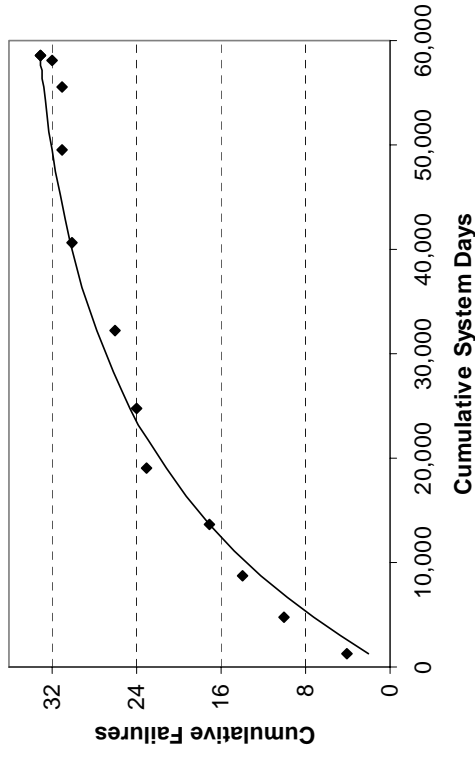
- 1) R3 Test Data in a non-operational profile environment
- 2) Field failure data for R1 and R2

Field Failure Data for R1

Month	System Days		Failures	
	Days	Cumulative	Month	Cumulative
1	1,249	1,249	4	4
2	3,472	4,721	6	10
3	4,065	8,786	4	14
4	4,883	13,669	3	17
5	5,425	19,094	6	23
6	5,656	24,750	1	24
7	7,549	32,299	2	26
8	8,295	40,594	4	30
9	8,882	49,476	1	31
10	6,120	55,596	0	31
11	2,465	58,061	1	32
12	527	58,588	1	33
13	45	58,633	0	33

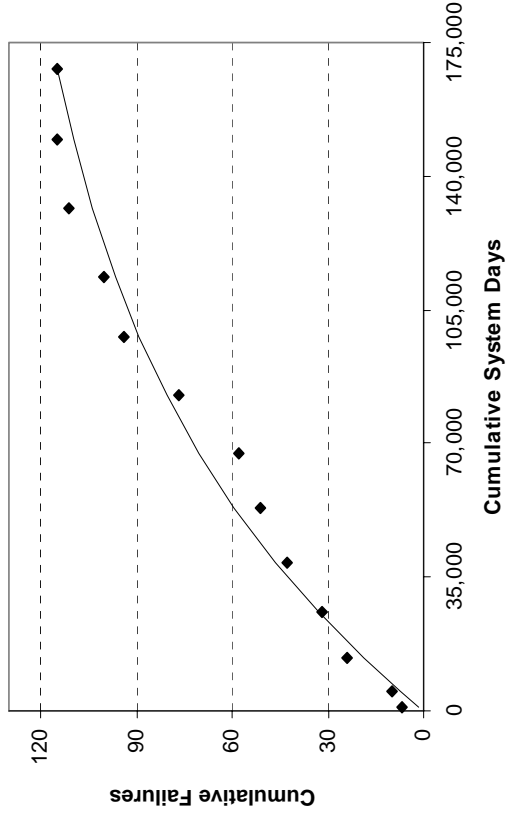
Field Failure Data Analysis

R1 Field



$T_1 = 58,633$ system-days (13 months, 167 systems)
 $\mu = 45$ days
 $\hat{c}_1 = 22.3$ faults
 $\hat{d}_1 = 0.0000746$ failures/day/fault
 (StdError: 3.07×10^{-5})

R2 Field



$T_2 = 167,900$ system-days (13 months, 370 systems)
 $\mu = 45$ days
 $\hat{c}_2 = 114.01$ faults
 $\hat{d}_2 = 0.0000128$ failures/day/fault
 (StdError: 3.07×10^{-6})

R3 Test Data

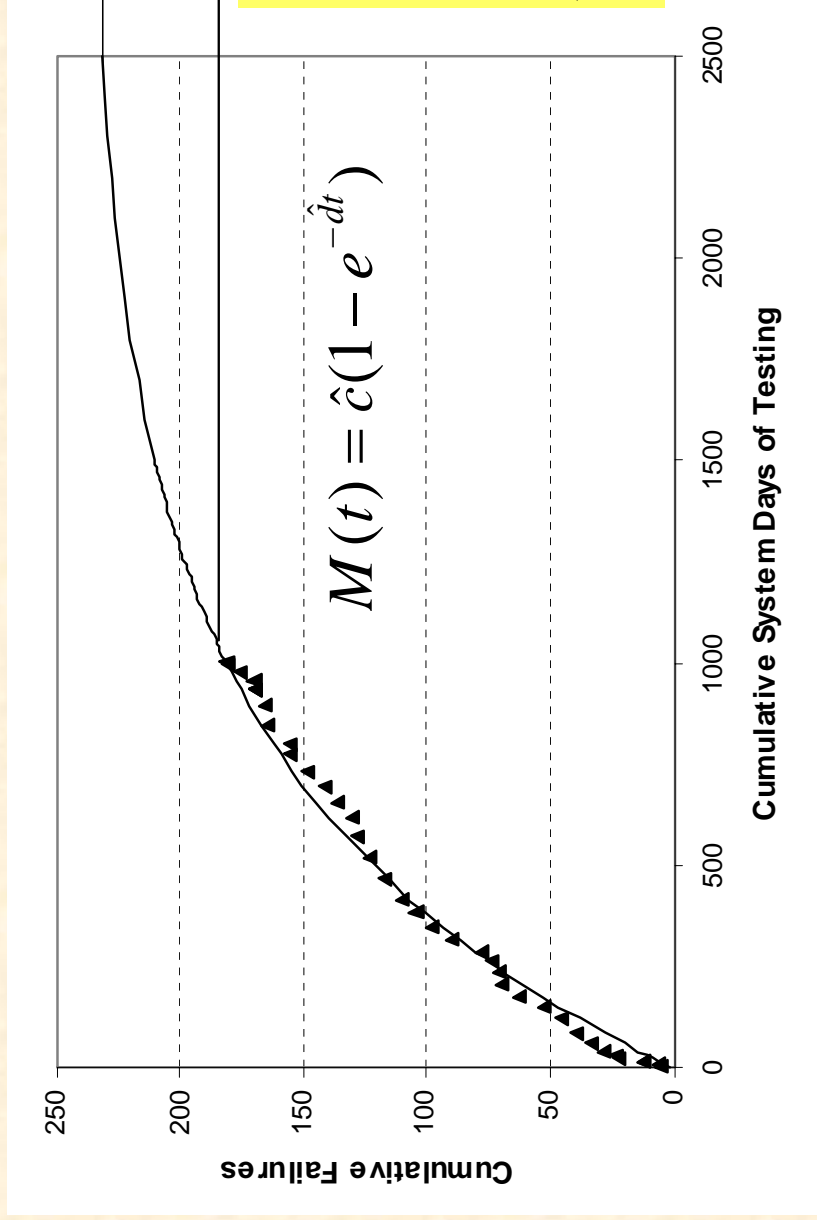
Calendar Week	Number of Days of Each Software Installation											Cumulative System-Days of Testing	Cumulative Failures
	1	2	3	4	5	6	7	8	9	10	11		
1				5								5	5
2				4								9	6
3				4								13	13
4				5								18	13
5		5		5								28	22
6				5								33	24
7		5		5								43	29
8	5	5		5				5				63	34
9	5	5		5			5	5	5			88	40
10	5	5		5			5	5	5	5		123	46

As many as 11 frames were being used in parallel during system test.
 'Failure' is defined as a severity 1, 2 or 3 MR.



33				5			5	5				6	955	170
34				5			5	6				6	977	176
35				5			5	6				6	999	180
36								2					1001	181

G-O Model Fit to the R3 Test Data

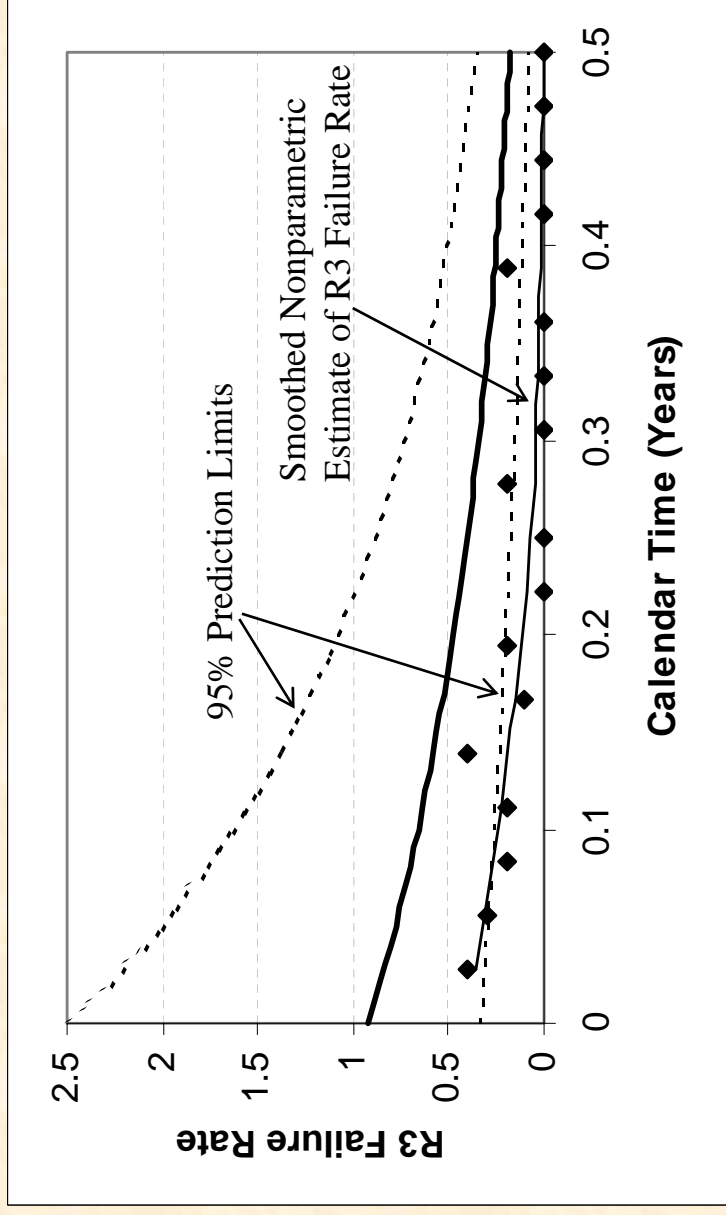


$$\begin{pmatrix} \hat{c} \\ \hat{d} \end{pmatrix} = \begin{pmatrix} 238.3 \text{ faults} \\ 0.00142 \text{ fails/day/fault} \end{pmatrix}$$

$$K = \frac{\hat{d}}{(\hat{d}_1 + \hat{d}_2)/2} = 31$$

$$V\hat{a}r \begin{pmatrix} \hat{c} \\ \hat{d} \end{pmatrix} = \begin{bmatrix} 27.02^2 & -0.00552 \\ -0.00552 & 0.00027^2 \end{bmatrix}$$

R3 Field Failure Rate Prediction



Possible Sources for Conservative Predictions

- Field Failures often under-reported
- Not all Severity 1/2/3 are failures
- Need to calibrate fault content

$$\hat{\lambda}_{adj}(t) = 58 \times (365 \times 0.00142 / 31) \times e^{-\frac{0.00142 / 31}{1 + 345 \times 45 \times 0.00142 / 31} \times 365 \times 345 t}$$

(fails/year, after t calendar-years of field exposure)

Needed Research

- SRGM fitting tools
- Small sample inference procedures
- Failure rate template
- Improved methods for linking SRGMs to architecture-based models

Clustering by Local Skewering*

David W. Scott

Department of Statistics
Rice University, Houston, Texas 77005
<http://www.stat.rice.edu/~scottdw>

March 31, 2005

Abstract

Clustering p -dimensional data by fitting a mixture of K normals has enjoyed renewed interest (for example, see Splus function “mclust”). However, the number of parameters for the model grows rapidly with dimension p . For example, even if all the covariance matrices are assumed to be equal, the number of parameters is $(K - 1) + K * p + p(p + 1)/2$ for the weights, means and covariance matrix. At ACAS in 2001, Scott introduced the partial mixture component algorithm which fits only one component of the mixture model at a time. This algorithm requires only $1 + p + p * (p + 1)/2$ parameters for the weight, mean vector, and covariance matrix. In this talk, we introduce a new algorithm which attempts to find the “best” line through individual clusters. This model requires only $2 * p - 1$ parameters. That is, the new algorithm is linear rather than quadratic in p . By repeatedly reinitializing the search algorithm, all clusters may be identified. Intuitively, the line found is approximately the largest eigenvector of the local covariance matrix. The GGobi visualization program will be used to illustrate the success of this algorithm on real and simulated data.

1 Introduction

Exploratory data analysis and its development owe much to problems and support of the Army scientists and the Army Research Office. One of the mainstays of exploratory analysis of multivariate data is the principal components technique for dimension reduction. For data $\mathbf{x}_k \in \mathfrak{R}^p$, the sample covariance matrix, S , is estimated and its eigenvalues and eigenvectors computed. The eigenvalues are examined in order to determine the number of dimensions, $p' \ll p$, to retain (through a scree plot, for example). Finally, the data vectors are projected onto the corresponding p' eigenvectors. If the data follow a multivariate normal distribution, even approximately, then investigation of the principal components (rather than the raw data) is extremely useful as a first step.

Even for large dimensions, p , estimation of S is no problem. However, even with today’s computing power, finding all of the eigenstructure often leads to software failure. As an extreme example, the data vector could represent a 1000×1000 gray scale image, so that the covariance matrix is a million by a million. Even a numerically stable approach of avoiding the formation of the covariance matrix

*Research supported in part by NSF grant DMS 02-04723 (non-parametric methodology) and NSF contract EIA-9983459 (digital government). Presented October 29, 2003, Army Conference on Applied Statistics, Napa Valley, California

by computing the singular value decomposition of the data matrix, \mathbf{X} , is not computationally feasible. (Recall that the SVD of $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ and that the eigenvectors of the covariance matrix are contained in the matrix, \mathbf{V} .) Fortunately, there exists specialized software that computes the first p' singular vectors (ARPACT; see Maschhoff and Sorensen, 1996). Since we only require the first p' singular vectors (or eigenvectors), the remaining $p - p'$ eigenvectors need not be computed.

However, this happy situation does not address a number of important problems in data reduction. In particular, multivariate data are much more likely in Army and applied settings to come from a mixture of normal densities, rather than just a single normal density. Of course, statisticians can use the EM algorithm (Dempster et al, 1977) to fit a mixture of normals. But with large problems, estimation of the covariance matrices cannot be avoided. Furthermore, we seem to need to simultaneously estimate the covariance matrices. The SVD approach will not help us here. Whereas ARPACK can find the p' singular vectors, each of length n , the EM approach requires the estimation of K covariance matrices, each of size $n \times n$.

In the following sections, we think about the unthinkable. Can we estimate individual components in a normal mixture without simultaneously having to estimate the other $K - 1$ components? Of course, we are still stuck estimating an $n \times n$ matrix. The second question we consider is the possibility of estimating a few singular vectors without the estimation of S at all. Affirmative answers are shown for both. Computational challenges still remain, but the framework for the optimization problem is provided.

2 Partial Mixture Estimation

Mixture estimation by EM is well-studied; see Titterton et al. (1985). General alternatives to likelihood criteria exist, for example, minimum distance estimation (Beran, 1977). The use of integrated squared error as an estimation criterion has also been considered by Terrell (1990), Basu et al. (1998), and Scott (2001). Given a model, $f_\theta(x)$, and data from the true but unknown density, $g(x)$, we seek to find θ which minimizes

$$\int_{-\infty}^{\infty} [f_\theta(x) - g(x)]^2 dx$$

or

$$\int_{-\infty}^{\infty} f_\theta(x)^2 dx - 2 \int_{-\infty}^{\infty} f_\theta(x) g(x) dx + \int_{-\infty}^{\infty} g(x)^2 dx.$$

An unbiased risk estimate is given by

$$\int_{-\infty}^{\infty} f_\theta(x)^2 dx - \frac{2}{n} \sum_{i=1}^n f_\theta(x_i),$$

where the final term is an unbiased estimate of $2 \int f_\theta(x)g(x)dx$. The integral, $\int g(x)^2 dx$, does not depend upon the unknown parameter, θ , and so may be ignored. If the L_2 norm of the model, $f_\theta(x)$, exists in closed form, then the criterion may easily be minimized numerically. Scott (2001) called the estimator the L_2E estimator, since integrated squared error is in fact the L_2 norm.

Recently, the estimation of normal mixture densities by L_2E was described by Scott (1999, 2004). For example, if the model is the 5-parameter mixture,

$$f_\theta(x) = wN(\mu_1, \sigma_1^2) + (1 - w)N(\mu_2, \sigma_2^2),$$

then the L_2E criterion is easily seen to be

$$\frac{w^2}{2\sqrt{\pi}\sigma_1} + \frac{(1-w)^2}{2\sqrt{\pi}\sigma_2} + 2w(1-w)\phi(0|\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2) - \frac{2}{n} \sum_{i=1}^n f_\theta(x_i).$$

Similar expressions exist in the multivariate normal case.

The L_2E technique has a number of interesting (and unique) features. First, it shares the robustness property of all minimum distance techniques. For example, in Figure 1, a single normal density is fitted to a 2-component mixture by L_2E . Rather than compromising over the two components as in MLE, the L_2E estimator focuses on the larger component, ignoring the smaller component. This practical behavior is our solution to the problem of finding individual components. Wojciechowski and Scott (1999) report a number of simulations comparing L_2E and other robust estimators of location.

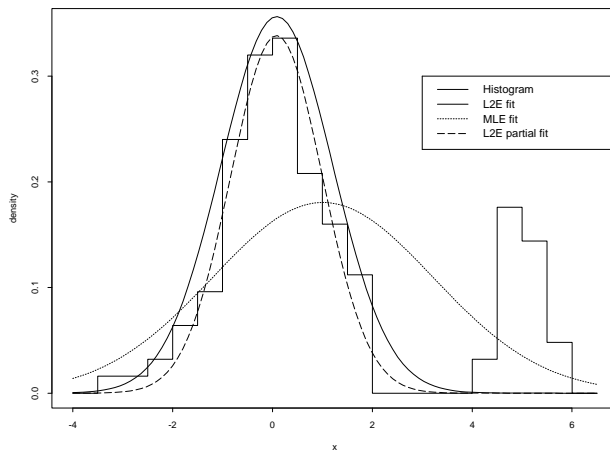


Figure 1: Histogram of 125 points from the mixture $0.8 N(0, 1) + 0.2 N(5, 1)$. Also shown are the maximum likelihood and L_2E fits using the incorrect model $N(\mu, \sigma^2)$. Finally, the L_2E fit of the 3-parameter model $w \cdot N(\mu, \sigma^2)$ is shown.

However, a second and unique feature of L_2E is also displayed in this figure. In the derivation of the L_2E criterion, the fact that $g(x)$ is a true (if unknown) density was of critical importance in order to estimate the integral $\int f_\theta(x) g(x) dx$ in the L_2E expression. However, the fact that the estimator, $f_\theta(x)$, is a true density is not used. Thus, we propose to use estimators that are not (complete) densities. For example, the second L_2E estimator in Figure 1 uses the 3-parameter normal model, $w N(\mu, \sigma^2)$. This equation is called a partial density component (PDC) model. Notice that the area of this PDC L_2E estimate is in fact less than 1.0, and very close to the true value of 0.8 for the left component. (Of special interest is the fact that L_2E can estimate the right component just as well. Which component L_2E converges to is a function of the initial guess for the parameter vector, θ . Since the value of \hat{w} is about 0.20, the usual robust theory about breakdown points never less than 0.50 must be relaxed.)

A similar example, but in two dimensions, is shown in Figure 2, together with the estimated value of \hat{w} . The 6 parameters of the MLE fit (with $w = 1$) were used as initial guesses in the L_2E iterations.

The PDC model can have more than one component. The L_2E estimate found depends entirely upon the initial guess for θ . In practice, a large number of guesses for θ are found by sampling, and the most commonly occurring solutions examined carefully. In Figure 3, we show eight such solutions for the Old Faithful Geyser dataset, which has been lagged and blurred to avoid rounding errors. Clearly these data have three components. Depending upon the choice for θ , the fits may find individual components or combine pairs. Thus we have provided a solution to the vexing problem of mixture estimation when the number of components is unknown. Useful estimates can be found when the number of components is

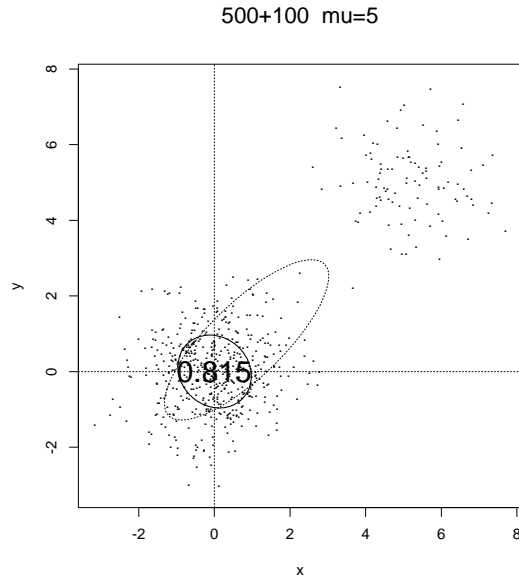


Figure 2: MLE and PDC L_2E contours.

underestimated, even severely. (Of course, if the number of components is overestimated, L_2E will suffer the same fate as MLE and overfitting will result.)

3 Skewers and Principal Components

As is well-known, the first principal component gives the univariate projection of the data with largest variance. In Figure 4, we look at an example of principal components for two variables of the Fisher/Ander-son Iris data. Notice that the axis for the principal components goes through the origin (of course).

Principal components also solves a related problem, which is not often used for motivation. Consider finding a set of points, constrained to lie on a line in \mathbb{R}^p , that are closest to the original data. The solution is provided by the points on the first principal component, where the line is shifted away from the origin to go through the sample mean, $\bar{\mathbf{x}}$. In Figure 5, we show the “skewered” version of the data shown in Figure 4.

Of course, there are 3 species of flowers in the Iris data, so that 3 skewers may be computed. The first principal component for each species, but centered at the mean for each species, is shown in Figure 6. Of course, it is instructive to visualize all four “skewers” for each of the 3 species; see Figure 7.

As instructive as these figures (and animated versions in *ggobi*) are, we are estimating the covariances matrices separately and then computing the eigenvectors of each. Can we find a criterion that is attracted to a skewer without going through the covariance calculation or estimation? Let us look closely at the line segments shown in Figure 5. Clearly, for the Iris Setosa species (for which the skewer was estimated), the distances between the raw data and their projections onto the skewer are quite small, compared to the projections of the Iris Versicolor and Iris Virginica species onto the Setosa skewer. A histogram of these 150 distances is shown in Figure 8. If we compute the Iris Setosa skewer using all 4 variables, we obtain the distance-to-skewer histogram shown in Figure 9.

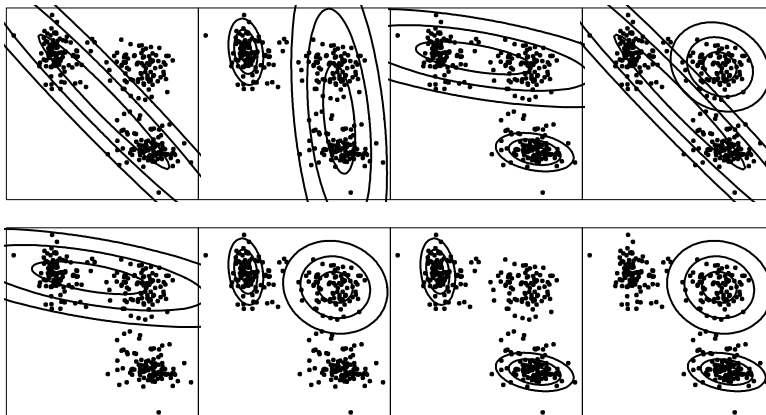


Figure 3: Examples of one- and two-component PDC L_2E fits. (top) The weights in each frame are (.78), (.25, .69), (.68, .28), and (.75, .30). (bottom) The weights in each frame are (.68), (.25, .32), (.25, .28), and (.32, .28).

What are the essentials of a skewer? Like the first principal component, a skewer has a direction, \mathbf{v} . While the principal component goes through the origin, the skewer goes through a general point, P . (If the data are labelled, we know that we can take P to be the sample mean of any single group.) For data in \mathbb{R}^p , the dimension of the point P is p , while the dimension of the direction vector, \mathbf{v} , is $p - 1$. Thus the dimension of the skewer in terms of unknown parameters is $2p - 1$. Thus, the dimension of the search for a skewer grows linearly with dimension, p , rather than quadratically as for the covariance matrix. Thus we have traded a computationally infeasible search for high-dimensional covariance matrices and associated eigenvectors to a linear search for a skewer. However, many random starts will be required in order to have a reasonable chance at finding some number of skewers.

We have not yet specifically stated what the criterion is for finding the skewer, only how we propose to parametrize the search for it. The answer lies in the bimodal structure of the histogram in Figure 9, which should be compared to the bimodal structure in Figure 1. To make our problem easier, imagine that we are more specific about the point, P , on the skewer. Suppose P is the point on the skewer closest to the origin. (Note, we do not advocate using this choice numerically, as instability may arise if the skewer happens to go through the origin, or nearly so.) We now have a vector, \mathbf{u} , which goes from the origin to the new point, P , on the skewer. We can use this vector, \mathbf{u} , in order to create an artificial “sign” on the distance from a data point, \mathbf{x}_k , to its projection onto the skewer, call it \mathbf{y}_k . We do so by taking the inner product of the vector from \mathbf{y}_k to \mathbf{x}_k with the vector \mathbf{u} . Thus the distance histogram as shown in Figure 9 will not have only positive values, but the signed distances of the points corresponding to the skewer will be almost exactly symmetric around the origin. The data points not coming from the skewer group (i.e. the Versicolor and Virginica data in our example) will still be farther away from 0, possibly all on one side, but not necessarily in general. By the robustness property of L_2E , we propose to model the distribution of points “in” the cluster by the PDC model, $wN(0, \sigma^2)$. Note that by fixing the mean at zero, we are asking the skewer to pierce a cluster of the data. The resulting value of w will indicate the rough size of the cluster the skewer has been attracted to. Note, however, that if the data contains 6 clusters, then depending upon the orientation of the first eigenvector of each cluster and the direction to other clusters, it may or may not be possible to isolate each cluster individually. Also, in high dimensions, the use of the normal model needs to be replaced by something closer to a chi-squared

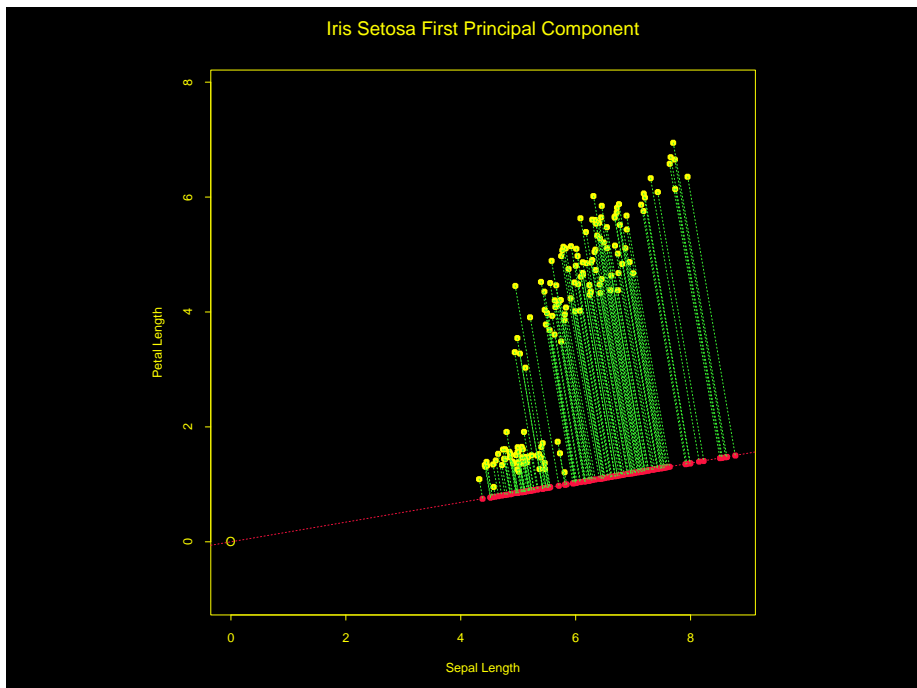


Figure 4: Example of principal components for two variables of the Iris data.

distribution. But if there is a gap in the histogram at the true skewer, then the robustness properties of the L_2E PDC algorithm suggest that a precise model for the PDC is not necessary. (However, the less precise the model, the less precise the estimated values of w and σ will be. Marking points as belonging to the skewer or not relies strongly on at least reasonable value for w and σ) Finally, note that the PDC density model only has 2 parameters, w and σ , no matter the dimension of the data. Of course, the skewer is also part of the estimation, so that the total number of parameters estimated simultaneously by the L_2E PDC skewer algorithm is $(2p - 1) + 2$ or $2p + 1$.

We implemented this algorithm in Splus. For the Iris data (in all 4 dimensions), we found only two skewers. They are shown in Figures 10 and 11 together with the first principal component of the Iris Setosa species. Clearly one estimated skewer is very close to that eigenvector. The other skewer is very close to the first eigenvector for the covariance matrix of all 150 data points (i.e., the overall covariance matrix with no group labels). While a skewer representing just the Versicolor and Virginica species (combined) might be expected, our algorithm always moved away from such an initial orientation to the skewer shown in Figure 11.

4 Extensions

We have limited our discussion to search for one-dimensional skewers. The search for a skewer “plane” or “hyperplane” is a straightforward extension of the algorithm described here. The only difference is that the skewer points, \mathbf{y}_k , lie on a hyperplane rather than a line. The criterion is still the distance from the raw data point, \mathbf{x}_k to the point \mathbf{y}_k on the skewer. Note that our “trick” of constructing a signed

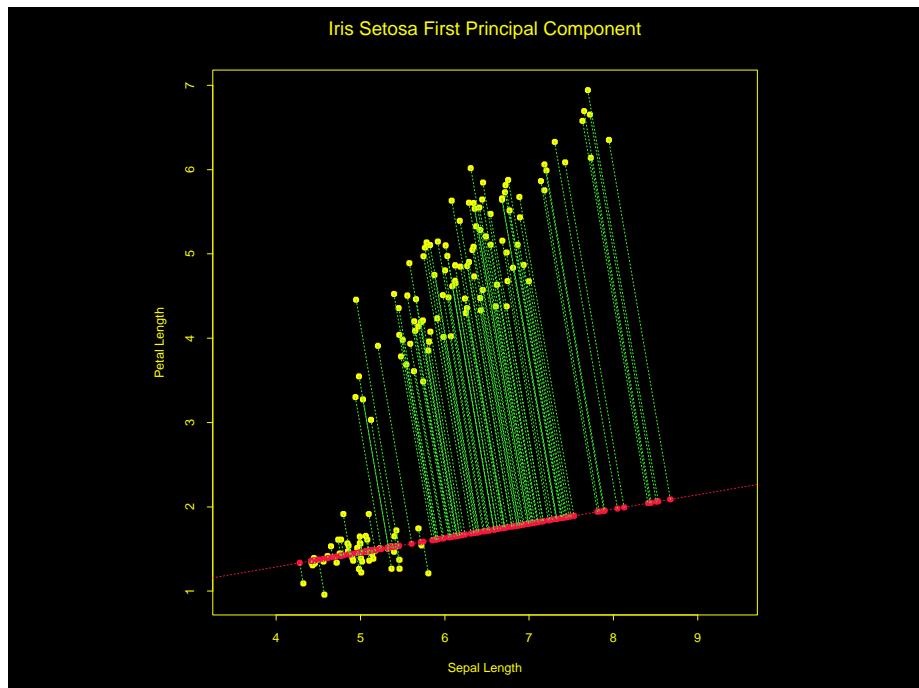


Figure 5: Skewer of the Iris data shown in Figure 4

distance so that the PDC model would be symmetric around the value 0 is now immediately obvious, as the hyperplane divides the space into two parts. The vector \mathbf{u} can be taken as any vector orthogonal to the skewer plane and used to take inner products to assign a sign to each distance computed.

5 Discussion

The ARPACK software allows principal components to be applied to enormously large datasets by avoiding computation of the covariance matrix in order to estimate its eigenvectors. However, if the data set is in fact a mixture of normals, a new attack is required.

In this paper, we have shown how individual mixtures may be estimated without having to estimate or identify all clusters using the L_2E criterion and PDC model. However, such an approach is still quadratic in the number of dimensions, p . But by utilizing a very simple 2-parameter PDC model on the distances from points to their projection onto the skewer, we have demonstrated the existence of a criterion that is linear in the number of dimensions, p . Many random initializations are suggested in order to obtain a reasonable coverage of interesting skewer solutions. But such a task is happily easily accomplished with a farm of parallel computers and requires almost no sophisticated programming tools.

Finally, the Army has a long history of supporting advanced statistical tools and visualization support, beginning with the PRIM9 work of John Tukey and colleagues. The high-dimensional data faced by researchers and workers today requires a whole array of new tools and out-of-the-box thinking. I have tried to illustrate how the use of minimum-distance criteria can free one from the usual set of behaviors into a new realm where seemingly impossible tasks may in fact be successfully addressed.

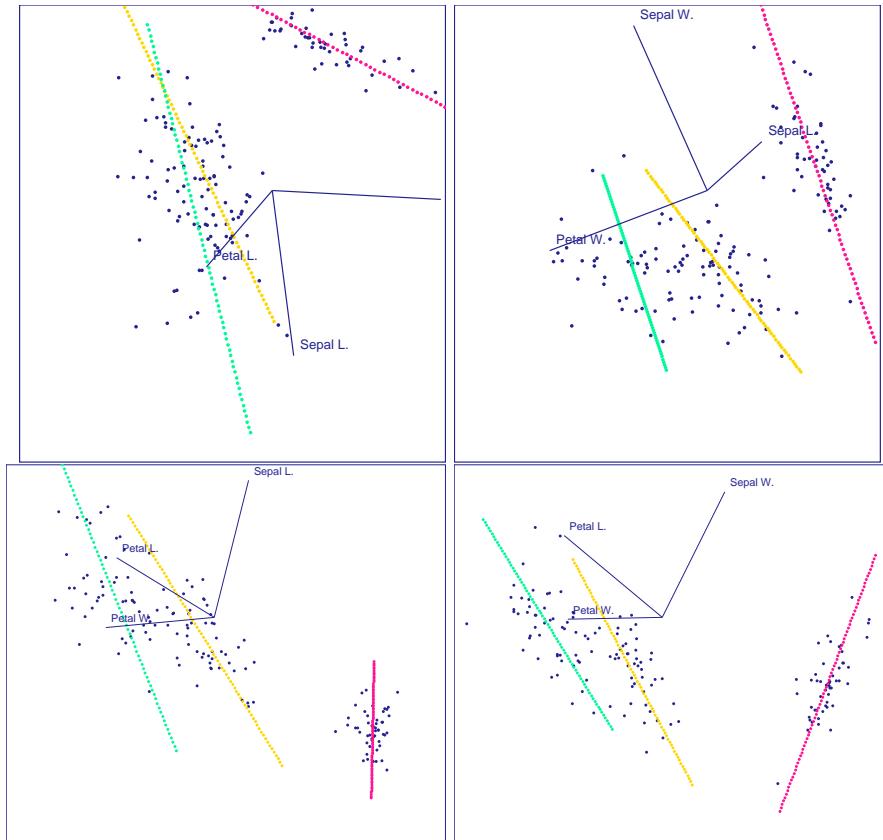


Figure 6: “4-D Skewers” of 3 Iris species in \mathfrak{R}^3 (top left: variables-123; top right: variables-124; bottom left: variables 134; bottom right: variables-234).

6 References

- Basu, A. and Harris, I.R., Hjort, N.L. and Jones, M.C. (1998), “Robust and Efficient Estimation by Minimising a Density Power Divergence,” *Biometrika*, 85, 549–560.
- Beran, R. (1977), “Robust Location Estimates,” *The Annals of Statistics*, 5, 431–444.
- Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977), “Maximum likelihood from incomplete data via the EM algorithm,” *Journal of the Royal Statistical Society, Series B*, 39, 1–37.
- Maschhoff, K.J. and Sorensen, D.C. (1996), “P-ARPACK: An Efficient Portable Large Scale Eigenvalue Package for Distributed Memory Parallel Architectures,” in *PARA*, pp. 478–486.
- Scott, D.W. (1999), “Remarks on Fitting and Interpreting Mixture Models,” *Computing Science and Statistics*, K. Berk and M. Pourahmadi, Eds., 31, 104–109.
- Scott, D.W. (2001), “Parametric Statistical Modeling by Minimum Integrated Square Error,” *Technometrics*, 43, 274–285.
- Scott, D.W. (2004), “Partial Mixture Estimation and Outlier Detection in Data and Regression,” in *Theory and Applications of Recent Robust Methods*, edited by M. Hubert, G. Pison, A. Struyf and S.

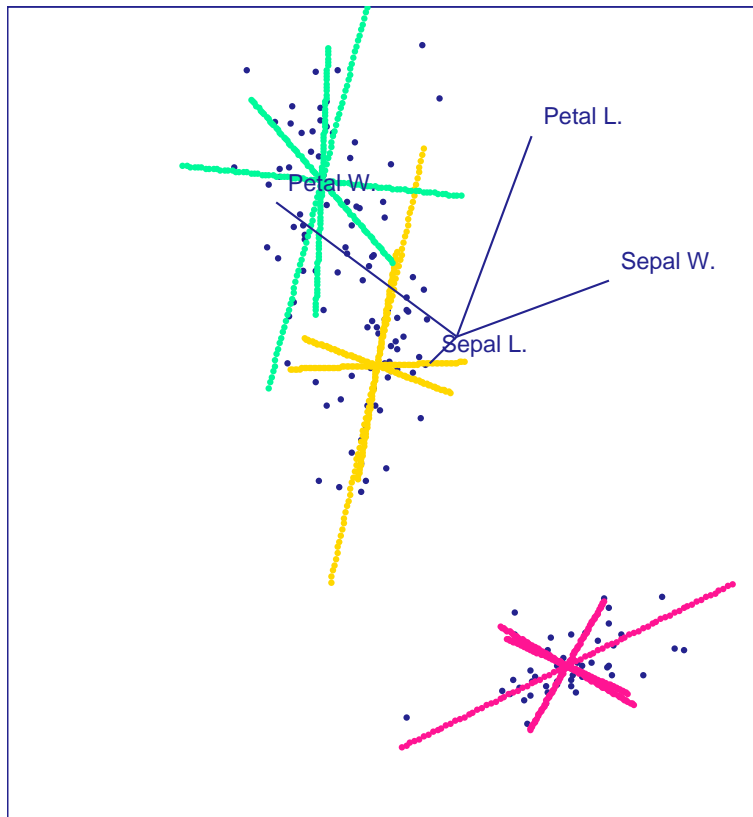


Figure 7: All four eigenvectors for each of the three Iris species in \mathbb{R}^4 (grand tour view in the *ggobi* package).

Van Aelst, Series: Statistics for Industry and Technology, Birkhauser, Basel, pp. 297–306.

Terrell, G.R. (1990), “Linear Density Estimates,” *Proceedings of the Statistical Computing Section*, American Statistical Association, 297–302.

Titterton, D.M., Smith, A.F.M. and Makov, U.E. (1985), *Statistical Analysis of Finite Mixture Distributions*, Wiley, Chichester.

Wojciechowski, W.C. and Scott, D.W. (1999), “Robust Location Estimation with L2 Distance,” *Computing Science and Statistics*, K. Berk and M. Pourahmadi, eds, 31, 292–295.

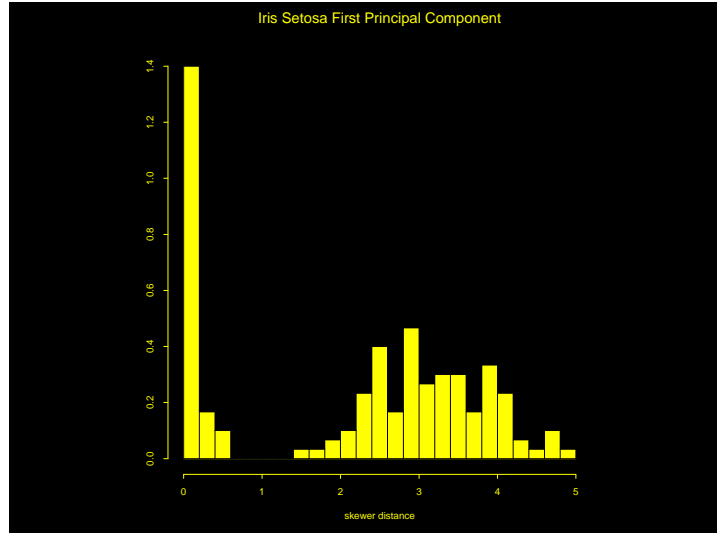


Figure 8: Histogram of distances from the 2-D Iris data to the Iris Setosa Skewer shown in Figure 7.

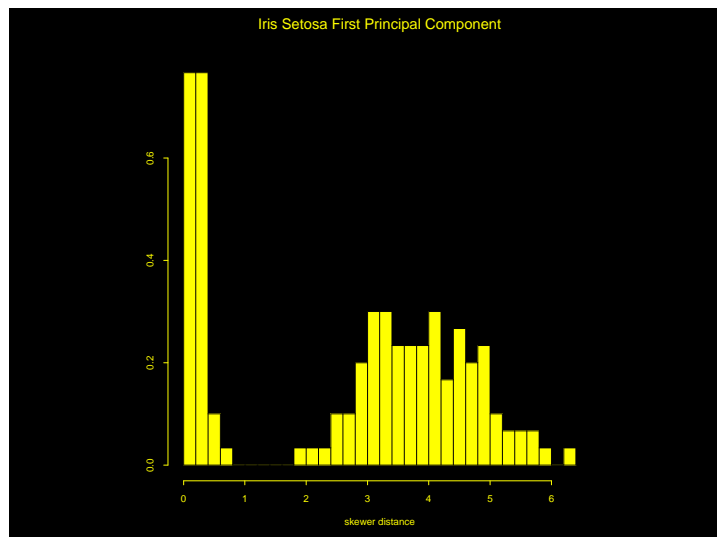


Figure 9: Histogram of distances from the full 4-D Iris data to the 4-D Iris Setosa Skewer.

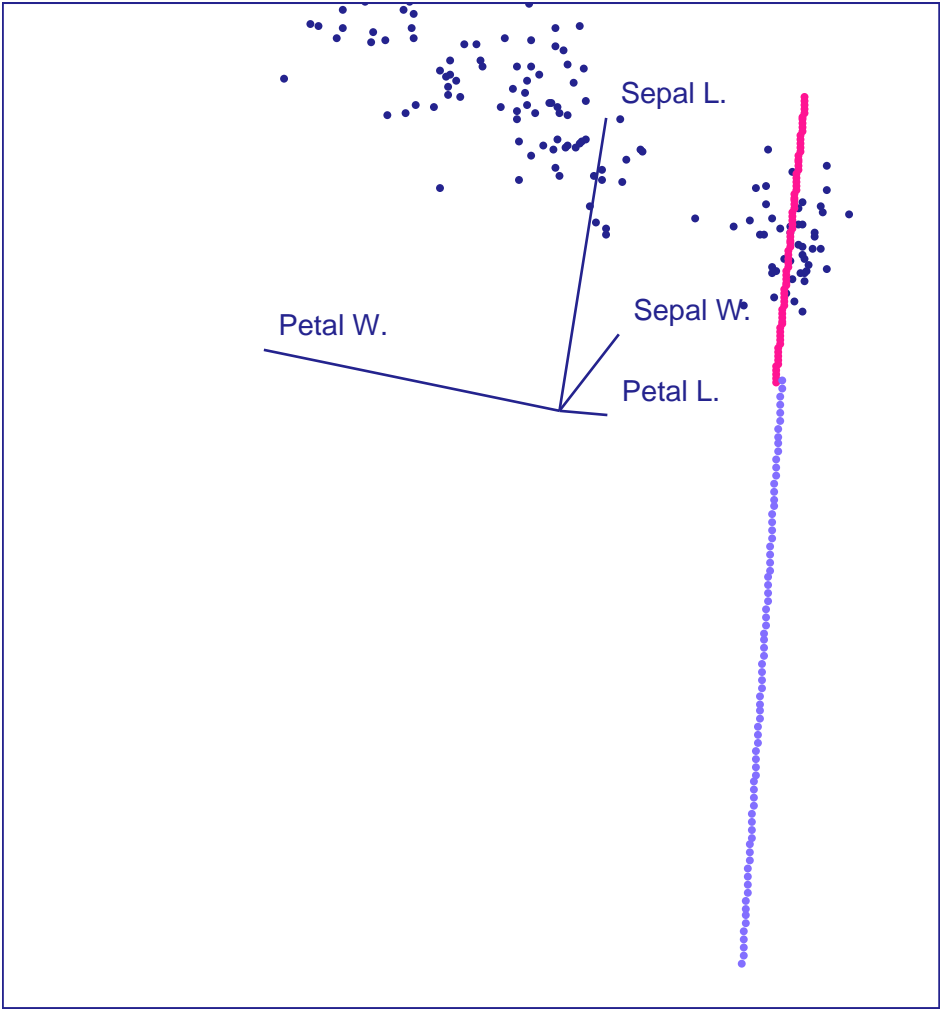


Figure 10: 4-D Skewer of Iris Data. The skewer is blue, while the principal component is red.

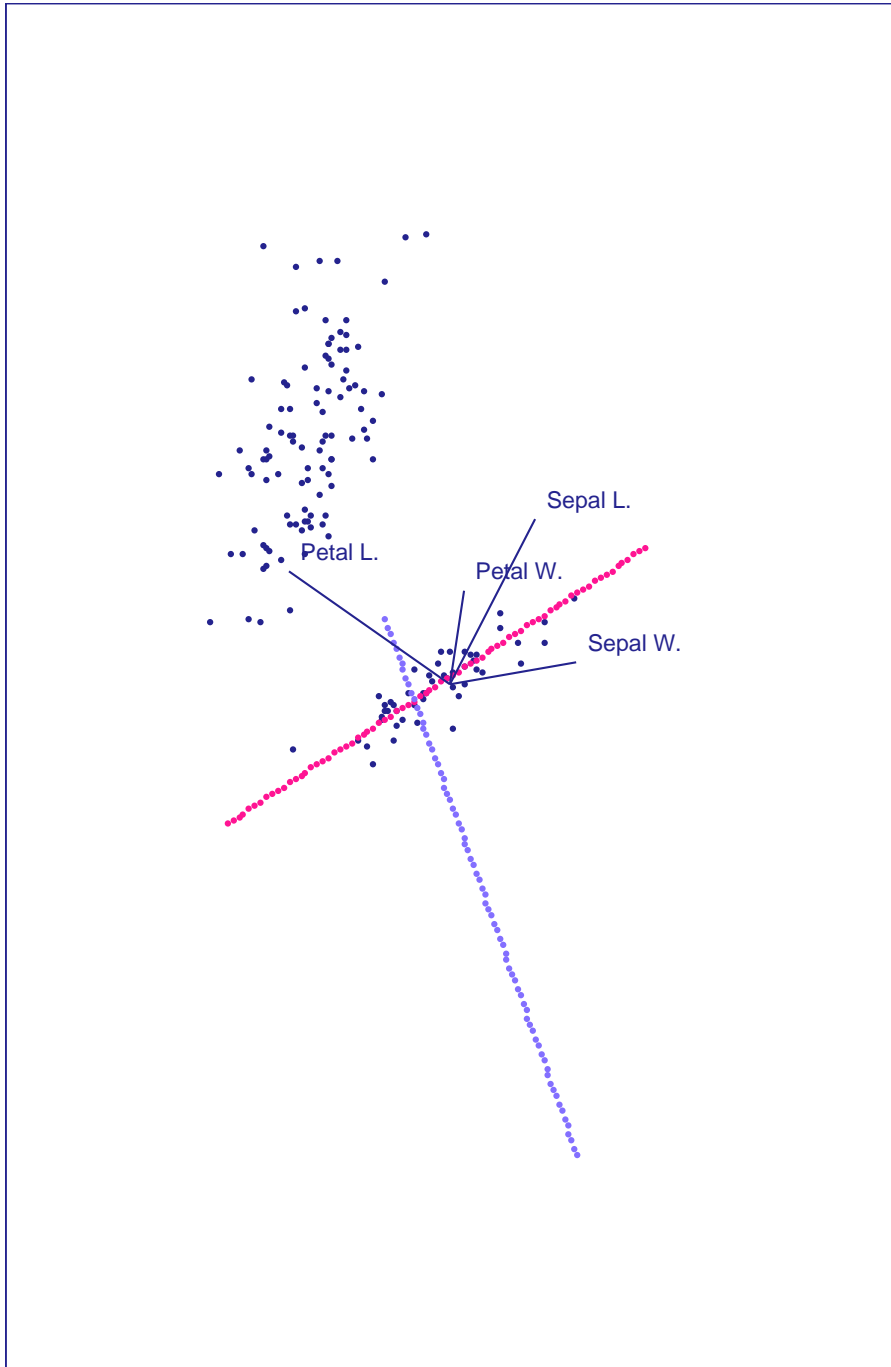


Figure 11: A second 4-D Skewer of Iris Data.

Hereditary Portfolio Optimization With Taxes and Fixed Plus Proportional Transaction Costs: A Quasi-Variational HJB Inequality

Mou-Hsiung Chang
Mathematics Division, U. S. Army Research Office
P. O. Box 12211, Research Triangle Park, North Carolina 27709
telephone: (919) 549-4229, email: mouhsiung.chang@us.army.mil

March 8, 2005

abstract

This paper outlines the problem formulation and results obtained for an infinite time horizon hereditary portfolio optimization problem in a market that consists of one *savings* account and one *stock* account. The *savings* account grows with a fixed interest rate and the unit price of the *stock* account satisfies a nonlinear stochastic hereditary differential equation. Within the *solvency* region the investor is allowed to consume from the *savings* account and can make transactions between the two assets subject to paying capital-gain taxes as well as a fixed plus proportional transaction cost. The *investor* is to seek an optimal consumption-investment strategy in order to maximize the expected utility from the total discounted consumption. The portfolio optimization problem is formulated as a stochastic classical-impulse control problem. A quasi-variational HJB inequality for the value function is derived and the verification theorem for the optimal investment-consumption strategy is obtained. The value function is also shown to be a viscosity solution of the quasi-variational HJB inequality in infinite dimensions.

key words: hereditary portfolio optimization, stochastic hereditary differential equation, optimal stochastic classical-impulse control, quasi-variational HJB inequality (QVHJBI), viscosity solution.

MSC 2000 subject classifications: Primary: 93E20, 91B28; Secondary: 60H30, 49L25, 35R45.

1 Introduction

This paper outlines the problem formulation and the results obtained for an infinite time horizon hereditary portfolio optimization problem in a financial market that consists of one *savings* account and one *stock* account. The full version of the paper containing the detailed derivations and proofs has been submitted elsewhere for publication.

In this paper, it is assumed that the *savings* account grows with a constant interest rate $r > 0$ and $\{S(t), t \geq 0\}$, the unit price process of the *stock* account, satisfies a nonlinear stochastic hereditary differential equation (see (1)) with an infinite but fading memory. In the price dynamics, we assume that both $f(S_t)$ (the mean rate of return) and $g(S_t)$ (the volatility coefficient) depend on the entire history of stock prices S_t over the time interval $(-\infty, t]$ instead of just the current stock price $S(t)$ alone. Within the solvency region \mathcal{S}_κ (see (12)) the *investor* is allowed to consume from his *savings* account in accordance with a consumption rate process $C = \{C(t), t \geq 0\}$ and can make transactions between his *savings* and *stock* accounts according to a transaction (or investment) strategy $\mathcal{T} = \{(\tau(i), \zeta(i)), i = 1, 2, \dots\}$. The proceeds for the sales of the *stock* minus the transaction costs and capital-gain taxes (if the shares of the *stock* are sold at a profit) shall be deposited in his *savings* account and the purchases of shares of the *stock* together with the associated transaction costs shall be financed from his *savings* account. If shares of the *stock* are sold at a loss, then the *investor* shall be given capital-loss credits at the time of the transactions. It is understood that the number of shares of the *stock* transacted can be either integral or fractional.

Throughout the end of the paper, the following rules regarding the transactions and the treatment of capital-gain taxes and capital-loss credits are to be followed.

Rule (1.1). The capital-gain tax and capital-loss credit rates $\beta > 0$ are the same and there is a transaction cost that consists of a fixed cost $\kappa > 0$ plus a proportional transaction cost with the cost rate $\mu > 0$ for buying and selling shares of the *stock*.

Rule (1.2). All the purchases and sales of any number of shares of the stock shall be considered one transaction if they are executed at the same time instance and therefore incurs only one fixed fee $\kappa > 0$ (in addition to a proportional transaction cost).

Rule (1.3). The amount of tax is proportional to the difference between the sale price and the base price with the fixed tax (credit) rate $\beta > 0$ regardless

of whether it is a long-term or short-term gain or loss. It is a capital-gain tax the *investor* has to pay at the time of transaction, if the sale price is higher than the base price. A negative amount of tax shall be interpreted as a capital-loss credit. In this paper, the base price is defined to be the price at which the shares of the *stock* were purchased. Therefore, if $n > 0$ shares of the *stock* are to be sold at time $t \geq 0$, then the unit sale price shall be the current *stock* price $S(t)$ and the unit base price $B(t)$ is set at $S(t - \tau)$ if shares of the *stock* were purchased at a previous time $t - \tau$ with $\tau > 0$. In this case, the total tax due equals $\beta n(S(t) - B(t))$, where

$$\beta n(S(t) - B(t)) = \beta n(S(t) - S(t - \tau)).$$

In the above, it is also assumed that $0 \leq \mu + \beta < 1$.

Rule (1.4). Within the solvency region \mathcal{S}_κ , the *investor* is allowed to borrow money for consumptions and/or *stock* purchases. He can also sell and/or buy-back at the current price shares of the *stock* he bought and/or short-sold at a previous time.

Rule (1.5). The *investor* shall also pay capital-gain taxes (respectively, be paid capital-loss credits) for the amount of profit (respectively, loss) by short-selling shares of the *stock* and then buying back the shares at a lower (respectively, higher) price. The tax shall be paid (or the credit shall be given) at the buying-back time in the amount given by

$$\beta n(S(t) - S(t - \tau))$$

for the number of shares, $n < 0$, of the *stock* the *investor* owed at the previous time $t - \tau$ and at the previous (base) price $S(t - \tau)$ and bought back at the current time t and at the current price $S(t)$.

Rule (1.6). Although the current IRS tax code prohibits wash sales of the *stock*, the *investor* is allowed to sell and/or buy back at a loss shares of the *stock* he owns and/or owes and concurrently buy and/or sell shares of the *stock* at the current price for the sole purpose of taking advantage of the capital-loss credits if it is still profitable after paying the transaction costs. However, the tax and/or credit should not exceed all other gross proceeds (for selling) and/or total costs (for buying) of the shares of the *stock* involved, *i.e.*,

$$n(1 - \mu)S(t) \geq \beta n|S(t) - B(t)| \text{ if } n \geq 0$$

and

$$n(1 + \mu)S(t) \leq \beta n|S(t) - B(t)| \text{ if } n < 0.$$

Under the above assumptions and Rules (1.1)-(1.6), the *investor's* objective is to seek an optimal consumption-investment strategy (C^*, T^*) in order to maximize

$$E\left[\int_0^\infty e^{-\delta t} \frac{C^\gamma(t)}{\gamma} dt\right],$$

the expected utility from the total discounted consumption over the infinite time horizon, where $\delta > 0$ represents the discount rate and $0 < \gamma < 1$ represents the *investor's* risk aversion factor.

Due to the fixed plus proportional transaction costs, the problem shall be formulated as a combination of a classical control (for consumptions) and an impulse control (for the transactions) problem. In this paper a quasi-variational Hamilton-Jacobi-Bellman inequality (QVHJBI) for the value function is derived and the verification theorem for the optimal investment-consumption strategy is obtained. The value function is also shown to be a viscosity solution of the QVHJBI (see QVHJBI (*) in §4.3(D)). Due to the complexity of the analysis involved, the uniqueness result and finite dimensional approximations for the viscosity solution of QVHJBI (*) shall be treated separately in a upcoming paper (see Chang (2004a)) in order to avoid further adding pages to the already lengthy paper.

In recent years there have been extensive amount of research on the optimal consumption-investment problems with proportional transaction costs (see e.g. Akian et al (1996), Akian et al (2001), Davis and Norman (1990), Shreve and Soner (1994) and references contained therein) and fixed plus proportional transaction costs (see e.g. Oksendal and Sulem (2002)) within the geometric Brownian motion financial market. In all these papers, the objective has been to maximize the expected utility from the total discounted or averaged consumption over the infinite time horizon without considering the issues of capital-gain taxes (respectively, capital-loss credits) when shares of the *stock* are sold at a profit (respectively, loss). In different contents, the issues of capital-gain taxes have been studied in Cadenillas and Pliska (1999), Constantinides (1983), Constantinides (1984), Dammon and Spatt (1996) and references contained therein. In particular, Constantinides (1983) and Constantinides (1984) considered the effect of capital-gain taxes and capital-loss credits on capital market equilibrium without consumption and transaction costs. These two papers illustrated that under some conditions, it may be more profitable to cut one's losses short and never to realize a gain because of capital-loss credits and capital-gain taxes as some conventional wisdom will suggest. In Cadenillas and Pliska (1999) the optimal

transaction time problem with proportional transaction costs and capital-gain taxes was considered in order to maximize the the long-run growth rate of the investment (or the so-called Kelley criterion), *i.e.*,

$$\lim_{t \rightarrow \infty} \frac{1}{t} E[\log V(t)],$$

where $V(t)$ is the value of the investment measured at time $t > 0$. This paper is quite different from ours in that the unit price of the *stock* is described by a geometric Brownian motion, and all shares of the stock owned by the investor are to be sold at a chosen transaction time and all of its proceeds from the sale are to be used to purchase new shares of the *stock* immediately after the sale without consumption. Fortunately due to the nature of the geometric Brownian motion market, the authors of that paper were able to obtain some explicit results. In Jouini et al (1999) and Jouini et al (2000), a nonclassical deterministic optimal control problem with *endogenous delay* was considered in maximizing a total cost function over a finite time horizon when the investor has a given deterministic income function for investment (without transaction costs) in a riskless asset that grows with a given variable interest rate function. The investor re-balances his portfolio (by buying or selling some financial assets) and spends the rest for consumptions. They obtained the existence and maximum principle for optimal choice of a deterministic purchasing plan under the first-come-first-out rule (*i.e.*, the shares that were purchased first should be sold first if the investor decide to sell some of his asset). Again, these two papers are quite different from ours in that the investment strategy is a absolutely continuous function of time and there is no randomness involved in the model.

To the best of the author's knowledge, this is the first paper that treats the consumption-investment problem in which the hereditary nature of the stock price dynamics and the issue of capital-gain taxes are taken into consideration. Due to drastically different nature of the problem and the techniques involved, the hereditary portfolio optimization problem with taxes and proportional transaction costs (*i.e.*, $\kappa = 0$ and $\mu > 0$) remains to be solved and is the subject of the author's upcoming research (see Chang (2004b)).

This paper is organized as follows. The description of the stock price dynamics, the admissible consumption-investment strategies, and the formulation of the hereditary portfolio optimization problem are given in section two. In section three, the behavior of the controlled state process is fur-

ther explored and corresponding infinite dimensional Markovian solution of the price dynamics is investigated. Section four contains the derivations of a Bellman-type dynamic programming principle and the QVHJBI together with its boundary conditions. A verification theorem for the optimal consumption-investment strategy is established in section five. It is demonstrated that the value function is discontinuous at interfaces of some parts of the boundary of the solvency region and hence can not satisfy the QVHJBI in the classical sense. A weaker concept of viscosity solution is introduced and defined. It is shown that the value function is a viscosity solution of the QVHJBI. All these are done in section six.

2 The Hereditary Portfolio Optimization Problem

2.1 Basic Notations and Preliminary Analysis

(A). The Past History Space $M_{\rho,+}^2$ for Price Dynamics.

Throughout the end of this paper, let $\rho : \mathfrak{R}_- \rightarrow \mathfrak{R}_+$ ($\mathfrak{R}_- \equiv (-\infty, 0]$ and $\mathfrak{R}_+ \equiv [0, \infty)$) be the *influence function with relaxation property* and satisfies the following conditions:

Condition (2.1.1). ρ is summable on \mathfrak{R}_- , i.e., $0 < \int_{-\infty}^0 \rho(\theta) d\theta < \infty$.

Condition (2.1.2). For every $\lambda \leq 0$ one has

$$\bar{K}(\lambda) = \text{ess sup}_{\theta \in \mathfrak{R}_-} \frac{\rho(\theta + \lambda)}{\rho(\theta)} \leq \bar{K} < \infty,$$

$$\underline{K}(\lambda) = \text{ess sup}_{\theta \in \mathfrak{R}_-} \frac{\rho(\theta)}{\rho(\theta + \lambda)} < \infty.$$

Condition (2.1.3). ρ is essentially bounded on \mathfrak{R}_- .

Condition (2.1.4). ρ is essentially strictly positive on $(-\infty, 0)$.

Condition (2.1.5). $\theta\rho(\theta) \rightarrow 0$ as $\theta \rightarrow -\infty$.

Note that Conditions (2.1.3)-(2.1.5) are consequences of Conditions (2.1.1)-(2.1.2) (see Coleman and Mizel (1966)).

Let $M_{\rho}^2 \equiv \mathfrak{R} \times L_{\rho}^2(\mathfrak{R}_-)$ be the past history space of the stock price dynamics, where $L_{\rho}^2(\mathfrak{R}_-)$ is the class of measurable functions $\phi : \mathfrak{R}_- \rightarrow \mathfrak{R}$ such that

$$\int_{-\infty}^0 |\phi(\theta)|^2 \rho(\theta) d\theta < \infty.$$

Note that \mathbf{M}_ρ^2 is a real separable Hilbert space of functions $(\phi(0), \phi) \in \mathfrak{R} \times L_\rho^2(\mathfrak{R}_-)$ equipped with the Hilbertian inner product $\langle \cdot, \cdot \rangle_M : \mathbf{M}_\rho^2 \times \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ defined by

$$\langle (\phi(0), \phi), (\varphi(0), \varphi) \rangle_M = \phi(0)\varphi(0) + \int_{-\infty}^0 \phi(\theta)\varphi(\theta)\rho(\theta)d\theta.$$

As usual, the Hilbertian norm $\| \cdot \|_M : \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}_+$ is defined by

$$\|(\phi(0), \phi)\|_M = \sqrt{\langle (\phi(0), \phi), (\phi(0), \phi) \rangle_M}.$$

Let $\mathbf{M}_{\rho,+}^2$ be the subspace of \mathbf{M}_ρ^2 defined by

$$\mathbf{M}_{\rho,+}^2 = \{(\phi(0), \phi) \in \mathbf{M}_\rho^2 \mid \phi(\theta) \geq 0 \forall \theta \in \mathfrak{R}_-\}.$$

If $t \geq 0$ and $\phi : \mathfrak{R} \rightarrow \mathfrak{R}$ is a measurable function, define $\phi_t : \mathfrak{R}_- \rightarrow \mathfrak{R}$ by $\phi_t(\theta) = \phi(t + \theta)$, $\theta \in \mathfrak{R}_-$.

Therefore, if $\phi : \mathfrak{R} \rightarrow \mathfrak{R}$ is such that

$$\int_{-\infty}^{\infty} |\phi(\theta)|^2 \rho(\theta) d\theta < \infty,$$

then $(\phi(t), \phi_t) \in \mathbf{M}_\rho^2$ for each $t \geq 0$.

If $\{S(t), t \in \mathfrak{R}\}$ is the unit stock price process that satisfies (1) with the initial historical price function $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$, then it can be shown that the $\mathbf{M}_{\rho,+}^2$ -valued process $\{(S(t), S_t), t \geq 0\}$ is a strong Markov process. Again, $S_t : \mathfrak{R}_-, \mathfrak{R}_+$ is defined by $S_t(\theta) = S(t + \theta)$ for each $\theta \in (-\infty, 0]$. Note that the *stock* price dynamics described by (1) is said to have an infinite but fading memory because, for each $t \geq 0$ the norm of $(S(t), S_t)$, $\|(S(t), S_t)\|_M$, depend on its entire past history up to time t in a weighted fashion by the function $\rho : \mathfrak{R}_- \rightarrow \mathfrak{R}_+$ satisfying Conditions (2.1.1)-(2.1.5).

(B). The Function Space \mathbf{N} for Stock Inventory.

Let \mathbf{N} denote the space of bounded functions $\xi : (-\infty, 0] \rightarrow \mathfrak{R}$ of the following form

$$\xi(\theta) = \sum_{k=0}^{\infty} n(\tau(-k)) \mathbf{1}_{\{\tau(-k)\}}(\theta), \quad \theta \in (-\infty, 0]$$

where $\{n(\tau(-k)), k = 0, 1, 2, \dots\}$ is a bounded sequence in \mathfrak{R} with

$$\sum_{k=0}^{\infty} |n(\tau(-k))| < \infty, \quad -\infty < \dots < \tau(-k) < \dots < \tau(-1) < \tau(0) = 0,$$

and $\mathbf{1}_{\{\tau(-k)\}}$ is the indicator function at $\tau(-k)$. In another words, $\xi(\theta) = n(\tau(-k))$ if $\theta = \tau(-k)$ and $= 0$ if $\theta \neq \tau(-k) \forall k = 0, 1, 2, \dots$. For notational simplicity, $\xi \in \mathbf{N}$ expressed above can and shall sometimes be represented by the pair of two sequences

$$\xi = \{(\tau(-k), n(\tau(-k))), k = 0, 1, 2, \dots\},$$

or simply $\xi = \{n(\tau(-k)), k = 0, 1, 2, \dots\}$ if there is no danger of ambiguity.

Let $\|\cdot\|_N$ (the norm of the space \mathbf{N}) be defined by

$$\|\xi\|_N = \sup_{\theta \in \mathfrak{R}_-} |\xi(\theta)| \quad \forall \xi \in \mathbf{N}.$$

If $\eta : (-\infty, \infty) \rightarrow \mathfrak{R}$ is a bounded function of the form similar to that of \mathbf{N} , *i.e.*,

$$\eta(t) = \sum_{k=-\infty}^{\infty} n(\tau(k)) \mathbf{1}_{\{\tau(k)\}}(t),$$

(or equivalently $\eta = \{(\tau(k), n(\tau(k))), k = \dots, -1, 0, 1, \dots\}$) such that

$$\sum_{k=-\infty}^{\infty} |n(\tau(k))| < \infty$$

and

$$-\infty < \dots < \tau(-k) < \dots < 0 = \tau(0) \leq \tau(1) < \dots < \tau(k) < \dots < \infty,$$

then for each $t \geq 0$ we define, with the same notation as in §2.1(A), *i.e.*, $\eta_t \in \mathbf{N}$ by $\eta_t(\theta) = \eta(t + \theta)$, $\theta \in (-\infty, 0]$. In this case,

$$\begin{aligned} \eta_t(\theta) &= \sum_{k=-\infty}^{Q(t)} n(\tau(k)) \mathbf{1}_{\{\tau(k)\}}(\theta) \\ &= \sum_{k=-\infty}^{\infty} n(\tau(k)) \mathbf{1}_{\{\tau(k)\}}(t + \theta), \end{aligned}$$

or simply

$$\begin{aligned} \eta_t &= \{n(\tau(k)), k = \dots, -1, 0, 1, \dots, Q(t)\} \\ &= \{n(\tau(Q(t) - k)), k = 0, 1, 2, \dots\}, \end{aligned}$$

where $Q(t) = \sup\{k \geq 0 \mid \tau(k) \leq t\}$. Again with a little abuse of notation, we shall write throughout the end of this paper that $n(\tau(k)) = n(k)$ when

there is no danger of ambiguity.

We assume the following convention and assumption:

Assumption (2.1.6). The sequence $\xi = \{n(-k), k = 0, 1, \dots\}$ is such that $n(0) = 0$ and $n(-k) = 0$ for all but finitely many k .

As illustrated in §2.3, \mathbf{N} is the space in which the *investor's* stock inventory lives. The Assumption (2.1.6) implies that the *investor* can only have shares of the same stock that were purchased or short-sold at a finite number of previous time instances. However, this finite number may increase from time to time if the *investor* does not sell all shares of what he owns and/or buy back all shares of what he owes.

2.2 Hereditary Price Structure with Infinite Memory

For $t \in (-\infty, \infty)$, let $S(t)$ denote the unit price of the *stock* at time t . It is assumed that the unit *stock* price process $\{S(t), t \in (-\infty, \infty)\}$ satisfy the following stochastic hereditary differential equation with an infinite but fading memory:

$$\frac{dS(t)}{S(t)} = f(S_t)dt + g(S_t)dW(t), \quad t \geq 0, \quad (1)$$

or

$$dS(t) = \tilde{f}(S(t), S_t)dt + \tilde{g}(S(t), S_t)dW(t), \quad t \geq 0, \quad (2)$$

and the initial price function $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$, where $\tilde{f}, \tilde{g} : \mathbf{M}_{\rho,+}^2 \rightarrow \mathfrak{R}_+$ are defined by

$$\tilde{f}(\phi(0), \phi) = \phi(0)f(\phi) \text{ and } \tilde{g}(\phi(0), \phi) = \phi(0)g(\phi). \quad (3)$$

In the above equations, the process $\{W(t), t \geq 0\}$ is an one-dimensional standard Brownian motion defined on a complete filtered probability space $(\Omega, \mathcal{F}, P; \mathbf{F})$, where $\mathbf{F} = \{\mathcal{F}(t), t \geq 0\}$ is the P -augmented natural filtration generated by the Brownian motion $\{W(t), t \geq 0\}$, *i.e.*,

$$\mathcal{F}(t) = \sigma(W(s), 0 \leq s \leq t) \vee \mathcal{N},$$

and

$$\mathcal{N} = \{A \subset \Omega \mid \exists B \in \mathcal{F} \text{ such that } A \subset B \text{ and } P(B) = 0\}.$$

In the above, $f(S_t)$ and $g(S_t)$ represent, respectively, the *mean growth rate* and the *volatility rate* of the *stock price* at time $t \geq 0$. Note that the *stock* is said to have a hereditary price structure with infinite memory because at time $t \geq 0$ both $f(S_t)$ and $g(S_t)$ explicitly depend on the entire past history of *stock prices* over the time interval $(-\infty, t]$ instead of the *stock price* $S(t)$ at time t alone.

Since security exchanges have only existed since a finite past, it is realistic but not technically required to assume that the initial historical price function $(\psi(0), \psi)$ to have the property that

$$\psi(\theta) = 0 \quad \forall \theta \leq \bar{\theta} < 0 \text{ for some } \bar{\theta} < 0.$$

Although the modelling of *stock prices* is still under intensive investigations, it is not the intention of this paper to address the validity of the model *stock price dynamics* treated in this paper but to illustrate the optimal consumption-investment problem that explicitly dependent upon the entire past history of the *stock prices* for computing capital-gain taxes or capital-loss credits. The term "hereditary portfolio optimization" is therefore coined in this paper for the first time.

We note here that the stochastic hereditary differential equation of the type described in (2) was studied in Mizel and Trutzer (1984) with its applications not in financial market but in modelling the behavior of some viscoelastic material in mind. Some special forms of hereditary *stock prices* with bounded memory have been studied in Chang and Youree (1999) and Arriojas et al (2003) for the pricing of European options.

It is assumed for simplicity that the functions $f, g : L^2_{\rho,+}(\mathfrak{R}_-) \rightarrow \mathfrak{R}_+$ are continuous, and satisfy the following Lipschitz and linear growth conditions in order to ensure the existence and uniqueness of a *strong* solution process $\{S(t), t \geq 0\}$ with the initial historical price function $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}^2_{\rho,+}$ (see e.g. Mizel and Trutzer (1984), Mohammed (1984), Mohammed (1996) and Arriojas (1997) for the theory of stochastic functional differential equations with an infinite or a bounded memory).

We make the following assumptions regarding the functions $f, g : L^2_{\rho,+}(\mathfrak{R}_-) \rightarrow \mathfrak{R}_+$.

Assumption (2.2.1). The functions \tilde{f} and \tilde{g} satisfy the following Lipschitz condition:

$$\begin{aligned}
|\tilde{f}(\phi(0), \phi) - \tilde{f}(\varphi(0), \varphi)| + |\tilde{g}(\phi(0), \phi) - \tilde{g}(\varphi(0), \varphi)| \\
\leq c_1 \|(\phi(0), \phi) - (\varphi(0), \varphi)\|_M;
\end{aligned}$$

Assumption (2.2.2). There exist positive constants α , and σ such that

$$0 < r < f(\phi) \leq \alpha \text{ and } 0 < \sigma \leq g(\phi); \text{ and}$$

Assumption (2.2.3). There exist constants $c_1, c_2 > 0$ such that

$$\begin{aligned}
0 \leq \tilde{f}(\phi(0), \phi) + \tilde{g}(\phi(0), \phi) \leq c_2(1 + \|(\phi(0), \phi)\|_M) \\
\forall (\phi(0), \phi), (\varphi(0), \varphi) \in \mathbf{M}_{\rho,+}^2.
\end{aligned}$$

Note that the lower bound of the *mean rate of return* $f : L_{\rho,+}^2(\mathfrak{R}_-) \rightarrow \mathfrak{R}$ in Assumption (2.2.2) is imposed to make sure that the *stock* account has a higher mean growth rate than the interest rate $r > 0$ for the *savings* account. Otherwise, it will be more profitable and less risky for the *investor* to put all his money in the *savings* account for the purpose of optimizing the expected utility from the total consumption.

For each initial historical price function $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$, the price process $\{S(t), t \geq 0\}$ is a positive, continuous, and \mathbf{F} -adapted process defined on $(\Omega, \mathcal{F}, P; \mathbf{F})$ (see Theorem (2.1) in Mizel and Trutzer (1984)). Using the notation adapted in §2.1(A), we frequently consider the corresponding $\mathbf{M}_{\rho,+}^2$ -valued process $\{(S(t), S_t), t \geq 0\}$, where $(S(0), S_0) = (\psi(0), \psi)$. It can be shown under Conditions (2.1.1)-(2.1.5) and Assumptions (2.2.1)-(2.2.3) (see §3 of Mizel and Trutzer (1984)) that the $\mathbf{M}_{\rho,+}^2$ -valued process $\{(S(t), S_t), t \geq 0\}$ is strong Markovian with respect to the filtration \mathbf{G} , where $\mathbf{G} = \{\mathcal{G}(t), t \geq 0\}$ is the filtration generated by $\{S(t), t \geq 0\}$, *i.e.*,

$$\mathcal{G}(t) = \sigma(S(s), 0 \leq s \leq t) (= \sigma((S(s), S_s), 0 \leq s \leq t)), \forall t \geq 0.$$

2.3 Consumption-Investment Strategies

Let $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$ be an initial historical price function of the *stock* over the interval $(-\infty, 0]$. It is assumed that immediately prior to the initial time $t = 0$ the investor inherited $x \in \mathfrak{R}$ dollars in his *savings* account and an *inventory* of the shares of the *stock* $\xi = \sum_{k=0}^{\infty} n(-k) \mathbf{1}_{\{\tau(-k)\}} \in \mathbf{N}$, where $n(-k) \equiv n(\tau(-k)) > 0$ (respectively, $n(-k) \equiv n(\tau(-k)) < 0$) denotes the number of shares the investor owns (respectively, owes) that were originally

purchased (respectively, short-sold) at the previous time $\tau(-k) < 0$ and at the base price $\psi(\tau(-k))$. Within the *solvency region* \mathcal{S}_κ (see (12)) the *investor* is allowed to consume from his *savings* account and can make transactions between his *savings* and *stock* account according to a consumption-investment strategy $\pi = (C, \mathcal{T})$.

Definition (2.3.1). The pair $\pi = (C, \mathcal{T})$ is said to be a consumption-investment strategy if the consumption rate process

(i) $C = \{C(t), t \geq 0\}$ is a non-negative \mathbf{G} -progressively measurable process such that

$$\int_0^T C(t)dt < \infty \quad P - a.s. \quad \forall T > 0;$$

and

(ii) $\mathcal{T} = \{(\tau(i), \zeta(i)), i = 1, 2, \dots\}$ is the transaction strategy with $\tau(i), i = 1, 2, \dots$, being an increasing sequence of \mathbf{G} -stopping times such that

$$\lim_{i \rightarrow \infty} \tau(i) = \infty \quad a.s.$$

and

$$\zeta(i), i = 1, 2, \dots,$$

being a sequence of \mathbf{N} -valued $\mathcal{G}(\tau(i))$ -measurable random variables. Note that

$$0 = \tau(0) \leq \tau(1) < \dots < \tau(i) < \dots$$

denotes the sequence of transaction times. The transaction amount at time $\tau(i), i = 1, 2, \dots$, is an \mathbf{N} -valued $\mathcal{G}(\tau(i))$ -measurable random vector given by

$$\zeta(i) = \{m(i-k), k = 0, 1, \dots, \},$$

where $m(i) > 0$ (respectively, $m(i) < 0$) is the number of shares of the stock the *investor* bought (respectively, sold) at the current time $\tau(i)$ and, for $k = 1, 2, \dots$, $m(i-k) > 0$ is the number of shares of the stock bought back at current time $\tau(i)$, and that were all or part of the shares short-sold at the previous time $\tau(i-k)$. Similarly, $m(i-k) < 0$ is the number of shares sold at the current time $\tau(i)$, and that were all or part of the shares purchased at the previous time $\tau(i-k)$.

The effect of an instantaneous transaction at time $\tau(i), i = 1, 2, \dots$ on the *investor's* current portfolio in his *savings* account and *stock* account is illustrated as follows. Taking into the account of stationarity of the

state equations (see (13)-(15), (16)-(17), and (21) or (1)), let us suppose without loss of generality the portfolio (or position) of the *investor* is at $(x, \xi, \psi(0), \psi) \in \mathfrak{R} \times \mathbf{N} \times \mathbf{M}_{\rho,+}^2$, where $x \in \mathfrak{R}$ denote the *investor's* holdings in his *savings* account, $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$ is the historical stock prices, and $\xi = \sum_{k=0}^{\infty} n(-k) \mathbf{1}_{\{\tau(-k)\}} \in \mathbf{N}$ (with $n(0) = 0$) denotes the *inventory* of the *investor's* holdings in the *stock* account in that $n(-k) \equiv n(\tau(-k)) > 0$ (respectively, < 0) represents the number of shares of the *stock* that were purchased (respectively, short-sold) at $\tau(-k)$ and that are still owned (respectively, owed). Under the transaction rules, costs, and taxes as described in section one, we define for each $\xi \in \mathbf{N}$ the constraint set at ξ , $\mathcal{R}(\xi) \subset \mathbf{N}$, by

$$\begin{aligned} \mathcal{R}(\xi) = \{ \zeta \in \mathbf{N} \mid \zeta = \sum_{k=0}^{\infty} m(-k) \mathbf{1}_{\{\tau(-k)\}}, -\infty < m(0) < \infty, \text{ and} \quad (4) \\ \text{either } n(-k) > 0, m(-k) \leq 0 \quad \& \quad n(-k) + m(-k) \geq 0 \\ \text{or } n(-k) < 0, m(-k) \geq 0 \quad \& \quad n(-k) + m(-k) \leq 0 \text{ for } k \geq 1 \}. \end{aligned}$$

An instantaneous transaction of the \mathbf{N} -valued quantity

$$\begin{aligned} \zeta = m(0) \mathbf{1}_{\{\tau(0)\}} + \sum_{k=1}^{\infty} m(-k) \mathbf{1}_{\{\tau(-k)\}} (\chi_{\{n(-k) < 0, 0 \leq m(-k) \leq -n(-k)\}} \\ + \chi_{\{n(-k) > 0, -n(-k) \leq m(-k) \leq 0\}}) \in \mathcal{R}(\xi) \end{aligned}$$

leads to a new state of the portfolio $(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi})$, where \hat{x} (the *investor's* new holdings in his *savings* account), $\hat{\xi}$ (the *investor's* new inventory in his *stock* account), and $(\hat{\psi}(0), \hat{\psi})$ (the new profile of stock prices) are given below.

$$\begin{aligned} \hat{x} = x - \kappa - (m(0) + \mu |m(0)|) \psi(0) - \sum_{k=1}^{\infty} \left[(1 - \mu - \beta) m(-k) \psi(0) \quad (5) \right. \\ \left. + \beta m(-k) \psi(\tau(-k)) \right] \chi_{\{n(-k) > 0, -n(-k) \leq m(-k) \leq 0\}} \\ - \sum_{k=1}^{\infty} \left[(1 + \mu - \beta) m(-k) \psi(0) + \beta m(-k) \psi(\tau(-k)) \right] \\ \times \chi_{\{n(-k) < 0, 0 \leq m(-k) \leq -n(-k)\}}, \end{aligned}$$

$$\hat{\xi} = \xi \oplus \zeta \quad (6)$$

and

$$(\hat{\psi}(0), \hat{\psi}) = (\psi(0), \psi) \quad (7)$$

where $\chi_{\{\dots\}}$ is the indicator function of the set (or event) $\{\dots\}$, and $\xi \oplus \zeta : (-\infty, 0] \rightarrow \mathfrak{R}$ is defined by

$$\begin{aligned} (\xi \oplus \zeta)(\theta) &= m(0)\mathbf{1}_{\{\tau(0)\}}(\theta) \\ &+ \sum_{k=1}^{\infty} \left[n(-k) + m(-k)(\chi_{\{n(-k) < 0, 0 \leq m(-k) \leq -n(-k)\}} \right. \\ &\left. + \chi_{\{n(-k) > 0, -n(-k) \leq m(-k) \leq 0\}}) \right] \mathbf{1}_{\{\tau(-k)\}}(\theta), \quad \forall \theta \in (-\infty, 0], \end{aligned} \quad (8)$$

or simply by the sequence

$$\xi \oplus \zeta = \{\hat{n}(-k), k = 0, 1, 2, \dots\}, \quad (9)$$

with

$$\hat{n}(0) = m(0),$$

and

$$\begin{aligned} \hat{n}(-k) &= n(-k) + m(-k)(\chi_{\{n(-k) < 0, 0 \leq m(-k) \leq -n(-k)\}} \\ &+ \chi_{\{n(-k) > 0, -n(-k) \leq m(-k) \leq 0\}}) \text{ for } k = 1, 2, \dots \end{aligned}$$

Remark (2.3.2). The reason that any instantaneous transaction of the amount $\xi \in \mathbf{N}$ has to be taken from the constraint set $\mathcal{R}(\xi)$ is due to the fact that the *investor* can purchase ($m(0) > 0$) or short-sell ($m(0) < 0$) new shares of the *stock* but can only be allowed to buy-back some or all shares ($n(-k) < 0$ & $0 \leq m(-k) \leq -n(-k)$) of what he short-sold and/or sell some or all shares ($n(-k) > 0$ & $-n(-k) \leq m(-k) \leq 0$) of what he purchased previously at the current price as specified in Rules (1.1)-(1.6).

We observe the following:

1. \hat{x} , the new holdings in his *savings* account, is obtained from his previous holding x minus the fixed cost κ , the cost for new purchases together with proportional transaction cost $(m(0) + \mu|m(0)|)\psi(0)$, the total cost for buying back some or all shares of what he owed with proportional transaction cost and tax paid *i.e.*,

$$\sum_{k=1}^{\infty} \left[(1 + \mu - \beta)m(-k)\psi(0) + \beta m(-k)\psi(\tau(-k)) \right] \chi_{\{n(-k) < 0, 0 \leq m(-k) \leq -n(-k)\}},$$

and plus the net proceeds (deducting the proportional transaction cost and tax) for selling some or all shares of what he owned, *i.e.*,

$$- \sum_{k=1}^{\infty} \left[(1 - \mu - \beta)m(-k)\psi(0) + \beta m(-k)\psi(\tau(-k)) \right] \chi_{\{n(-k) > 0, -n(-k) \leq m(-k) \leq 0\}};$$

2. $\hat{\xi} = \xi \oplus \zeta$ is the *investor's* new inventory in his *stock* account after making the transaction of amount $\zeta \in \mathcal{R}(\xi)$;
3. $(\hat{\psi}(0), \hat{\psi}) = (\psi(0), \psi)$ due to the fact that unit price of the *stock* is continuous with respect to time; and
4. $n(-k) > 0 \Rightarrow \hat{n}(-k) \equiv n(-k) + m(-k)\chi_{\{n(-k) \geq 0, -n(-k) \leq m(-k) \leq 0\}} \geq 0$ and $n(-k) < 0 \Rightarrow \hat{n}(-k) \equiv n(-k) + m(-k)\chi_{\{n(-k) < 0, 0 \leq m(-k) \leq -n(-k)\}} \leq 0$.

2.4 Solvency Region

If $(x, \xi, \psi(0), \psi) \in \mathfrak{R} \times \mathbf{N} \times \mathbf{M}_{\rho,+}^2$ is the current *portfolio* (or *state*) of the *investor*, he can borrow money for consumptions and/or for purchases of the *stock* and can also short-sell and/or buy back shares of the *stock* at the current price as long as the liquidated value of his portfolio remains non-negative. Again, writing $n(\tau(-k))$ as $n(-k)$ for simplicity, define the function $H_\kappa : \mathfrak{R} \times \mathbf{N} \times \mathbf{M}_{\rho,+}^2 \rightarrow \mathfrak{R}$ as follows.

$$H_\kappa(x, \xi, \psi(0), \psi) = \max \left\{ G_\kappa(x, \xi, \psi(0), \psi), \min\{x, n(-k), k = 0, 1, 2, \dots\} \right\}, \quad (10)$$

where

$$G_\kappa(x, \xi, \psi(0), \psi) = x - \kappa + \sum_{k=0}^{\infty} \left[(n(-k) - \mu|n(-k)|)\psi(0) - \beta n(-k)(\psi(0) - \psi(\tau(-k))) \right]. \quad (11)$$

Note that $G_\kappa(x, \xi, \psi(0), \psi)$ represents the cash value (if the assets can be liquidated at all) after selling all shares of the *stock* he owns and buying back all the shares of the *stock* he owes with all transaction costs (fixed plus proportional transactional costs) and taxes paid.

The *solvency region* \mathcal{S}_κ of the portfolio optimization problem is defined as

$$\begin{aligned} \mathcal{S}_\kappa &= \left\{ (x, \xi, \psi(0), \psi) \in \mathfrak{R} \times \mathbf{N}_+ \times \mathbf{M}_{\rho,+}^2 \mid H_\kappa(x, \xi, \psi(0), \psi) \geq 0 \right\} \\ &= \left\{ (x, \xi, \psi(0), \psi) \in \mathfrak{R} \times \mathbf{N}_+ \times \mathbf{M}_{\rho,+}^2 \mid G_\kappa(x, \xi, \psi(0), \psi) \geq 0 \right\} \\ &\quad \cup (\mathfrak{R}_+ \times \mathbf{N}_+ \times \mathbf{M}_{\rho,+}^2). \end{aligned} \quad (12)$$

Note that within the *solvency region* \mathcal{S}_κ there are shares of the *stock* that can not be liquidated at all, namely, those $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ such that

$$(x, \xi, \psi(0), \psi) \in \mathfrak{R}_+ \times \mathbf{N}_+ \times \mathbf{M}_{\rho,+}^2 \text{ and } G_\kappa(x, \xi, \psi(0), \psi) < 0.$$

The in-liquidity of these shares of the stock is due to the insufficient fund to pay for the transaction costs and/or taxes, etc. Observe that the *solvency* region \mathcal{S}_κ is an unbounded and non-convex subset of the state space $\mathfrak{R} \times \mathbf{N} \times \mathbf{M}_{\rho,+}^2$. The boundary $\partial\mathcal{S}_\kappa$ and other properties of the *solvency* region \mathcal{S}_κ shall be described in detail in §4.3.

2.5 Portfolio Dynamics and Admissible Strategies

At time $t \geq 0$, the investor's portfolio in the financial market shall be denoted by the triplet $(X(t), N_t, S(t), S_t)$, where $X(t)$ denotes the *investor's* holdings in his *savings* account, $N_t \in \mathbf{N}$ is the *inventory* of his *stock* account, and $(S(t), S_t)$ describes the profile of the unit prices of the *stock* over the past history $(-\infty, t]$ as described in §2.2.

Given the initial portfolio

$$(X(0-), N_{0-}, S(0), S_0) = (x, \xi, \psi(0), \psi) \in \mathfrak{R} \times \mathbf{N} \times \mathbf{M}_{\rho,+}^2$$

and applying a consumption-investment strategy $\pi = (C, \mathcal{T})$ (see Definition (2.3.1)), the portfolio dynamics of $\{(X(t), N_t, S(t), S_t), t \geq 0\}$ can then be described as follows.

Firstly, the *savings* account holdings $\{X(t), t \geq 0\}$ satisfies the following equations:

$$dX(t) = [rX(t) - C(t)]dt, \quad \tau(i) \leq t < \tau(i+1), \quad i = 0, 1, 2, \dots, \quad (13)$$

$$\begin{aligned} X(\tau(i)) &= X(\tau(i)-) - \kappa - (m(i) + \mu|m(0)|)S(\tau(i)) & (14) \\ &- \sum_{k=1}^{\infty} \left[(1 - \mu - \beta)m(i-k)S(\tau(i)) \right. \\ &\quad \left. + \beta m(i-k)S(\tau(i-k)) \right] \\ &\cdot \chi_{\{n(i-k) > 0, -n(i-k) \leq m(i-k) \leq 0\}} \\ &- \sum_{k=1}^{\infty} \left[(1 + \mu - \beta)m(i-k)S(\tau(i)) \right. \\ &\quad \left. + \beta m(i-k)S(\tau(i-k)) \right] \\ &\cdot \chi_{\{n(i-k) < 0, 0 \leq m(i-k) \leq -n(i-k)\}}, \quad \text{for } i = 0, 1, 2, \dots, \end{aligned}$$

and the initial holding in his *savings* account

$$X(\tau(0)-) = X(0-) = x. \quad (15)$$

As a reminder $m(-i) > 0$ (respectively, $m(-i) < 0$) means buying (respectively, selling) new shares of the *stock* and $m(i-k) > 0$ (respectively, $m(i-k) < 0$) means buying back (respectively, selling) of some or all of what he owed (respectively, owned).

Secondly, the inventory of the *investor's* stock account at time $t \geq 0$, $N_t \in \mathbf{N}$, evolves according to the following equations:

$$N_t = N_{\tau(i)} = \sum_{k=-\infty}^{Q(t)} n(k) \mathbf{1}_{\tau(k)} \quad \text{if } \tau(i) \leq t < \tau(i+1), i = 0, 1, 2, \dots, \quad (16)$$

where $Q(t) = \sup\{k \geq 0 \mid \tau(k) \leq t\}$, and

$$N_{\tau(i)} = N_{\tau(i)-} \oplus \zeta(i) \quad (17)$$

$N_{\tau(i)-} \oplus \zeta(i) : (-\infty, 0] \rightarrow \mathbf{N}$ is defined by

$$\begin{aligned} (N_{\tau(i)-} \oplus \zeta(i))(\theta) &= m(i) \mathbf{1}_{\{\tau(i)\}}(\theta) + \sum_{k=1}^{\infty} \left[n(i-k) \right. \\ &\quad \left. + m(i-k) (\chi_{\{n(i-k) < 0, 0 \leq m(i-k) \leq -n(i-k)\}} \right. \\ &\quad \left. + \chi_{\{n(i-k) > 0, -n(i-k) \leq m(i-k) \leq 0\}}) \right] \mathbf{1}_{\{\tau(i-k)\}}(\theta), \end{aligned} \quad (18)$$

or equivalently the sequence

$$N_{\tau(i)-} \oplus \zeta(i) = \{\hat{n}(i-k), k = 0, 1, 2, \dots\}$$

where

$$\hat{n}(i) = m(i), \text{ and} \quad (19)$$

$$\begin{aligned} \hat{n}(i-k) &= n(i-k) + m(i-k) (\chi_{\{n(i-k) < 0, 0 \leq m(i-k) \leq -n(i-k)\}} \\ &\quad + \chi_{\{n(i-k) \geq 0, -n(i-k) \leq m(i-k) \leq 0\}}), \end{aligned} \quad (20)$$

for $i = 0, 1, 2, \dots$ and $k = 0, 1, 2, \dots, i$.

Note that the sequence $-\infty < \dots < \tau(-k) < \dots < \tau(-1) < \tau(0) = 0$ is as previously given for the initial ξ and $n(i-k)$ is the number of shares of the *stock* owned at time $\tau(i)$ but were initially purchased (if $n(i-k) > 0$) or short-sold (if $n(i-k) < 0$) $\tau(k)$ at the previous time $\tau(i-k)$.

Thirdly, since the *investor* is small, the unit stock price process $\{S(t), t \geq 0\}$ will not be in anyway affected by the *investor's* action in the market and is again described as in (1) by

$$\frac{dS(t)}{S(t)} = f(S_t)dt + g(S_t)dW(t), \quad t \geq 0, \quad (21)$$

with an initial historical price function $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$.

Definition (2.5.1). If the *investor* starts with an initial portfolio

$$(X(0-), N_{0-}, S(0), S_0) = (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa.$$

The consumption-investment strategy $\pi = (C, \mathcal{T})$ defined in Definition (2.3.1) is said to be *admissible* at $(x, \xi, \psi(0), \psi)$ if

$$\zeta(i) \in \mathcal{R}(N_{\tau(i)-}) \quad \forall i = 1, 2, \dots$$

and

$$(X(t), N_t, S(t), S_t) \in \mathcal{S}_\kappa, \quad \forall t \geq 0.$$

The class of consumption-investment strategies admissible at $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ shall be denoted by $\mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$.

2.6 The Problem Statement

Given the initial state $(X(0-), N_{0-}, S(0), S_0) = (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$, the *investor's* objective is to find an admissible consumption-investment strategy $\pi^* \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$ that maximizes the following expected utility from the total discounted consumption:

$$J_\kappa(x, \xi, \psi(0), \psi; \pi) = E^{x, \xi, \psi(0), \psi; \pi} \left[\int_0^\infty e^{-\delta t} \frac{C^\gamma(t)}{\gamma} dt \right] \quad (22)$$

among the class of admissible consumption-investment strategies $\mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$, where $E^{x, \xi, \psi(0), \psi; \pi}[\cdot \cdot \cdot]$ is the expectation with respect to $P^{x, \xi, \psi(0), \psi; \pi}\{\cdot \cdot \cdot\}$, the probability measure induced by the controlled (by π) state process $\{(X(t), N_t, S(t), S_t), t \geq 0\}$ and conditioned on the initial state

$$(X(0-), N_{0-}, S(0), S_0) = (x, \xi, \psi(0), \psi).$$

In the above, $\delta > 0$ denotes the discount factor, and $0 < \gamma < 1$ indicates that the utility function $U(c) = \frac{c^\gamma}{\gamma}$, for $c > 0$, is a function of HARA (hyperbolic absolute risk aversion) type that were considered in most of optimal consumption-investment literature (see e.g. Davis and Norman (1990),

Akian et al (1996), Akian et al (2001), Shreve and Soner (1994), and Ok-sendal and Sulem (2002)) with or without a fixed transaction cost. The admissible (consumption-investment) strategy $\pi^* \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$ that maximizes $J_\kappa(x, \xi, \psi(0), \psi; \pi)$ is called an optimal (consumption-investment) strategy and the function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ defined by

$$\begin{aligned} V_\kappa(x, \xi, \psi(0), \psi) &= \sup_{\pi \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)} J_\kappa(x, \xi, \psi(0), \psi; \pi) \\ &= J_\kappa(x, \xi, \psi(0), \psi; \pi^*) \end{aligned} \quad (23)$$

is called the value function of the hereditary portfolio optimization problem.

The optimal consumption-investment problem (or the hereditary portfolio optimization problem) considered in this paper is then formalized as follows.

Problem (2.6.1). For each given initial state $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$, identify the optimal strategy π^* and its corresponding value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$.

3 The Controlled State Process

Given an initial state $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ and an admissible consumption-investment strategy $\pi = (C, \mathcal{T}) \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$, the \mathcal{S}_κ -valued controlled state process shall sometimes be denoted by $\{Z^{x, \xi, \psi(0), \psi; \pi}(t), t \geq 0\}$ (or simply $\{Z^\pi(t) = (X^\pi(t), N_t^\pi, S^\pi(t), S_t^\pi), t \geq 0\}$ or $\{Z(t) = (X(t), N_t, S(t), S_t), t \geq 0\}$ when there is no danger of ambiguity), where

$$Z^{x, \xi, \psi(0), \psi; \pi}(t) = (X^{x, \xi, \psi(0), \psi; \pi}(t), N_t^{x, \xi, \psi(0), \psi; \pi}, S^{x, \xi, \psi(0), \psi; \pi}(t), S_t^{x, \xi, \psi(0), \psi; \pi}).$$

The main purpose of this section is to establish the Markovian and other properties such as the Dynkin's formula for the controlled state process $\{Z^{x, \xi, \psi(0), \psi; \pi}(t), t \geq 0\}$. Note that the $\mathbf{M}_{\rho, +}^2$ -valued process $\{S(t), S_t, t \geq 0\}$ described by (1) is uncontrollable by the *investor* and is therefore independent of the consumption-investment strategy $\pi \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$ but is dependent on the initial historical price function $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}_{\rho, +}^2$.

3.1 The Holdings in the Savings Account

Given an initial state $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ and an admissible consumption-investment strategy $\pi = (C, T) \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$, it can be shown that the process $\{X(t), t \geq 0\}$ described by (13)-(14) is a Markov RCLL (right-continuous with a finite left-hand limit) real-valued strong Markov process with the change-of-variable formula given as follows (see Protter (1995) or Rogers and Williams (1987) for jumped Markov process):

$$\begin{aligned} e^{-\delta\tau}\Phi(X(\tau)) &= \Phi(x) + \int_0^\tau e^{-\delta t}[(rX(t) - C(t))\Phi_x(X(t))]dt \quad (24) \\ &\quad + \sum_{0 \leq t \leq \tau} e^{-\delta t}[\Phi(X(t)) - \Phi(X(t-))], \end{aligned}$$

for all $\Phi \in C_b^1(\mathfrak{R})$ (the space of bounded and continuously differentiable functions on \mathfrak{R}) and finite \mathbf{G} -stopping time τ , where Φ_x denotes the derivative of Φ (with respect to x) and $X(\tau(i)-) = \lim_{t \downarrow 0} X(\tau(i) - t)$.

3.2 The Inventory of the Stock Account

Similarly to the process $\{X(t), t \geq 0\}$, the \mathbf{N} -valued controlled inventory process $\{N_t, t \geq 0\}$ of the *investor's* stock account described by (16)- (17) also satisfies the following change-of-variable formula:

$$e^{-\delta\tau}\Phi(N_\tau) = \Phi(\xi) + \sum_{0 \leq t \leq \tau} e^{-\delta t}[\Phi(N_t) - \Phi(N_{t-})], \quad (25)$$

for all $\Phi \in C_b(\mathbf{N})$ (the space of bounded and continuous function from \mathbf{N} to \mathfrak{R}), and finite \mathbf{G} -stopping time τ , where $N_{\tau-} = \lim_{t \downarrow 0} N_{\tau-t}$.

3.3 The Properties of the Stock Prices

To study the Markovian properties of the $\mathbf{M}_{\rho,+}^2$ -valued solution process $\{(S(t), S_t), t \geq 0\}$ where $S_t(\theta) = S(t + \theta), \theta \in (-\infty, 0]$, and $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$, we need the following notation and ancillary results.

(A). Preliminary Results on (\mathbf{M}_ρ^2) .

Let $(\mathbf{M}_\rho^2)^*$ be the space of bounded linear functionals (or the topological

dual of the space \mathbf{M}_ρ^2) equipped with the operator norm $\|\cdot\|_M^*$ defined by

$$\|\Phi\|_M^* = \sup_{(\phi(0), \phi) \neq (0, \mathbf{0})} \frac{|\Phi(\phi(0), \phi)|}{\|(\phi(0), \phi)\|_M}.$$

For the benefit of the readers who are not familiar with the theory of infinite dimensional Hilbert space we note that $(\mathbf{M}_\rho^2)^*$ can be identified with \mathbf{M}_ρ^2 by the Riesz representation theorem re-stated below:

Theorem (3.3.1). $\Phi \in (\mathbf{M}_\rho^2)^*$ if and only if there exists a unique $(\varphi(0), \varphi) \in \mathbf{M}_\rho^2$ such that

$$\begin{aligned} \Phi(\phi(0), \phi) &= \langle (\varphi(0), \varphi), (\phi(0), \phi) \rangle_M \\ &\equiv \varphi(0)\phi(0) + \int_{-\infty}^0 \varphi(\theta)\phi(\theta)\rho(\theta)d\theta \quad \forall (\phi(0), \phi) \in \mathbf{M}_\rho^2. \end{aligned}$$

Let $(\mathbf{M}_\rho^2)^\dagger$ be the space of bounded bilinear functionals $\Phi : \mathbf{M}_\rho^2 \times \mathbf{M}_\rho^2 \rightarrow \Re$ (*i.e.*, $\Phi((\phi(0), \phi), (\cdot, \cdot)), \Phi((\cdot, \cdot), (\phi(0), \phi)) \in (\mathbf{M}_\rho^2)^*$ for each $(\phi(0), \phi) \in \mathbf{M}_\rho^2$), equipped with the operator norm $\|\cdot\|_M^\dagger$ defined by

$$\begin{aligned} \|\Phi\|_M^\dagger &= \sup_{(\phi(0), \phi) \neq (0, \mathbf{0})} \frac{\|\Phi((\cdot, \cdot), (\phi(0), \phi))\|_M^*}{\|(\phi(0), \phi)\|_M} \\ &= \sup_{(\phi(0), \phi) \neq (0, \mathbf{0})} \frac{\|\Phi((\phi(0), \phi), (\cdot, \cdot))\|_M}{\|(\phi(0), \phi)\|_M}. \end{aligned}$$

Let $\Phi : \mathbf{M}_\rho^2 \rightarrow \Re$. The function Φ is said to be Fréchet differentiable at $(\phi(0), \phi) \in \mathbf{M}_\rho^2$ if for each $(\varphi(0), \varphi) \in \mathbf{M}_\rho^2$,

$$\Phi((\phi(0), \phi) + (\varphi(0), \varphi)) - \Phi(\phi(0), \phi) = D\Phi(\phi(0), \phi)(\varphi(0), \varphi) + o(\|(\varphi(0), \varphi)\|_M),$$

where $D\Phi : \mathbf{M}_\rho^2 \rightarrow (\mathbf{M}_\rho^2)^*$ and $o : \Re \rightarrow \Re$ is a function such that

$$\frac{o(\|(\varphi(0), \varphi)\|_M)}{\|(\varphi(0), \varphi)\|_M} \rightarrow 0 \text{ as } \|(\varphi(0), \varphi)\|_M \rightarrow 0.$$

In this case, $D\Phi(\phi(0), \phi) \in (\mathbf{M}_\rho^2)^*$ is called the (first order) Fréchet derivative of Φ at $(\phi(0), \phi) \in \mathbf{M}_\rho^2$. The function Φ is said to be continuously Fréchet differentiable if its Fréchet derivative $D\Phi : \mathbf{M}_\rho^2 \rightarrow (\mathbf{M}_\rho^2)^*$ is continuous under the operator norm $\|\cdot\|_M^*$. The function Φ is said to be twice Fréchet differentiable at $(\phi(0), \phi) \in \mathbf{M}_\rho^2$ if its Fréchet derivative $D\Phi(\phi(0), \phi) : \mathbf{M}_\rho^2 \rightarrow \Re$

exists and there exists a bounded bilinear functional $D^2\Phi(\phi(0), \phi) : \mathbf{M}_\rho^2 \times \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ where for each $(\varphi(0), \varphi), (\varsigma(0), \varsigma) \in \mathbf{M}_\rho^2$

$$D^2\Phi(\phi(0), \phi)((\cdot, \cdot), (\varphi(0), \varphi)), D^2\Phi(\phi(0), \phi)((\varsigma(0), \varsigma), (\cdot, \cdot)) \in (\mathbf{M}_\rho^2)^*,$$

and where

$$\begin{aligned} & \left(D\Phi((\phi(0), \phi) + (\varphi(0), \varphi)) - D\Phi(\phi(0), \phi) \right) (\varsigma(0), \varsigma) \\ &= D^2\Phi(\phi(0), \phi)((\varsigma(0), \varsigma), (\varphi(0), \varphi)) + o(\|(\varsigma(0), \varsigma)\|_M, \|(\varphi(0), \varphi)\|_M). \end{aligned}$$

Here, $o : \mathfrak{R} \times \mathfrak{R} \rightarrow \mathfrak{R}$ is such that

$$\frac{o(\cdot, \|(\varphi(0), \varphi)\|_M)}{\|(\varphi(0), \varphi)\|_M} \rightarrow 0, \text{ as } \|(\varphi(0), \varphi)\|_M \rightarrow 0$$

and

$$\frac{o(\|(\varphi(0), \varphi)\|_M, \cdot)}{\|(\varphi(0), \varphi)\|_M} \rightarrow 0 \text{ as } \|(\varphi(0), \varphi)\|_M \rightarrow 0.$$

In this case, the bounded bilinear functional $D^2\Phi(\phi(0), \phi) : \mathbf{M}_\rho^2 \times \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ is the second order Fréchet derivative of Φ at $(\phi(0), \phi) \in \mathbf{M}_\rho^2$.

Throughout the end, we let $C^2(\mathbf{M}_\rho^2)$ be the space of functions $\Phi : \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ that are twice continuously Fréchet differentiable. The 2nd-order Fréchet derivative $D^2\Phi$ is said to be globally Lipschitz on \mathbf{M}_ρ^2 if there exists a constant $K > 0$ such that

$$\begin{aligned} \|D^2\Phi(\phi(0), \phi) - D^2\Phi(\varphi(0), \varphi)\|_M^\dagger &\leq K\|(\phi(0), \phi) - (\varphi(0), \varphi)\|_M, \\ \forall (\phi(0), \phi), (\varphi(0), \varphi) &\in \mathbf{M}_\rho^2. \end{aligned}$$

The space of $\Phi \in C^2(\mathbf{M}_\rho^2)$ with $D^2\Phi$ being globally Lipschitz will be denoted by $C_{lip}^2(\mathbf{M}_\rho^2)$.

If $\Phi \in C^2(\mathbf{M}_\rho^2)$, then the actions of the first order Fréchet derivative $D\Phi(\phi(0), \phi)$ and the 2nd order Fréchet $D^2\Phi(\phi(0), \phi)$ can be expressed as

$$D\Phi(\phi(0), \phi)(\varphi(0), \varphi) = \varphi(0)\partial_{\phi(0)}\Phi(\phi(0), \phi) + D_\phi\Phi(\phi(0), \phi)\varphi,$$

and

$$\begin{aligned} & D^2\Phi(\phi(0), \phi)((\varphi(0), \varphi), (\varsigma(0), \varsigma)) \\ &= \varphi(0)\partial_{\phi(0)}^2\Phi(\phi(0), \phi)\varsigma(0) + \varsigma(0)\partial_{\phi(0)}D_\phi\Phi(\phi(0), \phi)\varphi \\ &+ \varphi(0)D_\phi\partial_{\phi(0)}\Phi(\phi(0), \phi)\varsigma + D_\phi^2\Phi(\phi(0), \phi)(\varphi, \varsigma), \end{aligned}$$

where $\partial_{\phi(0)}\Phi$ and $\partial_{\phi(0)}^2\Phi$ are the first and 2nd order partial derivatives of Φ with respect to its first variable $\phi(0) \in \mathfrak{R}$, $D_\phi\Phi$ and $D_\phi^2\Phi$ are the first and 2nd order Fréchet derivatives with respect to its second variable $\phi \in L_\rho^2(\mathfrak{R}_-)$, $\partial_{\phi(0)}D_\phi\Phi$ is the second order derivative first with respect to ϕ in the Fréchet sense and then with respect to $\phi(0)$, *etc.*.

(B). The Weak Infinitesimal Generator Γ .

For each $\phi \in L_\rho^2(\mathfrak{R}_-)$, define $\tilde{\phi} : (-\infty, \infty) \rightarrow \mathfrak{R}$ by

$$\tilde{\phi}(t) = \begin{cases} \phi(0) & \text{for } t \in [0, \infty), \\ \phi(t) & \text{for } t \in (-\infty, 0). \end{cases}$$

Then for each $\theta \in (-\infty, 0]$ and $t \in [0, \infty)$,

$$\tilde{\phi}_t(\theta) = \tilde{\phi}(t + \theta) = \begin{cases} \phi(0) & \text{for } t + \theta \geq 0, \\ \phi(t + \theta) & \text{for } t + \theta < 0. \end{cases}$$

A bounded measurable function $\Phi : \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$, *i.e.*, $\Phi \in C_b(\mathbf{M}_\rho^2)$, is said to belong to $\mathcal{D}(\Gamma)$, the domain of the weak infinitesimal operator Γ , if the following limit exists for each fixed $(\phi(0), \phi) \in \mathbf{M}_\rho^2$:

$$\Gamma(\Phi)(\phi(0), \phi) \equiv \lim_{t \downarrow 0} \frac{\Phi(\phi(0), \tilde{\phi}_t) - \Phi(\phi(0), \phi)}{t}. \quad (26)$$

Remark (3.3.2). Note that $\Phi \in C_{lip}^2(\mathbf{M}_\rho^2)$ does not guarantee that $\Phi \in \mathcal{D}(\Gamma)$. For example, let $\bar{\theta} > 0$ and define a simple tame function $\Phi : \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ by

$$\Phi(\phi(0), \phi) = \phi(-\bar{\theta}) \quad \forall (\phi(0), \phi) \in \mathbf{M}_\rho^2.$$

Then it can be shown that $\Phi \in C_{lip}^2(\mathbf{M}_\rho^2)$ and yet $\Phi \notin \mathcal{D}(\Gamma)$.

It will be shown in the proof of Theorem (3.3.6), however, that any tame function of the above form can be approximated by a sequence of quasi-tame functions that are in $\mathcal{D}(\Gamma)$.

Again, consider the associated Markovian $\mathbf{M}_{\rho,+}^2$ -valued process $\{(S(t), S_t), t \geq 0\}$ (where $S_t(\theta) = S(t + \theta), \theta \in (-\infty, 0]$) described by the nonlinear stochastic hereditary differential equation (1) with the initial historical price function $(S(0), S_0) = (\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$. We have the following result for its weak infinitesimal generator $\mathbf{A} + \Gamma$ (see e.g. Arriojas (1997),

Mohammed (1984) and Mohammed (1996)):

Theorem (3.3.3). If $\Phi \in C_{lip}^2(\mathbf{M}_\rho^2) \cap \mathcal{D}(\Gamma)$ and $\{(S(t), S_t), t \geq 0\}$ is the $\mathbf{M}_{\rho,+}^2$ -valued solution process corresponding to (1) with an initial historical price function $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$, then

$$\lim_{t \downarrow 0} \frac{E[\Phi(S(t), S_t) - \Phi(\psi(0), \psi)]}{t} = (\mathbf{A} + \Gamma)\Phi(\psi(0), \psi), \quad (27)$$

where

$$\begin{aligned} \mathbf{A}\Phi(\psi(0), \psi) &= \frac{1}{2} \partial_{\psi(0)}^2 \Phi(\psi(0), \psi) \psi^2(0) g^2(\psi) \\ &+ \partial_{\psi(0)} \Phi(\psi(0), \psi) \psi(0) f(\psi), \end{aligned} \quad (28)$$

and $\Gamma(\Phi)(\psi(0), \psi)$ is as given in (26).

It seems from a glance at (28) that $\mathbf{A}\Phi(\psi(0), \psi)$ requires only the existence of the first and second order partial derivatives $\partial_{\psi(0)} \Phi$ and $\partial_{\psi(0)}^2 \Phi$ of $\Phi(\psi(0), \psi)$ with respect to its first variable $\psi(0) \in \mathfrak{R}$. However, detail derivations of the formula reveal that a stronger condition that $\Phi \in C_{lip}^2(\mathbf{M}_\rho^2)$ is required.

We have the following Dynkin's formula (see Mizel and Trutzer (1984) and Kolmanovskii and Shaikhet (1996)):

Theorem (3.3.4). Let $\Phi \in C_{lip}^2(\mathbf{M}_\rho^2) \cap \mathcal{D}(\Gamma)$. Then

$$\begin{aligned} E[e^{-\delta\tau} \Phi(S(\tau), S_\tau)] &= \Phi(\psi(0), \psi) \\ &+ E \left[\int_0^\tau e^{-\delta t} (\mathbf{A} + \Gamma - \delta I) \Phi(S(t), S_t) dt \right], \end{aligned} \quad (29)$$

for all $P - a.s.$ finite \mathbf{G} -stopping time τ .

The function $\Phi \in C_{lip}^2(\mathbf{M}_\rho^2) \cap \mathcal{D}(\Gamma)$ that has the following special form is referred to as a quasi-tame function

$$\Phi(\phi(0), \phi) = \Psi(m(\phi(0), \phi)), \quad (30)$$

where

$$\begin{aligned} m(\phi(0), \phi) &= \left(\phi(0), \int_{-\infty}^0 \eta_1(\phi(\theta)) \lambda_1(\theta) d\theta, \right. \\ &\left. \dots, \int_{-\infty}^0 \eta_n(\phi(\theta)) \lambda_n(\theta) d\theta \right) \forall (\phi(0), \phi) \in \mathbf{M}_\rho^2, \end{aligned} \quad (31)$$

for some positive integer n and some functions $m \in C(\mathbf{M}_\rho^2; \mathfrak{R}^{n+1})$, $\eta_i \in C^\infty(\mathfrak{R})$, $\lambda_i \in C^1((-\infty, 0])$ with

$$\lim_{\theta \rightarrow -\infty} \lambda_i(\theta) = \lambda_i(-\infty) = 0$$

for $i = 1, 2, \dots, n$, and $\Psi \in C^\infty(\mathfrak{R}^{n+1})$ of the form $\Psi(x, y_1, y_2, \dots, y_n)$.

We have the following Ito's formula in case $\Phi \in \mathbf{M}_\rho^2$ is a quasi-tame function in the sense defined above.

Theorem (3.3.5). Let $\{(S(t), S_t), t \geq 0\}$ be the $\mathbf{M}_{\rho,+}^2$ -valued solution process corresponding to (1) with an initial historical price function $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$. If $\Phi \in C(\mathbf{M}_\rho^2)$ is a quasi-time function, then $\Phi \in \mathcal{D}(\mathbf{A}) \cap \mathcal{D}(\Gamma)$ and

$$\begin{aligned} e^{-\delta\tau} \Phi(S(\tau), S_\tau) &= \Phi(\psi(0), \psi) \\ &+ \int_0^\tau e^{-\delta t} (\mathbf{A} + \Gamma - \delta I) \Phi(S(t), S_t) dt \\ &+ \int_0^\tau e^{-\delta t} \Phi_x(S(t), S_t) S(t) f(S_t) dW(t) \end{aligned} \quad (32)$$

for all finite \mathbf{G} -stopping time τ , where I is the identity operator. Moreover, if $\Phi \in C(\mathbf{M}_\rho^2)$ is the form described in (30)-(31), then

$$\begin{aligned} (\mathbf{A} + \Gamma) \Phi(\psi(0), \psi) &= \sum_{i=1}^n \Psi_{y_i}(m(\psi(0), \psi)) \\ &\times \left(\eta_i(\psi(0)) \lambda_i(0) - \int_{-\infty}^0 \eta_i(\psi(\theta)) \dot{\lambda}_i(\theta) d\theta \right) \\ &+ \Psi_x(m(\psi(0), \psi)) \psi(0) f(\psi) + \frac{1}{2} \Psi_{xx}(m(\psi(0), \psi)) \psi^2(0) g^2(\psi), \end{aligned} \quad (33)$$

where Ψ_x , Ψ_{y_i} and Ψ_{xx} denote the partial derivatives of $\Psi(x, y_1, \dots, y_n)$ with respect to its appropriate variables.

The above Ito's formula also holds for any tame function $\Phi : \mathfrak{R} \times \mathbf{C} \rightarrow \mathfrak{R}$ of the following form

$$\Phi(\phi(0), \phi) = \Psi(m(\phi(0), \phi)) = \Psi(\phi(0), \phi(-\theta_1), \dots, \phi(-\theta_n)) \quad (34)$$

where \mathbf{C} is the space continuous functions $\phi : (-\infty, 0] \rightarrow \mathfrak{R}$ equipped with uniform topology, $0 < \theta_1 < \theta_2 < \dots < \theta_n < \infty$, and $\Psi(x, y_1, \dots, y_n)$ is such that $\Psi \in C^\infty(\mathfrak{R}^{n+1})$.

Theorem (3.3.6). Let $\{(S(t), S_t), t \geq 0\}$ be the $\mathbf{M}_{\rho,+}^2$ -valued solution process corresponding to (1) with an initial historical price function $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$. If $\Phi : \mathfrak{R} \times \mathbf{C} \rightarrow \mathfrak{R}$ is a tame function defined by (34), then $\Phi \in \mathcal{D}(\mathbf{A}) \cap \mathcal{D}(\Gamma)$ and

$$\begin{aligned} & e^{-\delta\tau} \Psi(S(\tau), S(\tau - \theta_1), \dots, S(\tau - \theta_n)) \\ = & \Psi(\psi(0), \psi(-\theta_1), \dots, \psi(-\theta_n)) \\ & + \int_0^\tau e^{-\delta t} (\mathbf{A} - \delta I) \Psi(S(t), S(t - \theta_1), \dots, S(t - \theta_n)) dt \\ & + \int_0^\tau e^{-\delta t} \Psi_x(S(t), S(t - \theta_1), \dots, S(t - \theta_n)) S(t) f(S_t) dW(t) \end{aligned} \quad (35)$$

for all finite \mathbf{G} -stopping time τ , where

$$\begin{aligned} & \mathbf{A} \Psi(\psi(0), \psi(-\theta_1), \dots, \psi(-\theta_n)) \\ = & \Psi_x(\psi(0), \psi(-\theta_1), \dots, \psi(-\theta_n)) \psi(0) f(\psi) \\ & + \frac{1}{2} \Psi_{xx}(\psi(0), \psi(-\theta_1), \dots, \psi(-\theta_n)) \psi^2(0) g^2(\psi), \end{aligned} \quad (36)$$

with Ψ_x and Ψ_{xx} being the first and second order derivatives with respect to x of $\Psi(x, y_1, \dots, y_n)$.

3.4 Dynkin's Formula for the Controlled State Process

Combining the above three subsections, we have the following Dynkin's formula for the controlled (by the admissible strategy π) \mathcal{S}_κ -valued state process $\{Z(t) = (X(t), N_t, S(t), S_t), t \geq 0\}$ with the initial state

$$(X(0-), N_{0-}, S(0), S_0) = (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa :$$

$$\begin{aligned} & E[e^{-\delta\tau} \Phi(Z(\tau))] = \Phi(x, \xi, \psi(0), \psi) \\ + & E\left[\int_0^\tau e^{-\delta t} \mathcal{L}^{C(t)} \Phi(Z(t)) dt\right] \\ + & E\left[\sum_{0 \leq t \leq \tau} e^{-\delta t} (\Phi(Z(t)) - \Phi(Z(t-)))\right], \end{aligned} \quad (37)$$

for all $\Phi : \mathfrak{R} \times \mathbf{N} \times \mathbf{M}_{\rho,+}^2 \rightarrow \mathfrak{R}$ such that $\Phi(\cdot, \xi, \psi(0), \psi) \in C^1(\mathfrak{R})$ for each $(\xi, \psi(0), \psi) \in \mathbf{N} \times \mathbf{M}_\rho^2$ and $\Phi(x, \xi, \cdot, \cdot) \in C_{lip}^2(\mathbf{M}_\rho^2) \cap \mathcal{D}(\Gamma)$ for each $(x, \xi) \in$

$\mathfrak{R} \times \mathbf{N}$, where

$$\mathcal{L}^c \Phi(x, \xi, \psi(0), \psi) = (\mathbf{A} + \Gamma - \delta I + (rx - c)\partial_x)\Phi(x, \xi, \psi(0), \psi), \quad (38)$$

and \mathbf{A} and Γ are as define in (28) and (26).

Note that $E[\dots]$ in the above stands for $E^{x, \xi, \psi(0), \psi; \pi}[\dots]$, the expectation given the initial state $(x, \xi, \psi(0), \psi)$ and $\pi \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$.

In the case, $\Phi \in C(\mathfrak{R} \times \mathbf{N} \times \mathbf{M}_\rho^2)$ is such that $\Phi(x, \xi, \cdot, \cdot) : \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ is a quasi-tame (respectively, tame) function on \mathbf{M}_ρ^2 of the form described in (30)-(31) (respectively (34)), then the following Ito's formula for the controlled state process $\{Z(t) = (X(t), N_t, S(t), S_t), t \geq 0\}$ also holds true.

Theorem (3.4.1). If $\Phi \in C(\mathfrak{R} \times \mathbf{N} \times \mathbf{M}_\rho^2)$ is such that $\Phi(x, \xi, \cdot, \cdot) : \mathbf{M}_\rho^2 \rightarrow \mathfrak{R}$ is a quasi-tame function (respectively, tame) on \mathbf{M}_ρ^2 , then

$$\begin{aligned} e^{-\delta\tau}\Phi(Z(\tau)) &= \Phi(Z(0)) \\ &+ \int_0^\tau e^{-\delta t}\mathcal{L}^{C(t)}\Phi(Z(t))dt \\ &+ \int_0^\tau e^{-\delta t}\partial_{\psi(0)}\Phi(Z(t))S(t)f(S_t)dW(t) \\ &+ \left[\sum_{0 \leq t \leq \tau} e^{-\delta t}(\Phi(Z(t)) - \Phi(Z(t-))) \right], \end{aligned} \quad (39)$$

for every $P - a.s.$ finite \mathbf{G} -stopping time τ .

Moreover, if $\Phi(x, \xi, \psi(0), \psi) = \Psi(x, \xi, m(\psi(0), \psi))$ where $\Psi \in C(\mathfrak{R} \times \mathbf{N} \times \mathfrak{R}^{n+1})$ and $m(\psi(0), \psi)$ is given by (30)-(31) (respectively, (34)) then

$$\mathcal{L}^c \Phi(x, \xi, \psi(0), \psi) = (\mathbf{A} + \Gamma - kI + (rx - c)\partial_x)\Psi(x, \xi, m(\psi(0), \psi))$$

and $(\mathbf{A} + \Gamma)\Psi(x, \xi, m(\psi(0), \psi))$ is as given in (33)(respectively, (36)) for each fixed $(x, \xi) \in \mathfrak{R} \times \mathbf{N}$.

Notation (3.4.2). In the following, we shall use the convention that $C_{lip}^{1,0,2}(\mathcal{O}) \cap \mathcal{D}(\Gamma)$ as the collection of continuous functions $\Phi : \mathcal{O} \rightarrow \mathfrak{R}$ ($\mathcal{S}_\kappa \subset \mathcal{O}$) such that $\Phi(\cdot, \xi, \psi(0), \psi) \in C^1(\mathfrak{R})$ for each $(\xi, \psi(0), \psi)$, and $\Phi(x, \xi, \cdot, \cdot) \in C_{lip}^2(\mathbf{M}_\rho^2) \cap \mathcal{D}(\Gamma)$ for each (x, ξ) .

4 The Quasi-Variational HJB Inequality

The main objective of this section is to present the dynamic programming equation for the value function in the form of an infinite-dimensional quasi

variational Hamilton-Jacobi-Bellman (HJB) inequality (or QVHJBI) (see QVHJBI (*) in §4.3(D)).

4.1 The Dynamic Programming Principle

The following Bellman-type Dynamic Programming Principle (DPP) was established in Shreve and Soner (1994) and still holds true in our problem by combining with that obtained in Larssen (2002) and Larssen and Risebro (2003) for optimal classical control of stochastic functional differential equations with a bounded memory.

Proposition (4.1.1). (Bellman's Dynamic Programming Principle) Let $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ be given and let τ be a \mathbf{G} -stopping time. Then, for every $\pi = (C, \mathcal{T}) \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$, we have $P - a.s.$

$$E \left[\int_\tau^\infty e^{-\delta s} \frac{C(s)^\gamma}{\gamma} ds \mid \mathcal{G}(\tau) \right] \leq \chi_{\{\tau < \infty\}} e^{-\delta \tau} V_\kappa(X(\tau), N_\tau, S(\tau), S_\tau). \quad (40)$$

Moreover, for each $\epsilon > 0$, there is a consumption-investment policy $\pi^\epsilon = (C^\epsilon, \mathcal{T}^\epsilon) \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$ agreeing with π on the random interval $[0, \tau)$ for which

$$\chi_{\{\tau < \infty\}} e^{-\delta \tau} V_\kappa(X(\tau), N_\tau, S(\tau), S_\tau) < \epsilon + E \left[\int_\tau^\infty e^{-\delta s} \frac{C^\gamma(s)}{\gamma} ds \mid \mathcal{G}(\tau) \right], \quad (41)$$

is satisfied $P - a.s.$.

Note that $(X(\tau), N_\tau, S(\tau), S_\tau)$ in (41)-(42) is determined by $\pi = (C, \mathcal{T})$. The construction of $\pi^\epsilon = (C^\epsilon, \mathcal{T}^\epsilon)$ takes $(X(\tau-), N_{\tau-}, S(\tau), S_\tau)$ as the initial state, from which an initial jump may occur. Thus, the controlled state process $(X^\epsilon(\cdot), N^\epsilon, S^\epsilon(\cdot), S^\epsilon)$ associated with $\pi^\epsilon = (C^\epsilon, \mathcal{T}^\epsilon)$ agrees with $\{(X(t), N_t, S(t), S_t), t \in [0, \tau)\}$ and $(X^\epsilon(\tau), N_\tau^\epsilon, S^\epsilon(\tau), S_\tau^\epsilon)$ can be reached from $(X(\tau-), N_{\tau-}, S(\tau), S_\tau)$ by a transaction.

Corollary (4.1.2). Let $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ be given and let \mathcal{O} be an open subset of \mathcal{S}_κ containing $(x, \xi, \psi(0), \psi)$. For $\pi = (C, \mathcal{T}) \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$, let $\{(X(t), N_t, S(t), S_t), t \geq 0\}$ be given by (13)-(14), (16)-(17) and (1). Define

$$\tau = \inf \{t \geq 0 \mid (X(t), N_t, S(t), S_t) \notin \bar{\mathcal{O}}\}.$$

Then, for each $t \in [0, \infty)$, we have the following optimality equation:

$$\begin{aligned} V_\kappa(x, \xi, \psi(0), \psi) &= \sup_{\pi \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)} E \left[\int_0^{t \wedge \tau} e^{-\delta s} \frac{C^\gamma(s)}{\gamma} ds \right. \\ &\quad \left. + \chi_{\{t \wedge \tau < \infty\}} e^{-\delta(t \wedge \tau)} V_\kappa(X(t \wedge \tau), N_{t \wedge \tau}, S(t \wedge \tau), S_{t \wedge \tau}) \right]. \end{aligned} \quad (42)$$

4.2 Heuristic Derivation of the QVHJBI

In this subsection, we shall heuristically derive the Hamilton-Jacobi-Bellman (HJB) quasi-variational inequality (see QVHJBI (*) in §4.3(D)) based on the dynamic programming principle described in Proposition (4.1.1) and Corollary (4.1.2). We emphasize here that it is not our intension to rigorously verify every step involved in the derivations since the rigorous verification are to be done in §5 and §6.

To derive QVHJBI (*) in §4.3(D), we consider the effects on the value function when there is consumption but no transaction and when there is transaction but no consumption.

(A). Consumptions Without Transaction.

Assume first that there is no transaction then the corresponding state process $\{Z(t) = (X(t), N_t, S(t), S_t), t \geq 0\}$ satisfies the following set of equations:

$$dX(t) = [rX(t) - C(t)]dt, \quad t \geq 0; \quad (43)$$

$$\frac{dS(t)}{S(t)} = f(S_t)dt + g(S_t)dW(t), \quad t \geq 0; \quad \text{and} \quad (44)$$

$$N_t = \xi, \quad t \geq 0, \quad (45)$$

with the initial state $(X(0-), N_{0-}, S(0), S_0) = (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$.

In this case, $V_\kappa(X(t), N_t, S(t), S_t) = V_\kappa(X(t-), N_{t-}, S(t), S_t)$ for all $t \geq 0$, since there is no jump transaction. Assuming that the value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ is sufficiently smooth. From Corollary (4.1.2) and (38)-(39), we have

$$\begin{aligned} 0 &\geq \lim_{t \downarrow 0} \frac{E \left[e^{-\delta t} V_\kappa(X(t), N_t, S(t), S_t) - V_\kappa(x, \xi, \psi(0), \psi) \right]}{t} \\ &\quad + \lim_{t \downarrow 0} \frac{1}{t} E \left[\int_0^t e^{-\delta s} \frac{C^\gamma(s)}{\gamma} ds \right] \end{aligned}$$

$$\begin{aligned}
&= \lim_{t \downarrow 0} \frac{E \left[e^{-\delta t} (V_\kappa(X(t), N_t, S(t), S_t) - V_\kappa(x, \xi, \psi(0), \psi)) \right]}{t} \\
&\quad + \lim_{t \downarrow 0} \frac{\left[(e^{-\delta t} - 1) V_\kappa(x, \xi, \psi(0), \psi) \right]}{t} \\
&\quad + \lim_{t \downarrow 0} \frac{1}{t} E \left[\int_0^t e^{-\delta s} \frac{C^\gamma(s)}{\gamma} ds \right] \\
&= \lim_{t \downarrow 0} \frac{E \left[e^{-\delta t} \int_0^t (\mathbf{A} + \Gamma - (rX(t) - C(t)) \partial_x) V_\kappa(X(t), N_t, S(t), S_t) dt \right]}{t} \\
&\quad - \delta V_\kappa(x, \xi, \psi(0), \psi) + \lim_{t \downarrow 0} \frac{1}{t} E \left[\int_0^t e^{-\delta s} \frac{C^\gamma(s)}{\gamma} ds \right] \\
&= \left(\mathbf{A} + \Gamma + (rx - c) \partial_x - \delta I \right) V_\kappa(x, \xi, \psi(0), \psi) + \frac{c^\gamma}{\gamma}, \quad \forall c \geq 0.
\end{aligned}$$

This shows that

$$\begin{aligned}
0 &\geq \mathcal{A}V_\kappa(x, \xi, \psi(0), \psi) \equiv \sup_{c \geq 0} \left(\mathcal{L}^c V_\kappa(x, \xi, \psi(0), \psi) + \frac{c^\gamma}{\gamma} \right) \quad (46) \\
&= \left(\mathbf{A} + \Gamma + rx \partial_x - \delta \right) V_\kappa(x, \xi, \psi(0), \psi) \\
&\quad + \sup_{c \geq 0} \left(\frac{c^\gamma}{\gamma} - c \partial_x V_\kappa(x, \xi, \psi(0), \psi) \right) \\
&= \left(\mathbf{A} + \Gamma + rx \partial_x - \delta \right) V_\kappa(x, \xi, \psi(0), \psi) \\
&\quad + \frac{1 - \gamma}{\gamma} (\partial_x V_\kappa)^{\frac{\gamma}{\gamma-1}}(x, \xi, \psi(0), \psi),
\end{aligned}$$

since the maximum of the the above expression is achieved at

$$c^* = (\partial_x V_\kappa)^{\frac{1}{\gamma-1}}(x, \xi, \psi(0), \psi). \quad (47)$$

Note that the Fréchet differential operator \mathbf{A} and Γ are defined in (28) and (26), respectively.

(B). Transactions Without Consumption.

We next consider the case where there are transactions but no consumption. For each locally bounded $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ and each $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ define the *intervention operator*

$$\begin{aligned}
\mathcal{M}_\kappa \Phi(x, \xi, \psi(0), \psi) &= \sup \{ \Phi(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}) \mid \zeta \in \mathcal{R}(\xi) - \{\mathbf{0}\}, \quad (48) \\
&\quad (\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}) \in \mathcal{S}_\kappa \},
\end{aligned}$$

where $(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi})$ are as defined in (5)-(7). If $(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}) \notin \mathcal{S}_\kappa$ for all $\zeta \in \mathcal{R}(\xi) - \{\mathbf{0}\}$, we set $\mathcal{M}_\kappa \Phi(x, \xi, \psi(0), \psi) = 0$. If for all $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ there exists $(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}) \in \mathcal{S}_\kappa$ such that

$$\mathcal{M}_\kappa \Phi(x, \xi, \psi(0), \psi) = \Phi(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}),$$

then we set

$$\hat{\zeta}(x, \xi, \psi(0), \psi) = \hat{\zeta}_\Phi(x, \xi, \psi(0), \psi) = (\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}) \in \mathcal{R}(\xi). \quad (49)$$

Note we let $\hat{\zeta}(x, \xi, \psi(0), \psi)$ denote a measurable selection of the map

$$(x, \xi, \psi(0), \psi) \mapsto (\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi})$$

defined in (49).

We make the following technical assumption regarding the existence of a measurable selection

$$\hat{\zeta}(x, \xi, \psi(0), \psi) = \hat{\zeta}_{V_\kappa}(x, \xi, \psi(0), \psi)$$

for the value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$, *i.e.*, there exists a measurable function $\hat{\zeta}_{V_\kappa} : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ such that

$$V_\kappa(\hat{\zeta}(x, \xi, \psi(0), \psi)) = \mathcal{M}_\kappa V_\kappa(x, \xi, \psi(0), \psi) \quad \forall (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa. \quad (50)$$

Assumption (4.2.1). For each $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ There exists a measurable function $\hat{\zeta}_{V_\kappa} : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ such that (50) is satisfied for every $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$.

Assume without loss of generality that the *investor's* current portfolio is at $(X(t-), N_{t-}, S(t), S_t) = (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$. An immediate transaction of the amount $\zeta \in \mathcal{R} - \{\mathbf{0}\}$ without consumption (*i.e.*, $C(t) = 0$) yields $(X(t), N_t, S(t), S_t) = (\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi})$, where \hat{x} , $\hat{\xi}$, and $\hat{\psi}(0)$, $\hat{\psi}$ are as given in (5)-(7). It is clear that

$$V_\kappa(x, \xi, \psi(0), \psi) \geq \mathcal{M}_\kappa V_\kappa(x, \xi, \psi(0), \psi) \quad \forall (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa, \quad (51)$$

due to the following lemma.

Lemma (4.2.2). $V_\kappa(x, \xi, \psi(0), \psi) \geq \mathcal{M}_\kappa V_\kappa(x, \xi, \psi(0), \psi) \quad \forall (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$.

Combining §4.2(A) and §4.2(B), we have the following inequality:

$$\max \left\{ \mathcal{A}V_\kappa, \mathcal{M}_\kappa V_\kappa - V_\kappa \right\} \leq 0 \text{ on } \mathcal{S}_\kappa^\circ,$$

where \mathcal{S}_κ° denotes the interior of the *solvency region* \mathcal{S}_κ .

Using a standard technique in deriving the variational HJB inequality for stochastic classical-singular and classical-impulse control problems (see Bensoussan and Lions (1984) for stochastic impulse controls, Brekke and Oksendal (1998) and Oksendal and Sulem (2002) for stochastic classical-impulse controls, and Larssen (2002) and Larssen and Risero (2003) for classical and singular controls of stochastic delay equations), one can show that on the set

$$\{(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa^\circ \mid \mathcal{M}_\kappa V_\kappa(x, \xi, \psi(0), \psi) < V_\kappa(x, \xi, \psi(0), \psi)\}$$

we have $\mathcal{A}V_\kappa = 0$ and on the set

$$\{(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa^\circ \mid \mathcal{A}V_\kappa(x, \xi, \psi(0), \psi) < 0\}$$

we have $\mathcal{M}_\kappa V_\kappa = V_\kappa$. Therefore, we have the following QVHJBI on \mathcal{S}_κ° :

$$\max \left\{ \mathcal{A}V_\kappa, \mathcal{M}_\kappa V_\kappa - V_\kappa \right\} = 0 \text{ on } \mathcal{S}_\kappa^\circ, \quad (52)$$

where

$$\mathcal{A}\Phi = (\mathbf{A} + \Gamma + rx\partial_x - \delta)\Phi + \sup_{c \geq 0} \left(\frac{c^\gamma}{\gamma} - c\partial_x\Phi \right), \quad (53)$$

and $\mathcal{M}_\kappa\Phi$ is as given in (48).

The boundary values for the QVHJBI on $\partial\mathcal{S}_\kappa$ are given in the next subsection.

4.3 Boundary Values of the QVHJBI

(A). The Solvency Region and The Value Function for $\kappa = 0$ and $\mu > 0$.

Due to the drastic differences in their characteristics between the stochastic classical-singular control and the stochastic classical-impulse control problems, the hereditary portfolio optimization problem without a fixed transaction cost (*i.e.*, $\kappa = 0$ and $\mu > 0$) shall be treated in a separate paper.

However, we make the following observations for this case:

Remark (4.3.1). 1. When there is no fixed transaction cost (*i.e.*, $\kappa = 0$ and $\mu > 0$), the *solvency region* \mathcal{S}_0 becomes

$$\begin{aligned}\mathcal{S}_0 &= \{(x, \xi, \psi(0), \psi) \mid H_0(x, \xi, \psi(0), \psi) \geq 0\} \\ &= \{(x, \xi, \psi(0), \psi) \mid G_0(x, \xi, \psi(0), \psi) \geq 0\}\end{aligned}$$

due to the fact that

$$x \geq 0 \text{ and } n(-i) \geq 0 \forall i = 0, 1, 2, \dots \Rightarrow G_0(x, \xi, \psi(0), \psi) \geq 0.$$

Hence

$$\mathfrak{R}_+ \times \mathbf{N}_+ \times \mathbf{M}_{\rho,+}^2 \subset \{(x, \xi, \psi(0), \psi) \mid G_0(x, \xi, \psi(0), \psi) \geq 0\}.$$

In this case all shares of the *stock* owned or owed can be liquidated due the absence of a fixed transaction cost $\kappa = 0$.

2. For each $(\psi(0), \psi) \in \mathbf{M}_{\rho,+}^2$, let $\mathcal{S}_0(\psi(0), \psi)$ be the projection of the *solvency region* along $(\psi(0), \psi)$ defined by

$$\begin{aligned}\mathcal{S}_0(\psi(0), \psi) &\equiv \{(x, \xi) \in \mathfrak{R} \times \mathbf{N} \mid H_0(x, \xi, \psi(0), \psi) \geq 0\} \\ &= \{(x, \xi) \in \mathfrak{R} \times \mathbf{N} \mid G_0(x, \xi, \psi(0), \psi) \geq 0\}.\end{aligned}$$

The set $\mathcal{S}_0(\psi(0), \psi)$ is a convex subset of $\mathfrak{R} \times \mathbf{N}$ in the sense that

$$(x, \xi) \in \mathcal{S}_0(\psi(0), \psi) \Rightarrow \alpha(x, \xi) = (\alpha x, \alpha \xi) \in \mathcal{S}_0(\psi(0), \psi) \forall \alpha \geq 0.$$

3. The value function $V_0 : \mathcal{S}_0 \rightarrow \mathfrak{R}_+$ for the case $\kappa = 0$ and $\mu > 0$ has the following concavity property:

For each fixed $(\psi(0), \psi) \in \mathbf{M}_{\rho}^2$, $V_0(\cdot, \cdot, \psi(0), \psi) : \mathcal{S}_0(\psi(0), \psi) \rightarrow \mathfrak{R}_+$ is a concave function, *i.e.*, if $(x_1, \xi_1), (x_2, \xi_2) \in \mathcal{S}_0(\psi(0), \psi)$ and $0 \leq \lambda \leq 1$, then

$$\begin{aligned}&V_0(\lambda x_1 + (1 - \lambda)x_2, \lambda \xi_1 + (1 - \lambda)\xi_2, \psi(0), \psi) \\ &\geq \lambda V_0(x_1, \xi_1, \psi(0), \psi) + (1 - \lambda)V_0(x_2, \xi_2, \psi(0), \psi).\end{aligned}$$

The detail proof of this statement shall be provided in Chang (2004b).

(B). Decomposition of $\partial\mathcal{S}_\kappa$.

Note that $\partial\mathcal{S}_\kappa$ can be decomposed as follows.

Let $I = I(x, \xi, \psi(0), \psi) \subset \mathfrak{N} \equiv \{0, 1, 2, \dots\}$ be defined as

$$I(x, \xi, \psi(0), \psi) = \{i \in \mathfrak{N} \mid (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa \text{ and } n(-i) < 0\}.$$

Therefore, $\{\tau(-i), i \in I\}$ consists of those time instances at which the *investor* short-sold and $\{\tau(-i), i \notin I\}$ consists of that at which the *investor* purchased shares of the *stock*.

Note that the index set $I \subset \mathbb{N}$ defined above is a function of $(x, \xi, \psi(0), \psi)$. However, with an abuse of notation, we shall also interpret I as the collection of those states $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ with $n(-i) < 0$ for all $i \in I$ and $n(-i) \geq 0$ for all $i \notin I$, *i.e.*,

$$I = \{(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa \mid n(-i) < 0 \text{ for all } i \in I \\ \text{and } n(-i) \geq 0 \text{ for all } i \notin I\}.$$

With this interpretation,

$$\partial \mathcal{S}_\kappa = \bigcup_{I \subset \mathbb{N}} (\partial_{-,I} \mathcal{S}_\kappa \cup \partial_{+,I} \mathcal{S}_\kappa), \quad (54)$$

where

$$\partial_{-,I} \mathcal{S}_\kappa = \partial_{-,I,1} \mathcal{S}_\kappa \cup \partial_{-,I,2} \mathcal{S}_\kappa, \quad (55)$$

$$\partial_{+,I} \mathcal{S}_\kappa = \partial_{+,I,1} \mathcal{S}_\kappa \cup \partial_{+,I,2} \mathcal{S}_\kappa, \quad (56)$$

$$\begin{aligned} \partial_{+,I,1} \mathcal{S}_\kappa &= \{(x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) = 0, x \geq 0, \\ &n(-i) < 0 \text{ for all } i \in I \ \& \ n(-i) \geq 0 \text{ for all } i \notin I\}, \end{aligned} \quad (57)$$

$$\begin{aligned} \partial_{+,I,2} \mathcal{S}_\kappa &= \{(x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) < 0, x \geq 0, \\ &n(-i) = 0 \text{ for all } i \in I \ \& \ n(-i) \geq 0 \text{ for all } i \notin I\}, \end{aligned} \quad (58)$$

$$\begin{aligned} \partial_{-,I,1} \mathcal{S}_\kappa &= \{(x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) = 0, x < 0, \\ &n(-i) < 0 \text{ for all } i \in I \ \& \ n(-i) \geq 0 \text{ for all } i \notin I\}, \end{aligned} \quad (59)$$

and

$$\begin{aligned} \partial_{-,I,2} \mathcal{S}_\kappa &= \{(x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) < 0, x = 0, \\ &n(-i) = 0 \text{ for all } i \in I \ \& \ n(-i) \geq 0 \text{ for all } i \notin I\}. \end{aligned} \quad (60)$$

The interface (intersection) between $\partial_{+,I,1} \mathcal{S}_\kappa$ and $\partial_{+,I,2} \mathcal{S}_\kappa$ is denoted by

$$\begin{aligned} Q_{+,I} &= \{(x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) = 0, x \geq 0, \\ &n(-i) = 0 \text{ for all } i \in I \ \& \ n(-i) \geq 0 \text{ for all } i \notin I\}. \end{aligned} \quad (61)$$

Whereas the interface between $\partial_{-,I,1}\mathcal{S}_\kappa$ and $\partial_{-,I,2}\mathcal{S}_\kappa$ is denoted by

$$\begin{aligned} Q_{-,I} = \{ & (0, \xi, \psi(0), \psi) \mid G_\kappa(0, \xi, \psi(0), \psi) = 0, x = 0, \\ & n(-i) = 0 \text{ for all } i \in I \ \& \ n(-1) \geq 0 \text{ for all } i \notin I\}. \end{aligned} \quad (62)$$

For example, if $I = \mathbb{N}$, then $n(-i) < 0 \ \forall i = 0, 1, 2, \dots$ and

$$G_\kappa(x, \xi, \psi(0), \psi) \geq 0 \Rightarrow x \geq \kappa.$$

In this case, $\partial_{-, \mathbb{N}}\mathcal{S}_\kappa = \emptyset$ (the empty set),

$$\begin{aligned} \partial_{+, \mathbb{N}, 1}\mathcal{S}_\kappa = \{ & (x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) = 0, x \geq 0, \\ & \& \ n(-i) < 0 \text{ for all } i \in \mathbb{N}\}, \end{aligned}$$

and

$$\begin{aligned} \partial_{+, \mathbb{N}, 2}\mathcal{S}_\kappa &= \{(x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) < 0, x \geq 0, \\ & \& \ n(-i) = 0 \text{ for all } i \in \mathbb{N}\} \\ &= \{(x, \mathbf{0}, \psi(0), \psi) \mid 0 \leq x \leq \kappa\}. \end{aligned}$$

On the other hand, if $I = \emptyset$ (the empty set), *i.e.*, $n(-i) \geq 0$ for all $i \in \mathbb{N}$, then

$$\begin{aligned} \partial_{+, \emptyset, 1}\mathcal{S}_\kappa = \{ & (x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) = 0, x \geq 0, \\ & \& \ n(-i) \geq 0 \text{ for all } i \in \mathbb{N}\}, \end{aligned}$$

$$\begin{aligned} \partial_{+, \emptyset, 2}\mathcal{S}_\kappa = \{ & (x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) < 0, x \geq 0, \\ & \& \ n(-i) \geq 0 \text{ for all } i \in \mathbb{N}\}, \end{aligned}$$

$$\begin{aligned} \partial_{-, \emptyset, 1}\mathcal{S}_\kappa = \{ & (x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) = 0, x < 0, \\ & \& \ n(-i) \geq 0 \text{ for all } i \in \mathbb{N}\}, \end{aligned}$$

and

$$\begin{aligned} \partial_{-, \emptyset, 2}\mathcal{S}_\kappa = \{ & (x, \xi, \psi(0), \psi) \mid G_\kappa(x, \xi, \psi(0), \psi) < 0, x = 0, \\ & \& \ n(-i) \geq 0 \text{ for all } i \in \mathbb{N}\}. \end{aligned}$$

(C). Boundary Conditions for The Value Function.

Let us now examine the behavior of the value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ on the

boundary $\partial\mathcal{S}_\kappa$ of the *solvency region* \mathcal{S}_κ defined in (54)-(60).

We make the following observations regarding the behavior of the value function V_κ on the boundary $\partial\mathcal{S}_\kappa$.

Lemma (4.3.2). Let $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$, and let \hat{x} , $\hat{\xi}$, and $(\hat{\psi}(0), \hat{\psi})$ be as defined in (5)-(7). Then

$$G_0(\hat{x}, \hat{\xi}, \hat{\psi}(0), \hat{\psi}) = G_0(x, \xi, \psi(0), \psi) - \kappa. \quad (63)$$

Lemma (4.3.3). If there is no fixed transaction cost (*i.e.*, $\kappa = 0$ and $\nu > 0$) and if $(x, \xi, \psi(0), \psi) \in \partial_{I,1}\mathcal{S}_0$, *i.e.*,

$$G_0(x, \xi, \psi(0), \psi) = 0,$$

then the only admissible strategy is to do no consumption but buy back $n(-i)$ shares for $i \in I$ and sell $n(-i)$ shares for $i \in I^c$ of the stock in order to bring his portfolio to $\{0\} \times \{\mathbf{0}\} \times \mathbf{M}_{\rho,+}^2$ after paying proportional transaction costs and capital-gain taxes, *etc.* In other words, bring his portfolio from the position $(x, \xi, \psi(0), \psi) \in \partial_{I,1}\mathcal{S}_0$ to $(0, \mathbf{0}, \psi(0), \psi)$ by the quantity that satisfy the following equations:

$$\begin{aligned} 0 &= x + \sum_{i \in I^c} [n(-i)\psi(0)(1 - \mu - \beta) + \beta n(-i)\psi(\tau(-i))] \\ &+ \sum_{i \in I} [n(-i)\psi(0)(1 + \mu - \beta) - \beta n(-i)\psi(\tau(-i))]; \text{ and} \\ \mathbf{0} &= \xi \oplus \zeta. \end{aligned}$$

We have the following result.

Theorem (4.3.4). Let $\kappa > 0$ and $\mu > 0$. On $\partial_{I,1}\mathcal{S}_\kappa$ for $I \subset \mathbb{N}$, then the *investor* should not consume but buy back $n(-i)$ shares for $i \in I$ and sell $n(-i)$ shares for $i \in I^c$ of the stock in order to bring his portfolio to $\{0\} \times \{\mathbf{0}\} \times \mathbf{M}_{\rho,+}^2$ after paying transaction costs (fixed plus proportional) and capital-gain taxes and, *etc.* In other words, bring his portfolio from the position $(x, \xi, \psi(0), \psi) \in \partial_{I,1}\mathcal{S}_\kappa$ to $(0, \mathbf{0}, \psi(0), \psi)$ by the quantity that satisfy the following equations:

$$0 = x - \kappa + \sum_{i \in I^c} [n(-i)\psi(0)(1 - \mu - \beta) + \beta n(-i)\psi(\tau(-i))] \quad (64)$$

$$\begin{aligned}
& + \sum_{i \in I} [n(-i)\psi(0)(1 + \mu - \beta) + \beta n(-i)\psi(\tau(-i))]; \\
\mathbf{0} & = \xi \oplus \zeta.
\end{aligned} \tag{65}$$

In this case, the value function $V_\kappa : \partial_{I,1}\mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ satisfies the following equation:

$$(\mathcal{M}_\kappa \Phi - \Phi)(x, \xi, \psi(0), \psi) = 0. \tag{66}$$

We conclude from some simple observations and Theorem (4.3.4) that **Boundary Condition (i)**. On the hyper-plane

$$\partial_{-,0,2}\mathcal{S}_\kappa = \{(0, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa \mid G_\kappa(0, \xi, \psi(0), \psi) < 0, n(-i) \geq 0 \forall i\},$$

the only strategy for the *investor* is to do no transaction and no consumption, since $x = 0$ and $G_\kappa(0, \xi, \psi(0), \psi) < 0$ (hence there is no money to consume and not enough money to pay for the transaction costs, etc.), but to let the stock prices to grow according to (1). Thus, the value function V_κ on $\partial_{-,0,2}\mathcal{S}_\kappa$ satisfies the equation

$$\mathcal{L}^0 \Phi \equiv (\mathbf{A} + \Gamma - \delta + rx\partial_x)\Phi = 0 \tag{67}$$

provided that it is smooth enough;

Boundary Condition (ii). On $\partial_{I,1}\mathcal{S}_\kappa$ for $I \subset \mathfrak{N}$, then the *investor* should not consume but buy back $n(-i)$ shares for $i \in I$ and sell $n(-i)$ shares for $i \in I^c$ of the stock in order to bring his portfolio to $\{0\} \times \{\mathbf{0}\} \times \mathbf{M}_{\rho,+}^2$ after paying transaction costs and capital-gains taxes and etc. In other words, bring his portfolio from the position $(x, \xi, \psi(0), \psi) \in \partial_{I,1}\mathcal{S}_\kappa$ to $(0, \mathbf{0}, \psi(0), \psi)$ by the quantity that satisfy the following equations:

$$\begin{aligned}
0 & = x - \kappa + \sum_{i \in I^c} [n(-i)\psi(0)(1 - \mu - \beta) + \beta n(-i)\psi(\tau(-i))] \\
& + \sum_{i \in I} [n(-i)\psi(0)(1 + \mu + \beta) - \beta n(-i)\psi(\tau(-i))]; \\
\mathbf{0} & = \xi \oplus \zeta.
\end{aligned}$$

In this case, the value function $V_\kappa : \partial_{I,1}\mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ satisfies the following equation:

$$(\mathcal{M}_\kappa \Phi - \Phi)(x, \xi, \psi(0), \psi) = 0. \tag{68}$$

Note that this is a re-statement of Theorem (4.3.4).

Boundary Condition (iii). On $\partial_{+,I,2}\mathcal{S}_\kappa$ for $I \subset \mathbb{N}$, the only optimal strategy is to make no transaction but to consume optimally according to the optimal consumption rate function $c^*(x, \xi, \psi(0), \psi) = \left(\frac{\partial V_\kappa}{\partial x}\right)^{\frac{1}{\gamma-1}}(x, \xi, \psi(0), \psi)$ which is obtained via

$$c^*(x, \xi, \psi(0), \psi) = \arg \max_{c \geq 0} \left\{ \mathcal{L}^c V_\kappa(x, \xi, \psi(0), \psi) + \frac{c^\gamma}{\gamma} \right\},$$

where \mathcal{L}^c is the Frechet partial differential operator defined by

$$\mathcal{L}^c \Phi(x, \xi, \psi(0), \psi) \equiv (\mathbf{A} + \Gamma - \delta)\Phi + (rx - c)\partial_x \Phi. \quad (69)$$

This is because the cash in his *savings* account is not sufficient to buy back any shares of the *stock* but to consume optimally. In this case, the value function $V_\kappa : \partial_{+,I,2}\mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ satisfies the following equation provided that it is smooth enough.

$$\mathcal{A}\Phi \equiv (\mathbf{A} + \Gamma - \delta)\Phi + rx\partial_x \Phi + \frac{1-\gamma}{\gamma} \left(\partial_x \Phi\right)^{\frac{\gamma}{\gamma-1}} = 0. \quad (70)$$

Boundary Condition (iv). On $\partial_{-,I,2}\mathcal{S}_\kappa$, the only admissible consumption-investment strategy is to do no consumption and no transaction but to let the stock price grows as in the Boundary Condition (i).

Boundary Condition (v). On $\partial_{+,\mathbb{N},2}\mathcal{S}_\kappa = \{(x, \xi, \psi(0), \psi) \mid 0 \leq x \leq \kappa, n(-i) = 0 \forall i = 0, 1, \dots\}$, the only admissible consumption-investment strategy is to do no transaction but to consume optimally like in Boundary Condition (iii).

Remark (4.3.5). From Boundary Conditions (i)-(v), it is clear that the value function V_κ is discontinuous on the interfaces $Q_{+,I}$ and $Q_{-,I}$ for all $I \subset \mathbb{N}$.

(D). The QVHJBI With Boundary Conditions.

We conclude from the above subsections that the QVHJBI (together with the boundary conditions) should be expressed as follows.

$$QVHJBI(*) = \begin{cases} \max \left\{ \mathcal{A}\Phi, \mathcal{M}_\kappa \Phi - \Phi \right\} = 0 & \text{on } \mathcal{S}_\kappa^\circ; \\ \mathcal{A}\Phi = 0, & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{+,I,2}\mathcal{S}_\kappa; \\ \mathcal{L}^0 \Phi = 0, & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{-,I,2}\mathcal{S}_\kappa; \\ \mathcal{M}_\kappa \Phi - \Phi = 0 & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{I,1}\mathcal{S}_\kappa. \end{cases}$$

where $\mathcal{A}\Phi$, $\mathcal{L}^0\Phi$ ($\mathcal{L}^c\Phi$ with $c = 0$), and \mathcal{M}_κ are as defined in (82), (81) and (48).

5 The Verification Theorem

Let

$$\tilde{\mathcal{A}}\Phi = \begin{cases} \mathcal{A}\Phi & \text{on } \mathcal{S}_\kappa^\circ \cup \bigcup_{I \subset \mathbb{N}} \partial_{+,I,2}\mathcal{S}_\kappa; \\ \mathcal{L}^0\Phi & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{-,I,2}\mathcal{S}_\kappa. \end{cases}$$

We have the following verification theorem for the value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ for our hereditary portfolio optimization problem:

Theorem (5.1). (The Verification Theorem) (a). Let $U_\kappa = \mathcal{S}_\kappa - \bigcup_{I \subset \mathbb{N}} \partial_{I,1}\mathcal{S}_\kappa$. Suppose, there exists a locally bounded non-negative valued function $\Phi \in C_{lip}^{1,0,2}(\mathcal{S}_\kappa) \cap \mathcal{D}(\Gamma)$ (see Notation (3.4.2)) such that

$$\tilde{\mathcal{A}}\Phi \leq 0 \text{ on } U_\kappa; \text{ and} \quad (71)$$

$$\Phi \geq \mathcal{M}_\kappa\Phi \text{ on } U_\kappa. \quad (72)$$

Then $\Phi \geq V_\kappa$ on \mathcal{U}_κ .

(b). Define $D \equiv \{(x, \xi, \psi(0), \psi) \in U_\kappa \mid \Phi(x, \xi, \psi(0), \psi) > \mathcal{M}_\kappa\Phi(x, \xi, \psi(0), \psi)\}$. Suppose

$$\tilde{\mathcal{A}}\Phi(x, \xi, \psi(0), \psi) = 0 \text{ on } D \quad (73)$$

and that $\hat{\zeta}(x, \xi, \psi(0), \psi) = \hat{\zeta}_\Phi(x, \xi, \psi(0), \psi)$ exists for all $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ by Assumption (4.2.1). Define

$$c^* = \begin{cases} (\partial_x \Phi)^{\frac{1}{\gamma-1}} & \text{on } \mathcal{S}_\kappa^\circ \cup \bigcup_{I \subset \mathbb{N}} \partial_{+,I,2}\mathcal{S}_\kappa; \\ 0 & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{-,I,2}\mathcal{S}_\kappa, \end{cases} \quad (74)$$

and define the impulse control $\mathcal{T}^* = \{(\tau^*(i), \zeta^*(i)), i = 1, 2, \dots\}$ inductively as follows.

First put $\tau^*(0) = 0$ and inductively

$$\tau^*(i+1) = \inf\{t > \tau^*(i) \mid (X^{(i)}(t), N_t^{(i)}, S(t), S_t) \notin D\}, \quad (75)$$

$$\zeta^*(i+1) = \hat{\zeta}(X^{(i)}(\tau^*(i+1)-), N_{\tau^*(i+1)-}^{(i)}, S(\tau^*(i+1)), S_{\tau^*(i+1)}), \quad (76)$$

where $\hat{\zeta}$ is as defined in Assumption (4.2.1) and $\{(X^{(i)}(t), N_t^{(i)}, S(t), S_t), t \geq 0\}$ is the controlled state process obtained by applying the combined control

$$\pi^*(i) = (c^*, (\tau^*(1), \tau^*(2), \dots, \tau^*(i)); \zeta^*(1), \zeta^*(2), \dots, \zeta^*(i)), \quad i = 1, 2, \dots.$$

Suppose $\pi^* = (C^*, \mathcal{T}^*) \in \mathcal{U}_\kappa(x, \xi, \psi(0), \psi)$ and that

$$e^{-\delta t} \Phi(X^*(t), N_t^*, S(t), S_t) \rightarrow 0, \quad \text{as } t \rightarrow \infty \text{ a.s.}$$

and that the family

$$\{e^{-\delta \tau} \Phi(X^*(\tau), N_\tau^*, S(\tau), S_\tau) \mid \tau \text{ } \mathbf{G}\text{-stopping times}\} \quad (77)$$

is uniformly integrable. Then $\Phi(x, \xi, \psi(0), \psi) = V_\kappa(x, \xi, \psi(0), \psi)$ and π^* obtained in (86)-(88) is optimal.

6 The Viscosity Solution

It is clear that the value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ has discontinuity on the interfaces $Q_{I,+}$ and $Q_{I,-}$ (see Remark (4.3.5)) and hence it can not be a solution of QVHJBI (*) in the classical sense. The main purpose of this section is to show that it is a viscosity solution of the QVHJBI (*). See Ishii (1993) for connection of viscosity solutions of second order elliptic equations with stochastic classical control problems.

To give a definition of a viscosity solution, we first define the upper and lower semi-continuity concept as follows.

Let Ξ be a metric space, and let $\Phi : \Xi \rightarrow \mathfrak{R}$ be a Borel measurable function. Then the upper semi-continuous (*USC*) envelop $\bar{\Phi} : \Xi \rightarrow \mathfrak{R}$ and the lower semi-continuous (*LSC*) envelop $\underline{\Phi} : \Xi \rightarrow \mathfrak{R}$ of Φ are defined, respectively, by

$$\bar{\Phi}(\mathbf{x}) = \limsup_{\mathbf{y} \rightarrow \mathbf{x}, \mathbf{y} \in \Xi} \Phi(\mathbf{y}) \quad \text{and} \quad \underline{\Phi}(\mathbf{x}) = \liminf_{\mathbf{y} \rightarrow \mathbf{x}, \mathbf{y} \in \Xi} \Phi(\mathbf{y}).$$

We let $USC(\Xi)$ and $LSC(\Xi)$ denote the set of USC functions and LSC functions on Ξ , respectively.

Note that in general one has

$$\underline{\Phi} \leq \Phi \leq \bar{\Phi},$$

and that Φ is *USC* if and only if $\Phi = \bar{\Phi}$, Φ is *LSC* if and only if $\Phi = \underline{\Phi}$. In particular, Φ is continuous if and only if

$$\underline{\Phi} = \Phi = \bar{\Phi}.$$

Let $\mathcal{L}(\mathbf{M}_\rho^2)$ be the space of bounded linear operators from \mathbf{M}_ρ^2 to \mathfrak{R} equipped with the usual operator norm.

To define a viscosity solution, let us consider the following equation:

$$F(\mathbf{A}, \Gamma, \partial_x, V_\kappa, (x, \xi, \psi(0), \psi)) = 0 \quad \forall (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa, \quad (78)$$

where

$$F : (\mathbf{M}_\rho^2)^\dagger \times \mathcal{L}(\mathbf{M}_\rho^2) \times \mathfrak{R} \times \mathfrak{R}^{\mathcal{S}_\kappa} \times \mathcal{S}_\kappa \rightarrow \mathfrak{R}$$

is defined by

$$F = \begin{cases} \max \left\{ \Lambda(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)), \right. \\ \quad \left. (\mathcal{M}_\kappa \Phi - \Phi)(x, \xi, \psi(0), \psi) \right\}, & \text{on } \mathcal{S}_\kappa^\circ, \\ \Lambda(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)), & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{+, I, 2} \mathcal{S}_\kappa, \\ \Lambda^0(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)), & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{-, I, 2} \mathcal{S}_\kappa, \\ \left. (\mathcal{M} \Phi - \Phi)((x, \xi, \psi(0), \psi)), \right. & \text{on } \bigcup_{I \subset \mathbb{N}} \partial_{I, 1} \mathcal{S}_\kappa \end{cases} \quad (79)$$

where

$$\Lambda(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)) = \mathcal{A} \Phi(x, \xi, \psi(0), \psi),$$

and

$$\Lambda^0(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)) = \mathcal{L}^0 \Phi(x, \xi, \psi(0), \psi).$$

Note that

$$F(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)) = QVHJBI(*),$$

and

$$\begin{aligned} & \bar{F}(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)) \\ &= \max \left\{ \Lambda(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)), \right. \\ & \quad \left. (\mathcal{M}_\kappa \Phi - \Phi)(x, \xi, \psi(0), \psi) \right\} \quad \forall (x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa \end{aligned} \quad (80)$$

and that

$$\underline{F}(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)) = F(\mathbf{A}, \Gamma, \partial_x, \Phi, (x, \xi, \psi(0), \psi)).$$

Definition (6.1). (i) A function $\Phi \in USC(\mathcal{S}_\kappa)$ is said to be a viscosity sub-solution of (100) if for every function $\Psi \in C_{lip}^{1,0,2}(\mathcal{S}_\kappa) \cap \mathcal{D}(\Gamma)$ and for

every $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ such that $\Psi \geq \Phi$ on \mathcal{S}_κ and $\Psi((x, \xi, \psi(0), \psi)) = \Phi(x, \xi, \psi(0), \psi)$ we have

$$\bar{F}(\mathbf{A}, \Gamma, \partial_x, \Psi, (x, \xi, \psi(0), \psi)) \geq 0. \quad (81)$$

(ii) A function $\Phi \in LSC(\mathcal{S}_\kappa)$ is a viscosity super-solution of (100) if for every function $\Psi \in C_{lip}^{1,0,2}(\mathcal{S}_\kappa) \cap \mathcal{D}(\Gamma)$ and for every $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ such that $\Psi \leq \Phi$ on \mathcal{S}_κ and $\Psi(x, \xi, \psi(0), \psi) = \Phi(x, \xi, \psi(0), \psi)$ we have

$$\underline{F}(\mathbf{A}, \Gamma, \partial_x, \Psi, (x, \xi, \psi(0), \psi)) \leq 0. \quad (82)$$

(iii) A locally bounded function $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ is a viscosity solution of (100) if $\bar{\Phi}$ is viscosity sub-solution and $\underline{\Phi}$ is a viscosity super-solution of (100).

The following properties of the *intervention operator* \mathcal{M}_κ can be established similar to Lemma 3.2., Lemma 3.3. and Corollary 3.4. of Oksendal and Sulem (2002) with some modifications to fit our situation.

Lemma (6.2). The following statements hold true regarding \mathcal{M}_κ defined by (48).

- (i) If $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ is *USC*, then $\mathcal{M}_\kappa \Phi$ is *USC*.
- (ii) If $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ is *continuous*, then $\mathcal{M}_\kappa \Phi$ is *continuous*.
- (iii) Let $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$. Then $\overline{\mathcal{M}_\kappa \Phi} \leq \mathcal{M}_\kappa \bar{\Phi}$.
- (iv) Let $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ be such that $\Phi \geq \mathcal{M}_\kappa \Phi$. Then $\underline{\Phi} \geq \mathcal{M}_\kappa \underline{\Phi}$.
- (v) Suppose $\Phi : \mathcal{S}_\kappa \rightarrow \mathfrak{R}$ is *USC* and $\Phi(x, \xi, \psi(0), \psi) > \mathcal{M}_\kappa \Phi(x, \xi, \psi(0), \psi) + \epsilon$ for some $(x, \xi, \psi(0), \psi) \in \mathcal{S}_\kappa$ and $\epsilon > 0$. Then

$$\Phi(x, \xi, \psi(0), \psi) > \overline{\mathcal{M}_\kappa \Phi}(x, \xi, \psi(0), \psi) + \epsilon.$$

Theorem (6.3). Suppose $\delta > r\gamma$. Then the value function $V_\kappa : \mathcal{S}_\kappa \rightarrow \mathfrak{R}_+$ defined by (23) is a viscosity solution of the QVHJBI (*).

References

- [1] M. Akian, J. L. Menaldi, and A. Sulem (1996), *On an investment-consumption model with transaction costs*, SIAM J. Control & Optimization, **34**, 329-364.
- [2] M. Akian, A. Sulem, and M. I. Taksar (2001), *Dynamic optimization of long term growth rate for a portfolio with transaction costs and logarithmic utility*, Mathematical Finance, **11**, 153-188.
- [3] M. Arriojas (1997), *A Stochastic Calculus for Functional Differential Equations*, Doctoral Dissertation, Department of Mathematics, Southern Illinois University at Carbondale.
- [4] M. Arriojas, Y. Hu, S.-E. Mohammed, and G. Pap (2003), *A delayed Balck and Scholes formula*, a preprint.
- [5] A. Bensoussan and J.-L. Lions (1984), *Impulse Control and Quasi-Variational Inequalities*, Gauthier-Villars, Paris.
- [6] K. A. Brekke and B. Oksendal (1998), *A verification theorem for combined stochastic control and impulse control*, in Stochastic Analysis and Related Topics, Vol. 6, Progress in Probability **42**, pp. 211-220, Birkhauser Boston, Cambridge, MA.
- [7] A. Cadenillas and S. R. Pliska (1999), *Optimal trading of a security when there are taxes and transaction costs*, Finance and Stochastics, **3**, 137-165.
- [8] A. Cadenillas and F. Zapataro (2000), *Classical and impulse stochastic control of the exchange rate using interest rates and reserves*, Mathematical Finance, **10**, 141-156.
- [9] M. H. Chang and R. K. Youree (1999), *The European option with hereditary price structures: Basic theory*, Applied Mathematics and Computation, **102**, 279-296.
- [10] M. H. Chang (2004), *Hereditary Portfolio Optimization with Taxes and Fixed Plus Proportional Transaction Costs II: Uniqueness and Approximations of The Viscosity Solution*, in preparation.
- [11] M. H. Chang (2004), *Hereditary Portfolio Optimization with Taxes and Proportional Transaction Costs*, in preparation.

- [12] B. D. Coleman and V. J. Mizel (1966), *Norms and semigroups in the theory of fading memory*, Arch. Rational Mech. Anal, **23**, 87-123.
- [13] G. M. Constantinides (1983), *Capital Market Equilibrium with Personal Tax*, Econometrica, **51**, 611-636.
- [14] G. M. Constantinides (1984), *Optimal stock trading with personal taxes: implications for prices and the abnormal January returns*, J. Financial Economics, **13**, 65-89.
- [15] R. Dammon and C. Spatt (1996), *The optimal trading and pricing of securities with asymmetric capital gains taxes and transaction costs*, Reviews of Financial Studies, **9**, 921-952.
- [16] M. H. A. Davis and A. Norman (1990), *Portfolio selection with transaction costs*, Math. Operations Research, **15**, 676-713.
- [17] K. Ishii (1993), *Viscosity solutions of nonlinear second order elliptic PDEs associated with impulse control problems*, Funkcial. Ekvac., **36**, 123-141.
- [18] E. Jouini, P.-F. Koehl, and N. Touzi (1999), *Optimal investment with taxes: an optimal control problem with endogeneous delay*, Nonlinear Analysis, Theory, Methods and Applications, **37**, 31-56.
- [19] E. Jouini, P.-F. Koehl, and N. Touzi (2000), *Optimal investment with taxes: an existence result*, J. Mathematical Economics, **33**, 373-388.
- [20] I. Karatzas and S. E. Shreve (1991), *Brownian Motion and Stochastic Calculus*, 2nd edition, Springer-Verlag, New York.
- [21] V. B. Kolmanovskii, and L. E. Shaikhet (1996), *Control of Systems with Aftereffect*, Translations of Mathematical Monographs Vol. 157, American Mathematical Society.
- [22] B. Larssen (2002), *Dynamic programming in stochastic control of systems with delay*, Stochastics & Stochastics Reports, **74**, 651-673.
- [23] B. Larssen and N. H. Risebro (2003), *When are HJB-equations for control problems with stochastic delay equations finite dimensional*, Stochastic Analysis and Applications, **21**, 643-661.
- [24] V. J. Mizel and V. Trutzer (1984), *Stochastic hereditary equations: existence and asymptotic stability*, J. of Integral Equations, **7**, 1-72.

- [25] S.-E. A. Mohammed (1984), *Stochastic Functional Differential Equations*, Research Notes in Mathematics **99**, Pitman Advanced Publishing Program, Boston, London, Melbourne.
- [26] S.-E. A. Mohammed (1996), *Stochastic differential systems with memory: theory, examples, and applications* in L. Decreasefond, J. Gjerde, B. Oksendal, and A. S. Ustunel, editors, *Stochastic Analysis and Related Topics VI, The Geilo Workshop 1996*, Progress in Probability, Birkhauser.
- [27] B. Oksendal (2000), *Stochastic Differential Equations*, 5th edition, Springer-Verlag, Berlin, New York.
- [28] B. Oksendal, & A. Sulem (2002), *Optimal consumption and portfolio with both fixed and proportional transaction costs*, SIAM J. Control & Optimization **40**, 1765-1790.
- [29] P. Protter (1995), *Stochastic Integration and Differential Equations*, Springer-Verlag.
- [30] L. C. G. Rogers and D. Williams (1987), *Diffusions, Markov Processes, and Martingales*, Vol. 2, New York: John Wiley & Sons.
- [31] S. E. Shreve and H. M. Soner (1994), *Optimal investment and consumption with transaction costs*, Ann. Appl. Probab., **4**, 609-692.

US Army
TRADOC Analysis Center (TRAC)
Objective Force Urban Operations
Agent Based Simulation Experiment



LTC Tom Cioppa, Ph.D., U.S. Army

Tom.Cioppa@trac.nps.navy.mil

Major Lloyd P. Brown, U.S. Marine Corps

Lloyd.Brown@trac.nps.navy.mil

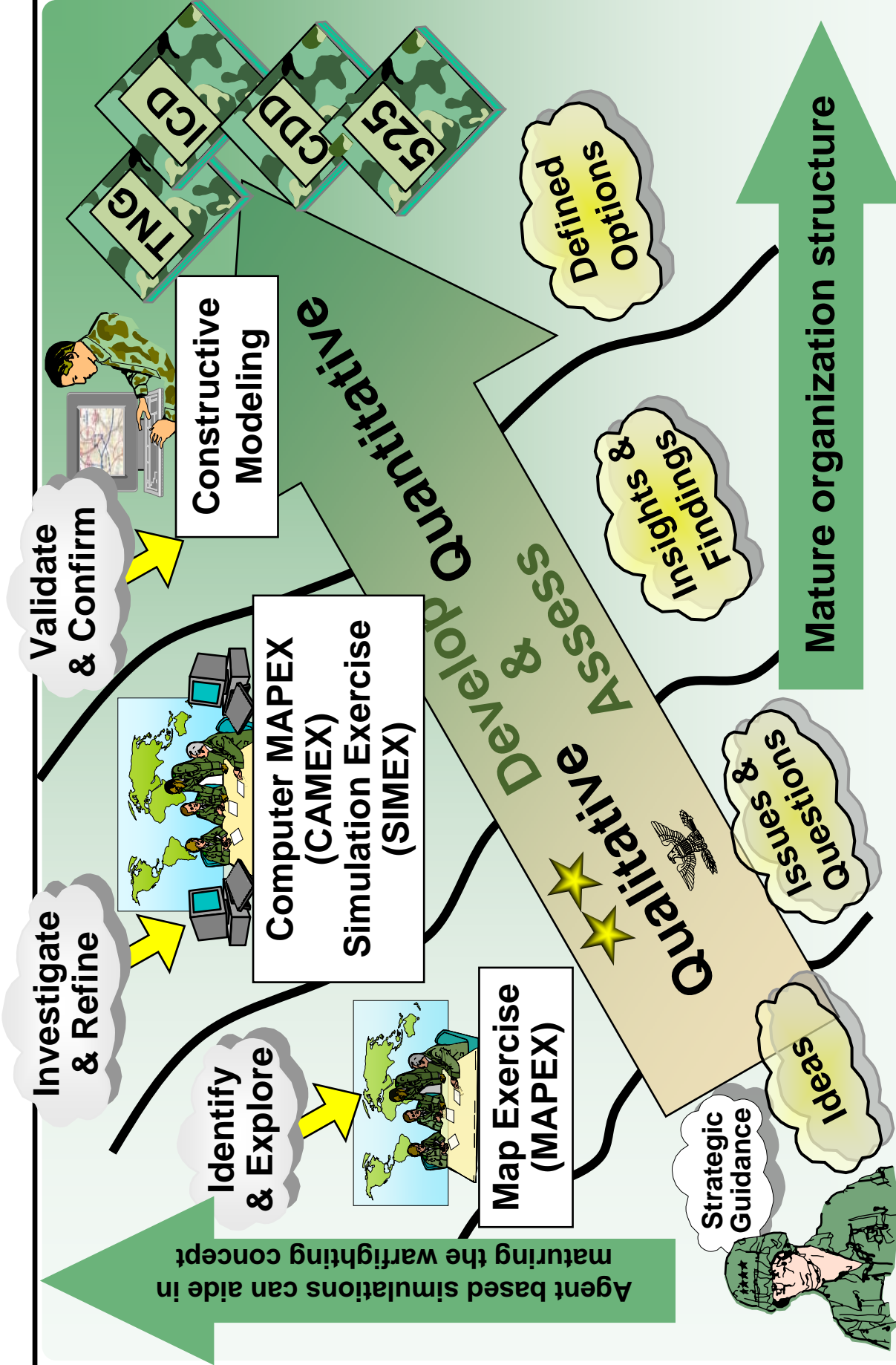
Purpose

Determine the suitability of agent based simulations (ABS) and gain insights into Objective Force small unit operations in an urban environment.

Intermediate goals included:

- Study of ABS (in general)
- Assessment of ABS as exploratory tool and precursor to high-resolution runs

Concepts & Requirements Development



Agent based simulations can aide in maturing the warrighting concept

Participating Organizations

- TRAC-Monterey
- US Army Infantry Center Dismounted Battlespace Battle Lab (USAIC DBBL)
- US Marine Corps Warfighting Lab (Project Albert)
- US Naval Postgraduate School

Agent Based Simulations

Strengths:

- Quick scenario set up time and fast run times
- Non-scripted simulation runs
- Ability to rapidly consider many alternatives

Limitations:

- Some data does not correlate well to real world data
- Computation demands increase rapidly with large entity counts and large battlefield representations

Urban Experiment

- Use a series of new models/analytic tools developed under Marine Corps Warfighting Lab's Project Albert and existing high-resolution simulations.
 - Across Simulations or Distillations
 - MANA
 - Pythagoras
 - Between levels of resolution
 - JANUS
- **Exploit advances in computing power and visualization tools.**
 - Maui High Performance Computer Center (MHPCC)
- **Look at questions from the perspective of conducting exploratory analysis using many data points from a robust design of experiments with replications.**
 - 16 variable / 65 run design and 7 variable / 17 run design

Urban Experiment EEAs

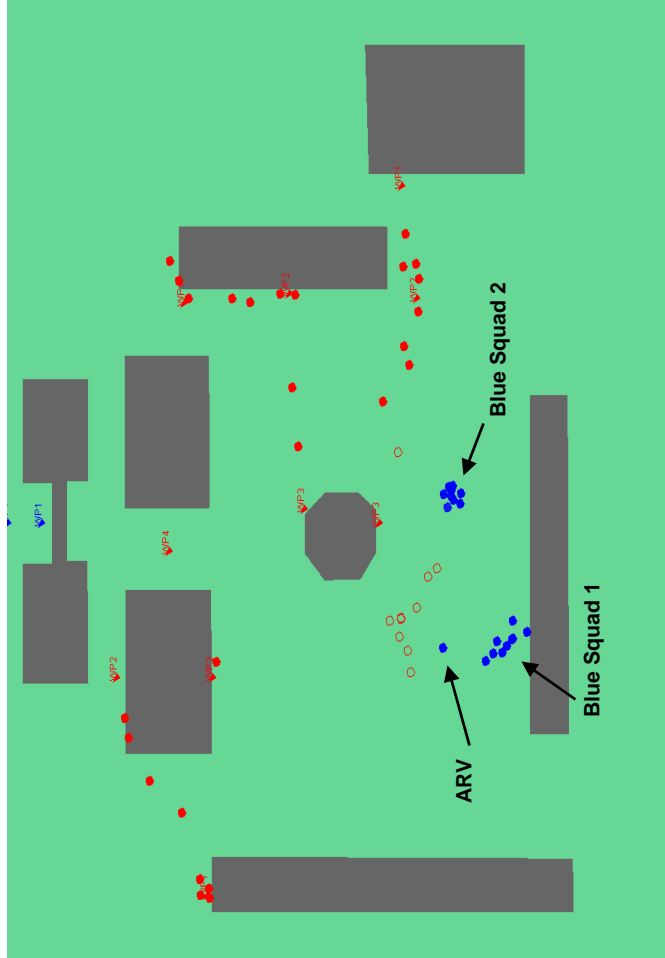
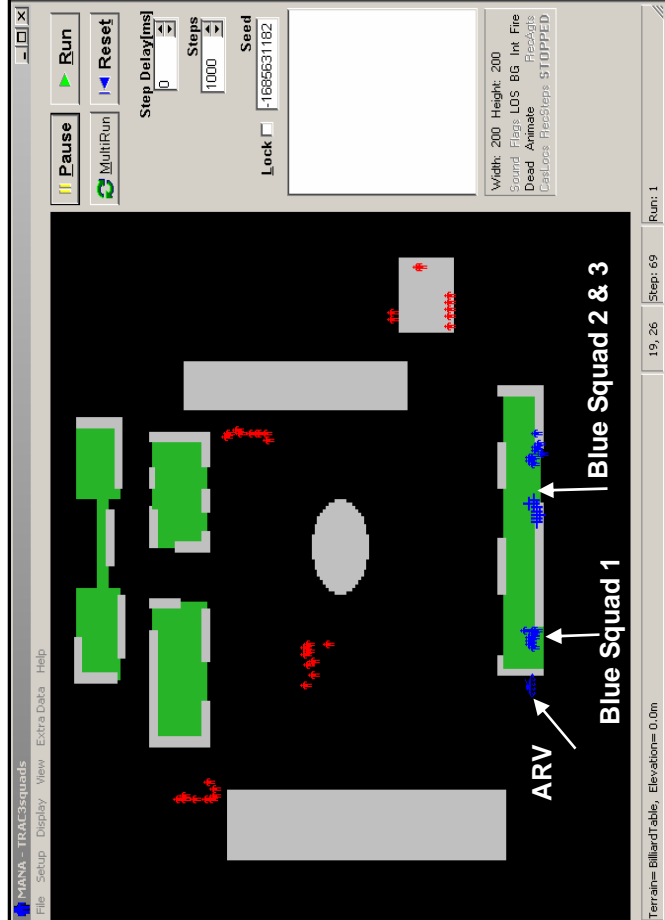
- **Objective Force (OF):**
 - What is the appropriate squad size and number of squads?
- **Armed Robotic Vehicle (ARV) FCS Issues**
 - What is the best Operational Employment Concept?

Basic Scenario

Basic scenario developed with guidance from DBBL's Chief of Analytical Simulations: Reference document – Army FCS Unit of Action Systems Book ver 1.2 - 1.6

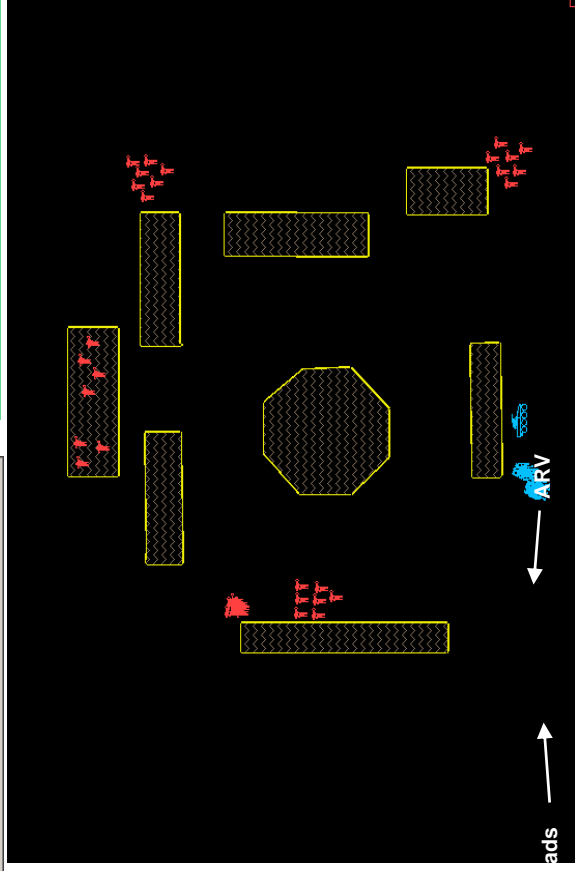
- **Objective Force infantry platoon with Armed Robotic Vehicle (ARV) against an organized threat in an urban environment. Blue to red force ratio (~1:1)**
- **Blue forces maneuver through urban environment to an seize an objective**
 - Blue squad size varied (7, 9, 12 per squad)
 - Blue number of squads varied (2, 3, and 4 squads)
 - Blue soldiers maneuver to avoid red contact
 - ARV operates in coordination with blue squads
 - ARV maneuvers to engage red forces when red is within specified operating ranges of blue forces
- **Red forces organized into 4 groups of 9 which patrol the urban area**
 - Red engages blue forces when blue is within red sensor range

Three Simulation Scenarios



MANA Scenario

Pythagoras Scenario

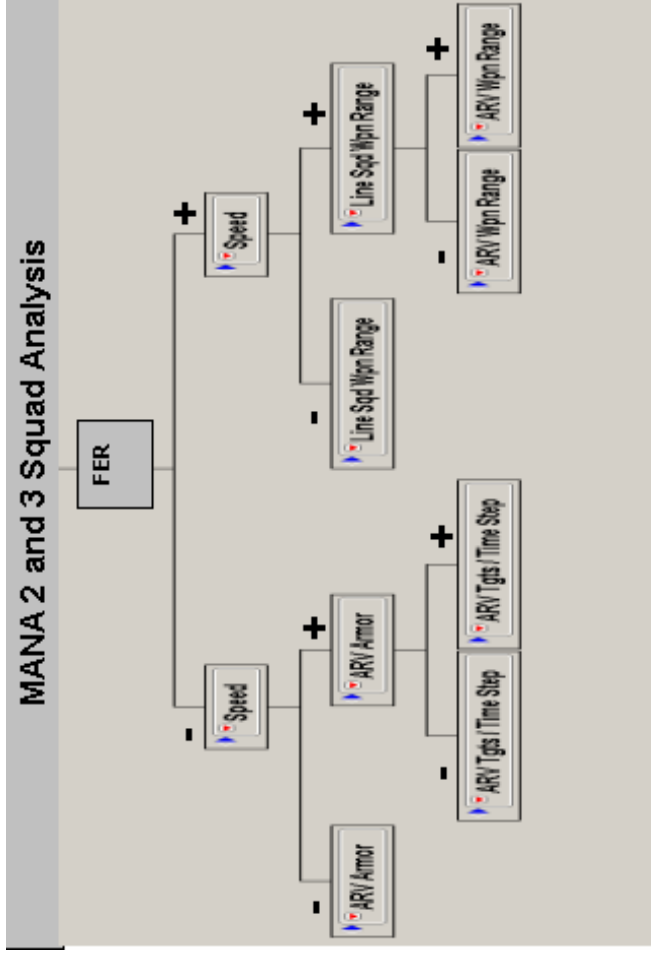
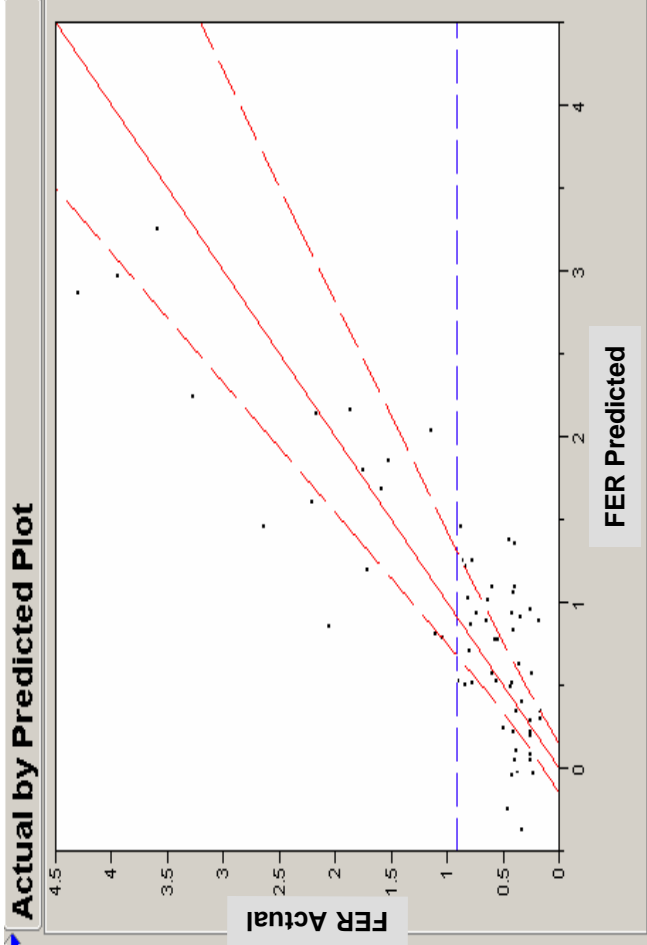


JANUS Scenario

MANA Data Analysis – (2 and 3 Squads)

Regression equation found using stepwise approach to identify significant effects. Results below show differences in predicted response (FER) and actual response (FER).

The classification tree illustrates the significant relationships in the identified effects.



• Armed robotic vehicle (ARV) speed

• Armor thickness of the ARV

• Squad size

• Combined effect of the ARV's primary weapon and firing rate of line squad

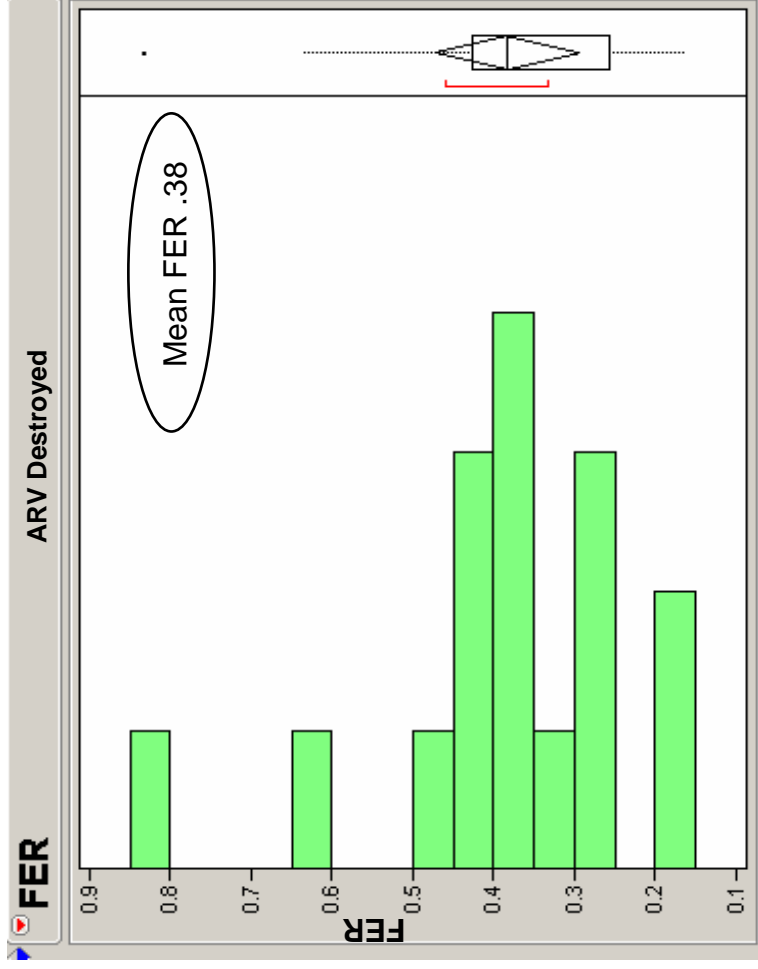
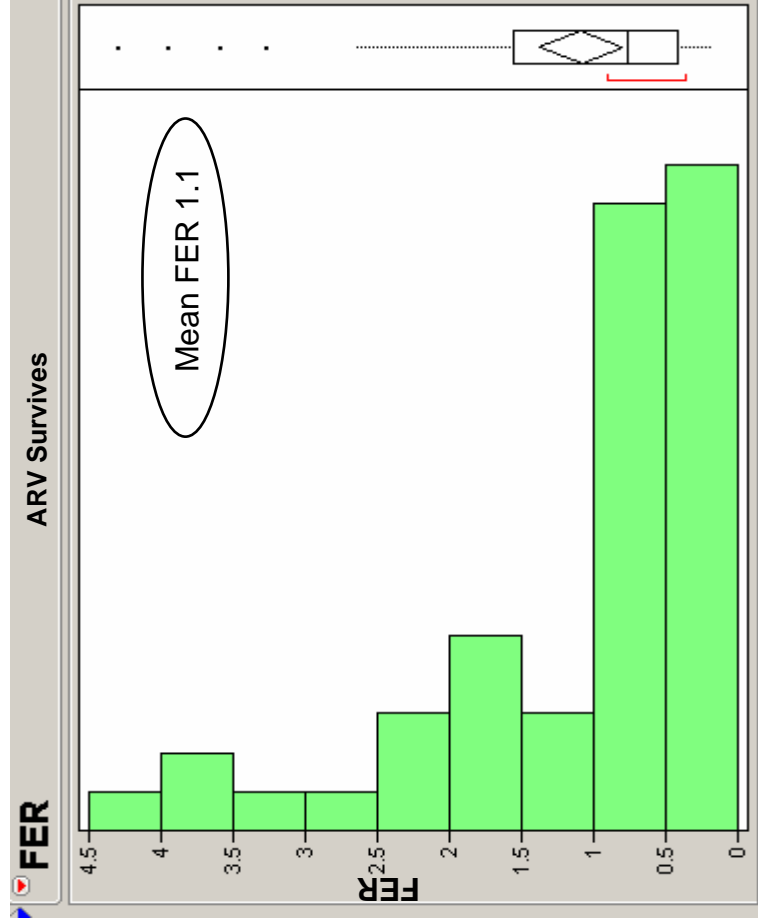
• Combined effect of M240 maximum effective ranges (MER's) and line squad's firing accuracy

• When ARV operating at low speeds, ARV armor thickness key to increased blue survivability.

• When ARV operating at high speeds, MERs of line squad(s) and ARV primary weapon key to increased blue survivability.

MANA Data Analysis – FER (2 and 3 Squads)

Comparison of FERs when ARV survives and when ARV is destroyed



Blue survivability is most related to the survivability of the ARV in this urban scenario. When the ARV survives, the Force Exchange Ratio (FER) is nearly 3 times greater.

Potential Insights from MANA Data Analysis

- **Objective Force (OF):**

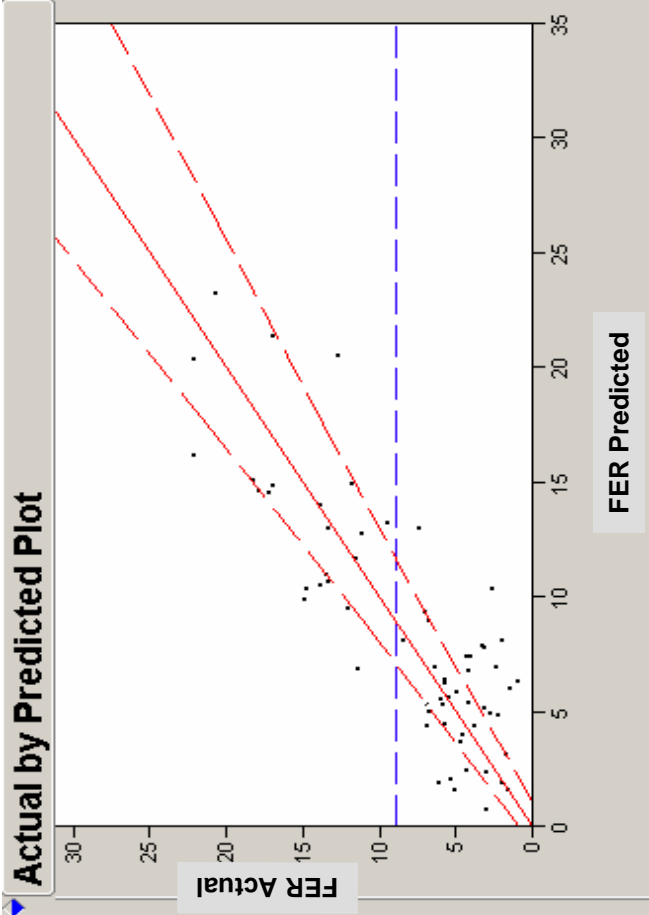
- What is the appropriate squad size and number of squads?
- When the ARV survives, no significant differences in survivability for squad sizes of 9 and 12 exist.
- When the ARV is destroyed, squad sizes of 12 offer significantly improved survivability.

- **Armed Robotic Vehicle (ARV) FCS Issues**

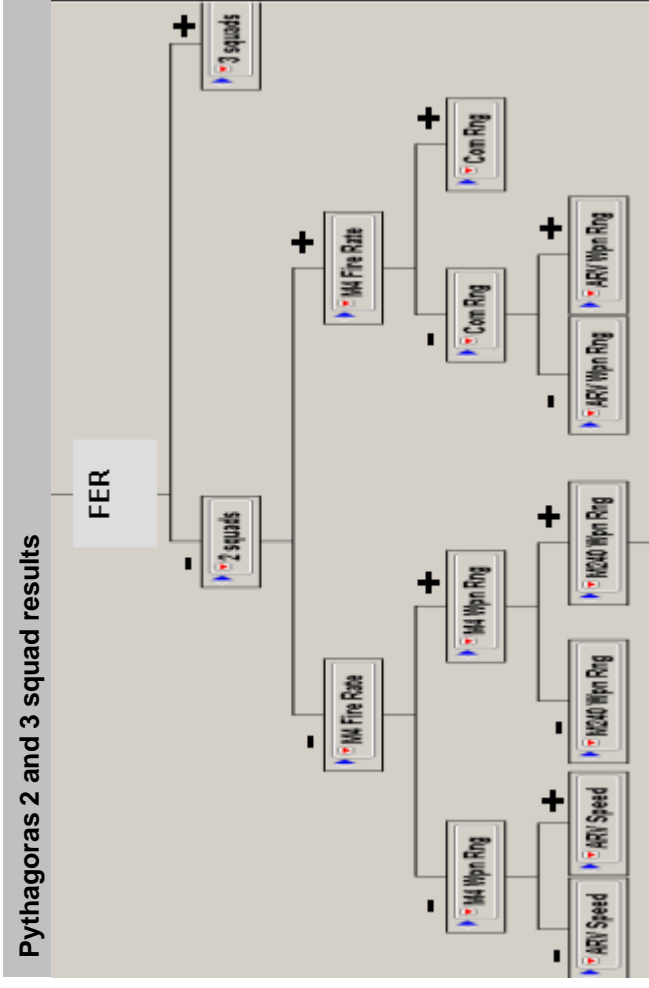
- What is the best Operational Employment Concept?
 - If the ARV operates at roughly the same speed and in close proximity to the dismounted infantry, then the ARV's armor thickness and ability to engage multiple targets are significantly important to improved FER.
 - As the ARV operates forward of and faster than the speed of the dismounted infantry, the ARV's 25mm cannon is the primary contributor to an improved FER.
- Higher FER {

Pythagoras Data Analysis – (2 and 3 Squads)

Regression equation found using stepwise approach to identify significant effects. Results below show differences in predicted response (FER) and actual response (FER).



The classification tree illustrates the significant relationships in the identified effects.

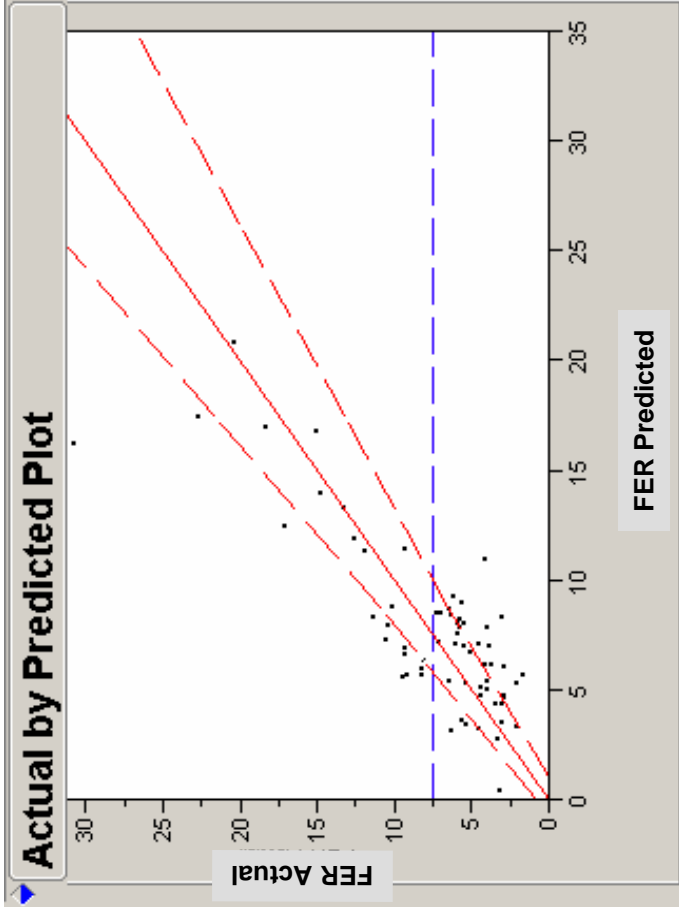


- Number of squads and squad size
- Weapon squad's automatic weapon (M4) firing rate
- Weapon squad's automatic weapon (M240) MER
- ARV speed
- Combined effect of number of squads and M240 MER
- With only 2 squads, then weapon squad firing rates, weapon squad MERs, and ARV speed are keys to blue survivability.

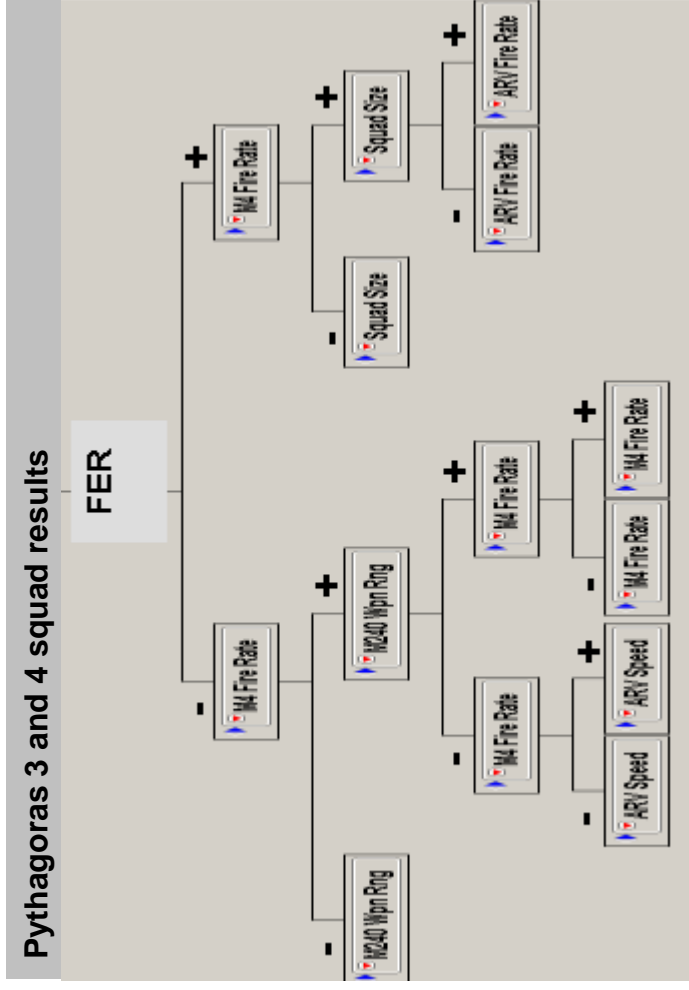
Pythagoras Data Analysis – (3 and 4 Squads)

Regression equation found using stepwise approach to identify significant effects. Results below show differences in predicted response (FER) and actual response (FER).

The classification tree illustrates the significant relationships in the identified effects.



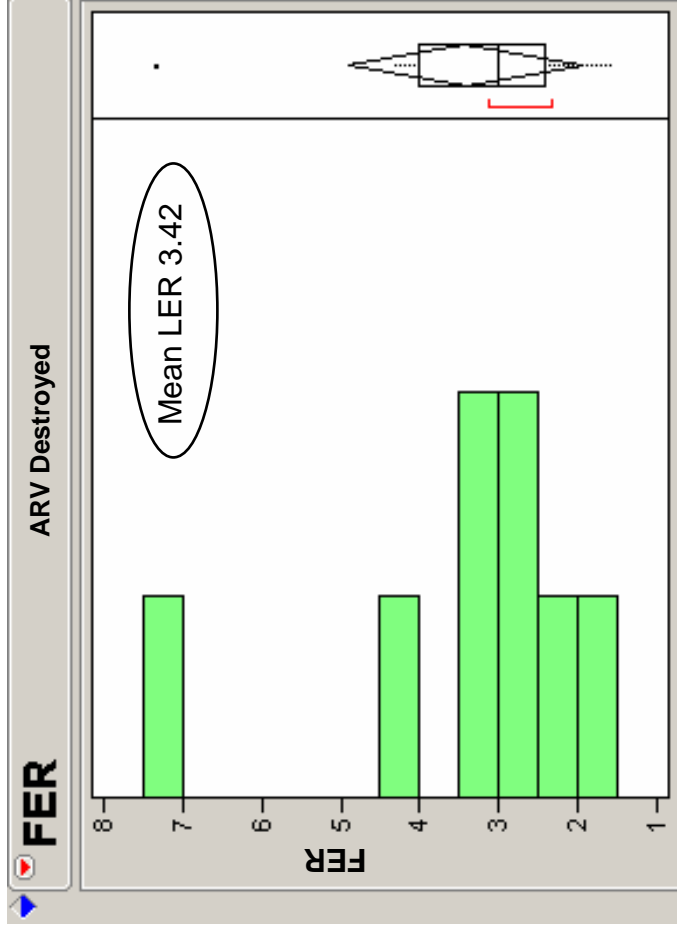
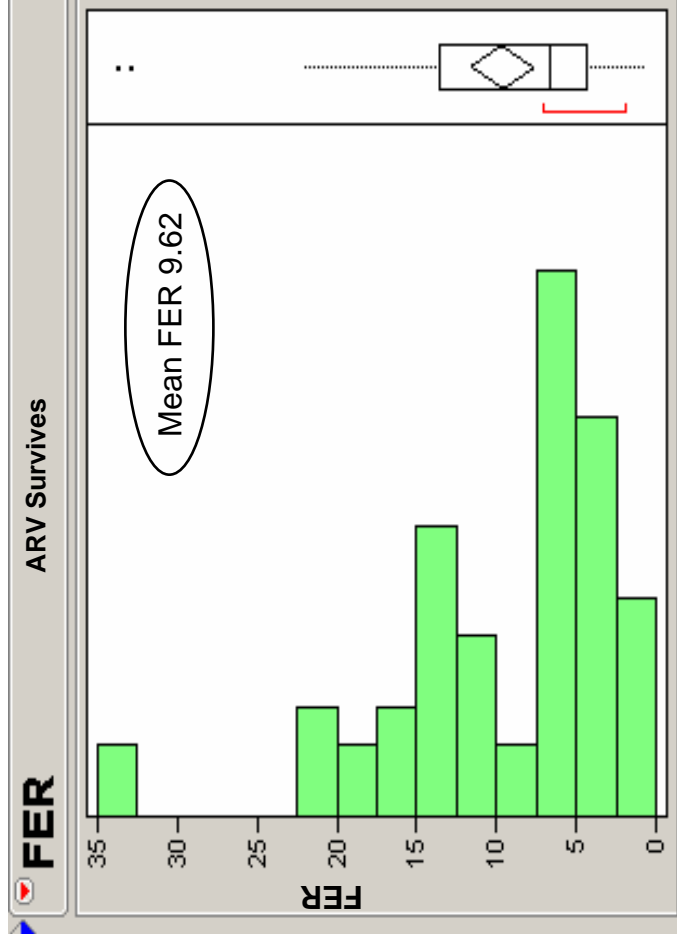
- Weapon squad M4 firing rate, M4 MER, and M240 MER
- Squad size and number of squads
- ARV speed and ARV primary weapon firing rate
- Combined effect of squad size and wpn sqd (M4) firing rate
- Combined effect of number of squads and line sqd firing rate



- When M4 firing rates are high, effect of larger squad sizes on blue survivability is mitigated.
- When M4 firing rates decrease, M240 MER and ARV speed improve blue survivability.

Pythagoras Data Analysis - FER (2 and 3 Squads)

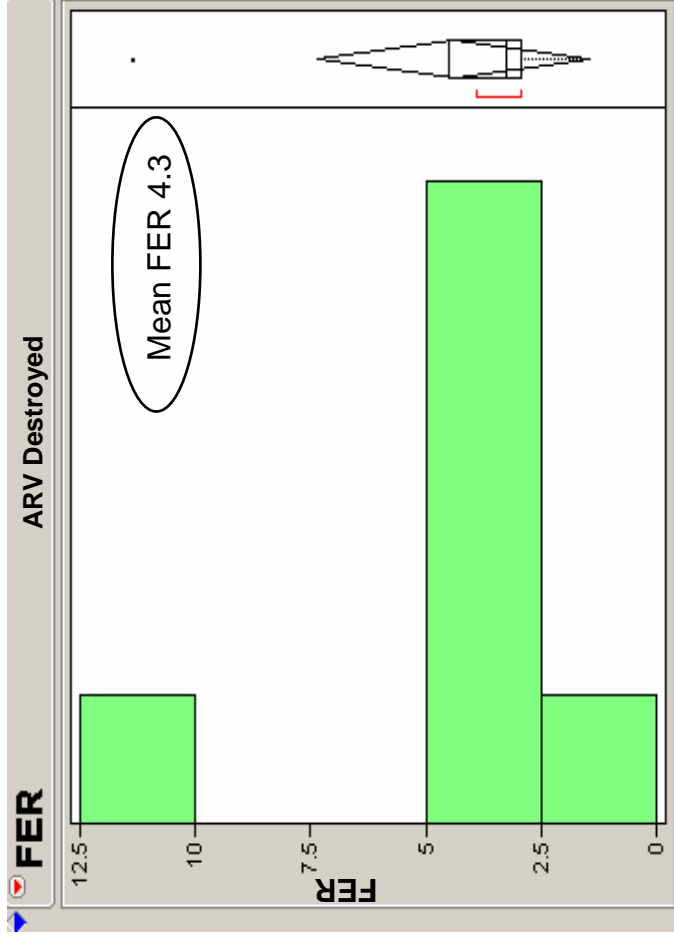
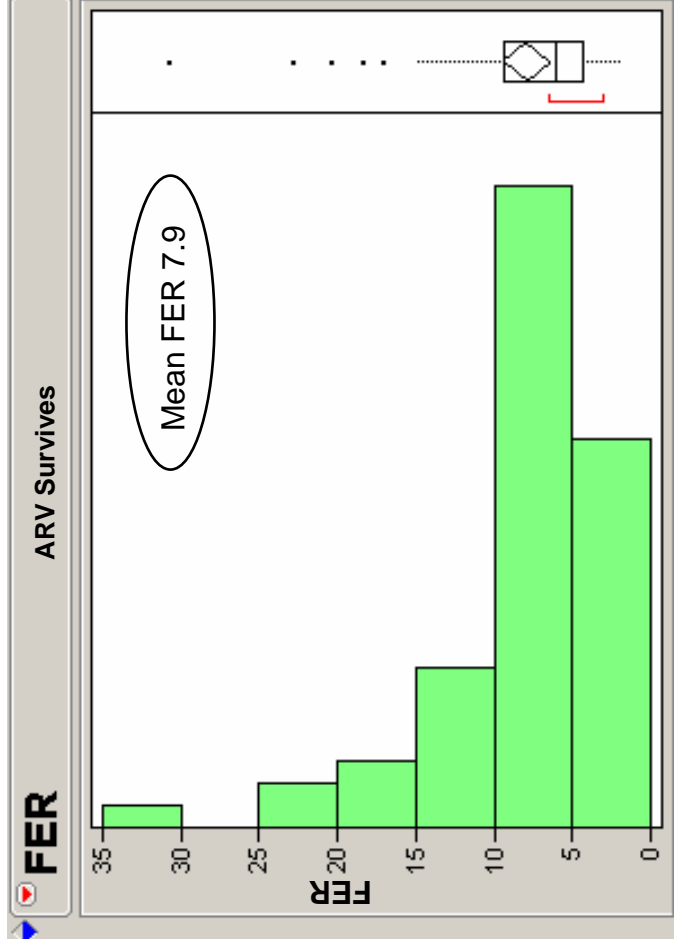
Comparison of FERs when ARV survives and when ARV is destroyed



Similar to previous result, Blue survivability is most related to the survivability of the ARV in this urban scenario. When the ARV survives, the Force Exchange Ratio (FER) is nearly 3 times greater.

Pythagoras Data Analysis - FER (3 and 4 Squads)

Comparison of FERs when ARV survives and when ARV is destroyed



Even with an increased number of squads, Blue survivability is most related to the survivability of the ARV in this urban scenario. When the ARV survives, the Force Exchange Ratio (FER) is nearly 2 times greater.

Potential Insights From Pythagoras Data Analysis

- **Objective Force (OF):**

- What is the appropriate squad size and number of squads?
 - With 4 squads of 12 soldiers each, the survivability of the ARV does not significantly impact the FER.
 - When the ARV is destroyed, squad sizes of 12 offer significantly improved survivability.
 - Improved marksmanship with the M240 and M4 offsets the need for increased squad sizes.

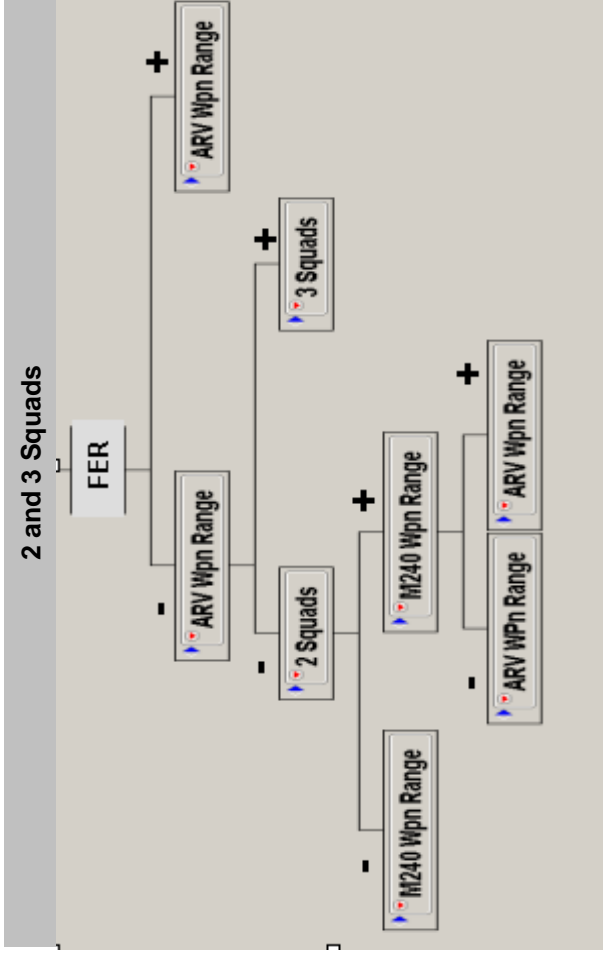
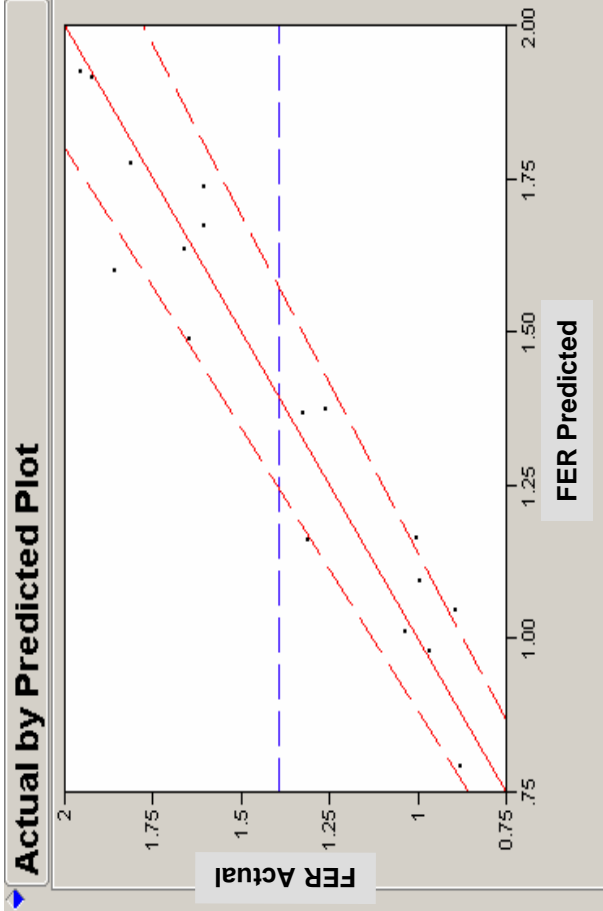
- **Armed Robotic Vehicle (ARV) FCS Issues**

- What is the best Operational Employment Concept?
 - As the ARV operates forward of and faster than the speed of a dismounted infantry platoon with 2 or 3 squads, the increased speed of the ARV and its ability to provide fire support offsets the reduced organic firepower of the platoon.

JANUS Data Analysis – (2 and 3 Squads)

Regression equation found using stepwise approach to identify significant effects. Results below show differences in predicted response (FER) and actual response (FER).

The classification tree illustrates the significant relationships in the identified effects.



- ARV primary weapon MER
- Number of squads and squad size
- Weapon squad's M240 MER
- Combined effect of squad size and ARV primary wpn range
- Combined effect of Scheme of Maneuver (SoM) and single shot kill probability of ARV

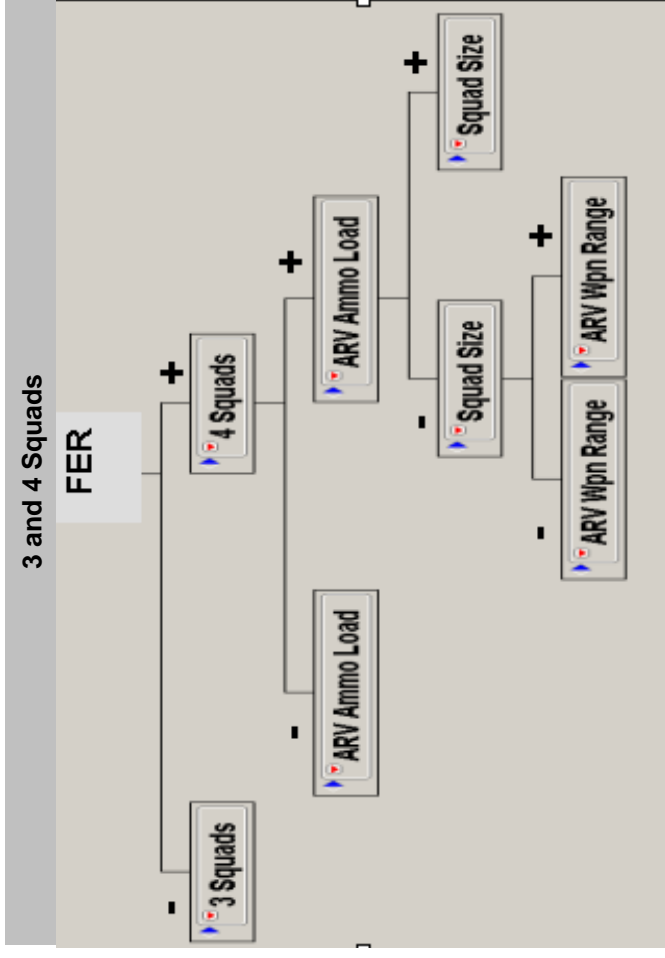
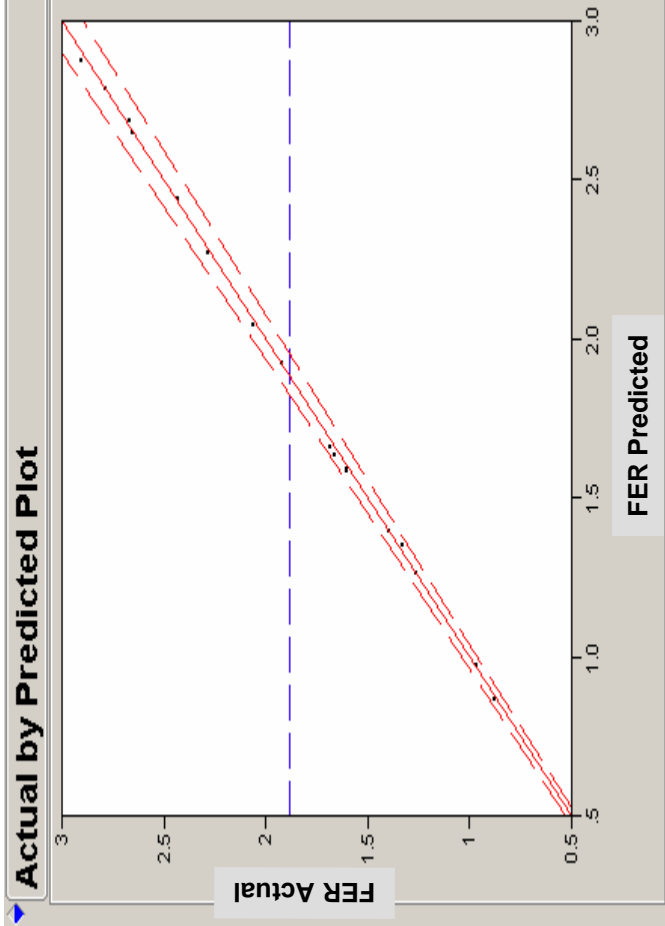
- ARV primary weapons MER is critical regardless of number of squads
- With only 2 squads present, weapons squad's MERs and ARV weapon range are keys to blue Survivability

* JANUS runs executed in closed form

JANUS Data Analysis – (3 and 4 Squads)

Regression equation found using stepwise approach to identify significant effects. Results below show differences in predicted response (FER) and actual response (FER).

The classification tree illustrates the significant relationships in the identified effects.

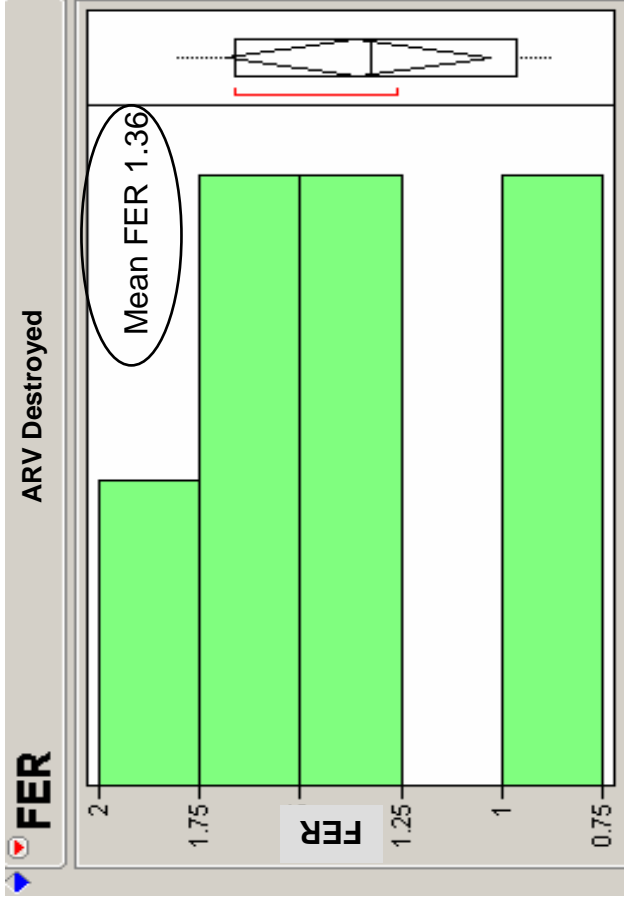
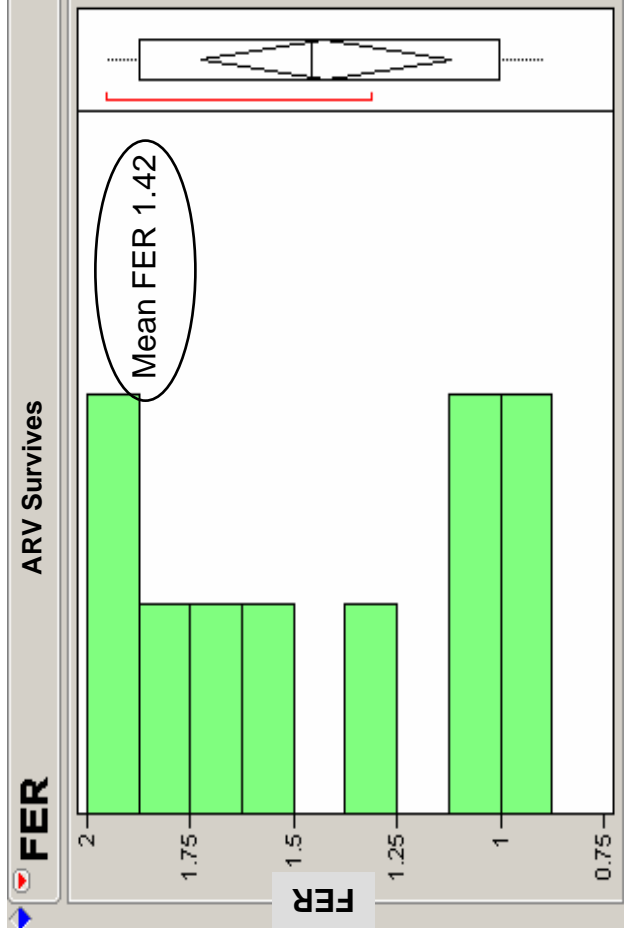


- Number of squads
- ARV primary weapon range
- Combined effect of squad size and M240 MER
- Combined effect of number of squads, squad size and Scheme of Maneuver

- Number of squads is primary factor in blue survivability.
- With 4 squads present:
 - ammunition usage increases significantly
 - squad sizes less than 12 indicate need for greater ARV MER

JANUS Data Analysis- FER (2 and 3 Squads)

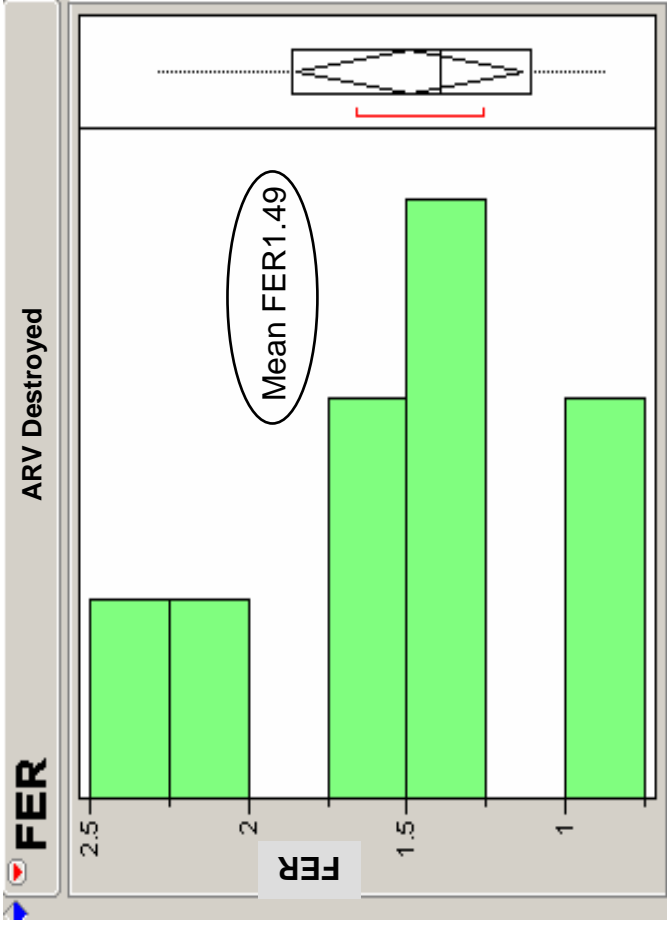
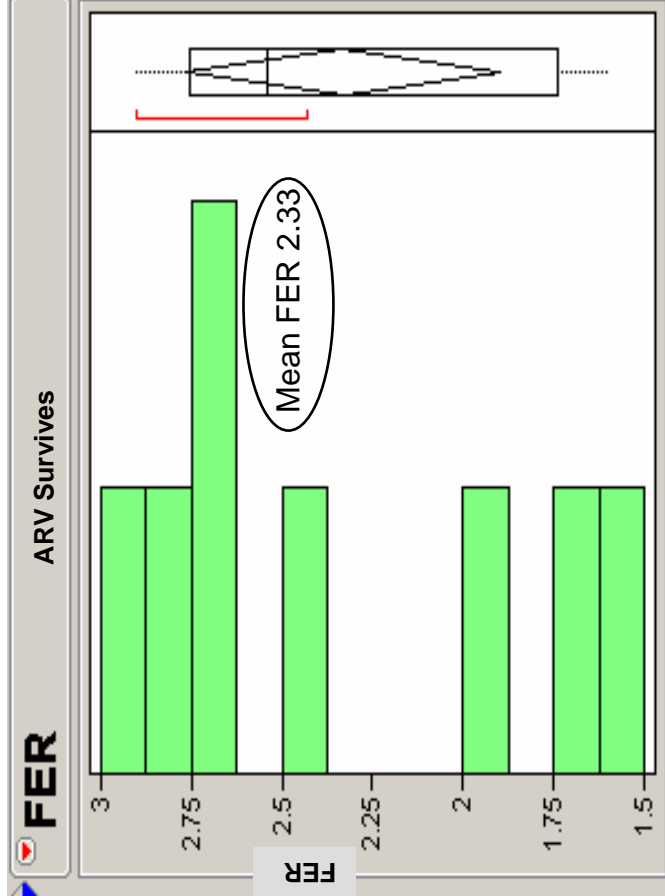
Comparison of FERs when ARV survives and when ARV is destroyed



These results indicate there is not a significant difference in FER's between 2 and 3 squads and ARV survivability.

JANUS Data Analysis - FER (3 and 4 Squads)

Comparison of FERs when ARV survives and when ARV is destroyed



Even with an increased number of squads, Blue survivability is most related to the survivability of the ARV in this urban scenario.

Potential Insights From JANUS Data Analysis

- **Objective Force (OF):**
 - What is the appropriate squad size and number of squads?
 - 4 squads had the highest FERs, but consumed 30% more ammunition.
 - If ARV is destroyed, squad sizes of 12 suffered minimal decrease in FER (in these cases, Blue had a 1.3:1 force advantage).
 - Improved marksmanship with the M240 and M4 offsets the need for increased squad sizes.
- **Armed Robotic Vehicle (ARV) FCS Issues**
 - What is the best Operational Employment Concept?
 - When the ARV had greater stand-off advantage and better capability withstanding anti-armor weapons, FERs were improved.

Note: JANUS Scheme of Maneuver and scripting of routes for both soldiers and ARV can greatly influence outcome.

Analysis Summary

Significant Factors Potential Insights

<ul style="list-style-type: none"> - ARV armor thickness -Wpn's sqd max weapon ranges and firing rates -ARV speed -ARV wpn max range and sqd firing rates - Squad size 	<ul style="list-style-type: none"> - ARV survives, no difference in squad survivability - ARV destroyed, squad size of 12 improved survivability - ARV operates at low speeds, armor thickness key - ARV operates at high speeds, weapon range key
<ul style="list-style-type: none"> -Squad size -Number of squads -Number of sqds and sqd firing rates -Sqd firing rates and squad size Wpn's sqd max weapon ranges and sqd firing rates - ARV speed 	<ul style="list-style-type: none"> - With 4 squads of 12, ARV does not significantly impact FER - Improved wpn's sqd firepower can offset need for increased squad sizes - ARV destroyed, squad size of 12 improved survivability - ARV operates at high speeds, offsets reduced sqd organic firepower
<ul style="list-style-type: none"> - Number of squads - Squad size - ARV wpn max range - Scheme of Maneuver (SOM) -Wpn's sqd max weapon ranges -ARV armor thickness 	<ul style="list-style-type: none"> - With 4 squads, highest FER but large ammo consumption - ARV destroyed, squad size of 12 improved survivability - Improved wpn's sqd firepower can offset need for increased squad sizes - Increased ARV wpn max range improved FER

MANA

Pythagoras

JANUS

Simulation Comparison

	MANA	Pythagoras	JANUS
Squad Size	S	M	M
Number Squads	NS	M	M
ARV Armor Thickness	M	NA	S
ARV Speed	S	S	S
ARV Wpn Max Range	S	S	S
Wpn Max Ranges and Sqd Firing Rates	S	S	NA
Wpn Max Ranges	NS	S	S
Sqd Firing Rates and Squad Size	NS	S	NA
Number of Sqds and Sqd Firing Rates	NS	S	NA
Scheme of Maneuver	NA	NA	S

Most Similar Least Similar
 Simulations were consistent

M = Most Significant **S** = Significant **NS** = Not Significant **NA** = Not Applicable

Insights Summary

OF Platoon and FCS Summary:

- ARV survivability has a significant impact upon blue squad survivability in this type of urban scenario. ARV survivability is enhanced by its weapon's max effective range and protection against anti-armor weapons.
- If the ARV survives, there is little difference in LER for squad sizes of 9 and 12. However, if the ARV is destroyed, squad sizes of 12 have better survivability.
- In general, if the ARV meets future specifications, 3 squads are sufficient.

ABS Summary:

- Agent based simulations allow a more robust exploration of the possible outcomes for urban operations. They allow the generation of ideas and insights prior to bringing significant resources to bear on an issue.
- In this example, agent based simulation results are relatively consistent with high resolution combat simulations.

Future Research

- **Use of agent based simulations for screening analysis**
 - FCS Key Performance Parameter analysis
- **ABS representing millions of agents in an urban environment.**
 - Drafting DARPA Proposal with Los Alamos National Laboratory
- **Use of new experimental designs for high resolution simulations.**
 - AMSAA's FCS System-of-Systems metrics
- **Facilitate Project Albert's research into ABS for urban operations.**
 - Command and Control research (1 and 2 levels higher than basic agent)
 - More robustly capturing sensor and targeting information
 - Logistics and humanitarian assistance operations

Historical Reference - Squad Size Studies

- **WWII** – 12 man squad which consisted of 3 teams (2 ,3,5) plus a squad leader
- **1946** Infantry Conference – Recommended reduction to a 9 man squad of 2 teams(4,4) and a squad leader. 1953 T/O &E changed to 2 teams of 4 ([Main reason improved controllability](#))
- **1956** “A Study of the Infantry Rifle Squad TOE” (ASIRS) – 11 man squad of teams (5,5) plus squad leader adopted ([Main reasons improved fire power, ability to move and cover and resiliency](#))
- **1960** Ft Benning, Infantry School study – conclusion: 11 man squad worked so well it should not be changed.
- **1961** US Army conducted 2 studies:
 - “Optimum Composition of the Rifle Squad and Platoon” (OCRSP) – Ft Ord, California
 - Conclusions: 11 man squad of 2 teams (5,5) plus squad leader was optimal. Test was qualitative and quantitative. ([Main reasons fire power, resiliency, staying\(absorb 25% losses\), good leadership ratio](#))
 - ”Rifle Squad and Platoon Evaluation Program” (RSPEP) – Ft Benning, Infantry School
 - Conclusion: 10 man squad of 2 teams (4,5) plus squad leader was optimum. No objective data used from other study. ([Main reasons optimal leader to lead ratio, flexibility](#))
 - Infantry School accepted 10 man squad. Reorganization of Active Army Division (ROAD) T/O&E 7-18E. Took into Vietnam.
- **1966-1969** US Army conducted a series of studies “Infantry Rifle Unit Study” (IRUS)for implementation in the 1970-1975 force. Very exhaustive study looking at 10, 11, and 13 man squads with various weapons and included physiological factors.
 - 1973 Army acted on the IRUS studies and increased squad size to 11 man squad of 2 teams (5,5) plus squad leader ([Main reasons fire power, resiliency, staying\(absorb 25% losses\), good leadership ratio. Close on 11/13-cost was factor](#))
- **1973** IRUS recommendation increased squad size to an 11 man squad.
- **1979-1980** Army conducted Division 86 study series with transition to take place in 1983. Also, 1983 Army Commanders Conference drove another study, Army of Excellence.
- **1984** Army implemented changes based on above. The Army fielded a 9 man squad of 2 (4,4) teams plus a squad leader. Decision driven by resource constraints and not by validated field test. This structure used in Panama (1989), Persian Gulf War (1990-91) and Somalia (1993)

Parsimonious survival analysis models: Studying early attrition in the armed services by frailty and time-dependent survival analysis

Yuanzhang Li, Timothy Powers, and Margot Krauss

Introduction

Early attrition among new enlistees is a costly problem for the federal military. It is of interest to recruiters and planners to have an idea of how much attrition will take place, as well as when it will occur. It is also of interest to know what factors are related to likelihood of attrition so that the best prospects for retention can be sought. Survival analysis is one way of addressing these needs, and the current paper will examine different types of survival models for suitability to this task.

A popular approach to survival analysis problems with multiple covariates is Cox proportional hazards modeling. Also known as “semi-parametric” modeling, this approach does not require such strenuous assumptions as to the functional form of the survival curve as do parametric models. Two critical assumptions with regard to functional form are 1) a multiplicative relation between the underlying hazard function and the log-linear function of the covariates; and 2) the hazard associated with any particular combination of factor levels is proportional over time to that associated with any other combination of factor levels. Quite conveniently, this assumption holds true for many of the more commonly used parametric models, including the exponential and the Weibull.

However, while the assumption of proportional hazards is less restrictive than those needed for parametric modeling, even this lesser assumption is sometimes of dubious merit. In particular, the first assumption stated above outcome of interest may be related to factors beyond those included in the model, resulting at best in bias of effects estimates. Another possible violation of the proportional hazards assumption is that the effects of some factors on the outcome of interest might not be constant over time. Frailty models were proposed by Vaupel et al (1979) to deal with the problem of underspecification or misspecification of covariates in a proportional hazards setting. In essence, these models allow for the survival among individuals with the same levels of predictor factors to have survival that differs according to a distribution, rather than simply by random error.

The second assumption listed above (time independence) might also not be true. For example, a military recruit with a pre-existing medical condition might be at increased risk of early attrition, although presumably more so at the beginning of training than once he/she has successfully undergone rigorous training. In such a case, the impact of the covariate is clearly dependent on time. Non-proportional hazards models are used to extend the proportional model to analyze such data. The user can specify arithmetic expressions to define covariates as functions of several variables and survival time.

The aim of this study is to develop more accurate depictions of early military attrition using these survival model variations.

Subjects and Methods

All first-time enlistees beginning active duty enlisted service in the Army during January 1998 - December 2001 were included in the analyses. Accession records on these individuals were linked with military personnel records to determine whether or not a subsequent early attrition occurred. In addition to the accession data, the demographic factors, such as sex, age, race, education, Armed Forces Qualification Test scores (AFQT), medical qualify, Body Mass Index (BMI), married status, the number of dependents as well as the geographic factor, the Military Entrance Processing Station (MEPS) are included in the model. These factors have been found in previous studies to be strongly related to likelihood of attrition.

Three types of survival models are used to relate attrition to these multiple factors:

Cox-Proportional Model-

$$h_i(t) = h_0(t) \exp(\beta_1 x_{i1} + \dots + \beta_k x_{ik}) \quad 1$$

where h is the hazard function, which is the instantaneous probability of failure given survival up to t and h_0 is the baseline hazard.

Frailty Model-

Unexplained variability, that not accounted for by including covariates, is known as overdispersion. Overdispersion is caused either by misspecification or omitted covariates, and makes the assumption of the proportional model invalid. The hazard is

$$h_i(t | z) = h_i(t) z_i$$

A frailty model attempts to account for variability that is not accounted for by the included covariates (overdispersion) by the following model:

$$h_i(t) = h_0(t) z_i \exp(\beta_1 x_{i1} + \dots + \beta_k x_{ik}) \quad 2.1$$

Here, z_i varies across the individuals. However, frailty models are also used to model survival times in the presence of group-specific random effects. Such models are termed "shared" frailty models, and depicted as follows:

$$h_{ij}(t_{ij}) = h_0(t_{ij}) z_j \exp(\beta_1 x_{ij1} + \dots + \beta_k x_{ijk}) \quad 2.2$$

In general, the random unobservable frailty effects are often assumed to follow either a gamma or inverse-Gaussian distribution. For the shared group-specific frailty, the model is constrained to be equal over those observations from a given group or panel.

Time-Dependent Proportional Model:

As discussed above, the hazard ratios for different combinations of factor levels may vary according to the time. Hence we should consider a time-dependent hazard model such as the one depicted below:

$$h_{ij}(t_{ij}) = h_0(t_{ij}) z_j \exp(\beta_1 x_{ij1}(t_{ij}) + \dots + \beta_k x_{ijk}(t_{ij}) + f(t_{ij})\eta) \quad 3.1$$

where $f(t)$ is a function of the survival time t . In this study, we consider a special case of 3.1.

$$h_i(t) = h_0(t) \exp\left(\sum_{j=1}^K x_{ij} \beta_j [1 + \gamma_j \ln(t) + \eta_j \ln(t)^2]\right) \quad 3.2$$

Note that the natural logarithm of survival time, $\ln(t)$, is used rather than survival time t , as suggested by many previous investigators. The square of $\ln(t)$ is used to measure non-linear effect of time.

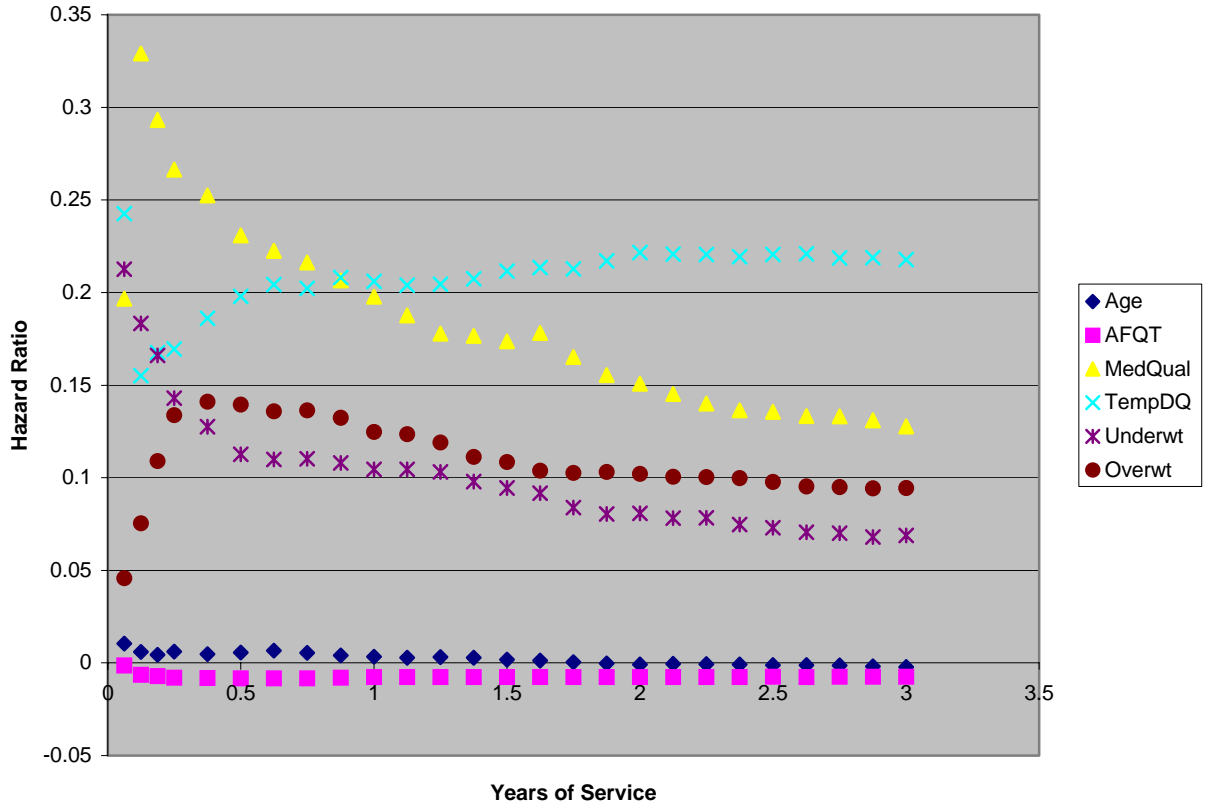
Some special cases of model 3.2 should be noted. If $\eta_j=0$ for $j=1, \dots, K$, then model 3.2 becomes the time-dependent linear model. If $\gamma_j=\gamma$, $\eta_j=\eta$ for $j=1, \dots, K$, then model 3.2 becomes a special case frailty model. Finally, if $\gamma_j=\eta_j=0$ for $j=1, \dots, K$, then model 3.2 becomes the Proportional Hazards Model.

A basic principal in model estimation is to include only those prediction factors that have some influence on the dependent variable. Hence, in the estimation process, we will delete all terms with non-significant γ or η . If none of these terms is significant, the proportional hazards model can be taken as the appropriate model.

Results

First we applied a proportional hazards model within different years of services to examine the effects of the various factors on hazard of attrition over the first three years of service. Figure 1 shows changes in some of the estimated factor effects over the examined time period. It is clear that the effects of several factors vary over time, meaning that the proportional hazards assumption is not tenable. For example, the effect of being overweight increases over the first half-year of service, then diminishes gradually over the next 2 ½ years. Similar observations of changing effects can be seen for being underweight, being medically qualified, and being temporarily medically disqualified. The effects of age were a little higher for the first year of service, then relatively stable around the same value afterwards. The effect of AFQT score was relatively small and appeared more constant over time, with a minor wavering.

Hazard Ratios in The Army



Given that the proportional hazards assumption is not met by the Army attrition data, we proceed to consider a frailty model. We consider the frailty to be due to the geographic regions from which military applicants come, and the demographic distributions of applicants from the various regions. Table 1 shows the differences in factor effect estimates between the frailty model and the simpler Cox proportional hazards model.

It is seen that the inclusion of frailty has virtually no effect on factor effect estimates for attrition within the first year of service. However, the consideration of frailty does result in altered effect estimates at longer periods of service. For example, the effect of being black (and of being white, for that matter) on attrition within the first three years of service is estimated to be much less when frailty is considered than when it is not considered. Other factor effect estimates are seen to show such differences at the longer time range.

Table 1. Factor Coefficients With and Without Frailty Consideration

Factor	Within 1 Year		Within 2 Years		Within 3 Years	
	No Frailty	With Frailty	No Frailty	With Frailty	No Frailty	With Frailty
Age	0.0033	0.0033	-0.001	-0.001	-0.002	-0.002
AFQT	-0.0076	-0.0076	-0.008	-0.008	-0.007	-0.011
Black	0.0813	0.0811	0.092	0.094	0.106	0.083
Dependents	0.0374	0.0373	0.063	0.065	0.059	0.071
Less than HS	-0.0478	-0.0479	-0.063	-0.064	-0.049	-0.050
Married	-0.1533	-0.1528	-0.161	-0.167	-0.175	-0.178
MedQual	0.1977	0.1975	0.151	0.147	0.128	0.119
Single	-0.3551	-0.3539	-0.341	-0.360	-0.352	-0.373
White	0.4810	0.4811	0.463	0.470	0.448	0.376
Underwt	0.1045	0.1046	0.081	0.079	0.069	0.065
TempDQ	0.2060	0.2059	0.222	0.217	0.218	0.211
Overwt	0.1248	0.1248	0.102	0.104	0.095	0.128
Sex	0.6662	0.6667	0.628	0.605	0.601	0.479

The relation of several of the considered factors to attrition has been seen to differ according to the length of time served. We therefore consider a time-dependent model to account for this time dependency. The effect of being overweight showed a non-linear pattern, so a non-linear time effect is considered for it. The effects of sex, medical qualification status, temporary medical disqualification, being underweight, having less than a high school education, being black and being white showed a linear relation with time, and are modeled accordingly. Finally, the effects of age and AFQT score did not show a relation to time, therefore no time component is considered for these factors.

Table 2 shows hazard ratio estimates for three factors (sex, permanent medical disqualification, and temporary medical disqualification) from both the Cox proportional hazards model at the indicated time cut-points, and the Time-Dependent model. It is seen that the effect estimates for these three factors differ over time within each model, and across models. Further work is needed to account for factors influencing attrition that are not yet accounted for.

Table 2. Hazard Ratio Estimates for Selected Factors:
Cox Proportional Hazards Model vs. Time-Dependent Model

Months	Proportional Hazards Model			Time-Dependent Model		
	Sex	Perm't DQ	Temp DQ	Sex	Perm't DQ	Temp DQ
1.5	2.73	1.34	1.34	1.86	1.15	1.19
3	2.13	1.25	1.33	1.70	1.12	1.19
4.5	1.94	1.22	1.32	1.61	1.10	1.19
6	1.88	1.19	1.33	1.55	1.09	1.19
7.5	1.85	1.18	1.32	1.50	1.08	1.19
9	1.84	1.17	1.32	1.46	1.08	1.18
10.5	1.82	1.16	1.32	1.43	1.07	1.18
12	1.80	1.14	1.31	1.41	1.07	1.18
13.5	1.79	1.13	1.30	1.39	1.06	1.18
15	1.77	1.12	1.30	1.37	1.06	1.18
16.5	1.76	1.12	1.29	1.35	1.06	1.18
18	1.74	1.11	1.29	1.33	1.05	1.18
19.5	1.72	1.12	1.29	1.32	1.05	1.18
21	1.70	1.10	1.29	1.31	1.05	1.18
22.5	1.68	1.08	1.29	1.29	1.05	1.18
24	1.66	1.08	1.29	1.28	1.04	1.18
30	1.59	1.05	1.27	1.24	1.04	1.18
36	1.54	1.04	1.25	1.21	1.03	1.18

Conclusions

This study has found that the assumptions of the Cox proportional hazards model are not adequately met by Army attrition data. Several of the factors seen in many studies to be related to early attrition exhibit effects that are not constant over time in service, casting doubt on the proportional hazards model estimates. Consideration of a frailty model did not have an appreciable impact on factor effect estimates, although a time-dependent model did.

The time-dependent model considered here was derived in subjective fashion. Functional form of time-dependency for the various factors was decided by examining factor effects from Cox proportional hazards models applied at various time points. Future work will focus on a more objective way to develop the time-dependent model. Application to other services' attrition data will also be pursued.

Reference

Vaupel JW, Manton KG and Stallard E. The impact of heterogeneity in individual frailty on the dynamics of mortality. *Demography*. 16(1979)3: 439-454.

The Effect of Dosage Errors and Step Selection Method on the Performance of the Up-and-Down Method

Ninth US Army Conference on Applied Statistics
29 to 31 October 2003
Napa Valley Marriott, Napa Valley, CA

Douglas R. Sommerville, PE
Modeling Simulation & Analysis Team

DISCLAIMER: The findings presented in this briefing are not to be construed as an official Department of the Army position unless so designated by other authorizing documents.

Edgewood Chemical Biological Center
5183 Blackhawk Road, ATTN: AMSRD-ECB-RT-IM, Bldg. E5951
Aberdeen Proving Ground, Maryland, USA 21010-5424



Email: Douglas.Sommerville@us.army.mil
Phone: (410) 436-4253
FAX: (410) 436-2742



Background

- Up-and-Down (UaD) Method is a common approach for estimating median effective stress
 - Originally developed by the Explosives Research Laboratory, Bruceston, PA in early 1940's, and the Bruceston Method was basis for work of Dixon and Mood (1948)
 - Extensive documentation shows that the method provides an accurate estimate while keeping number of trials to a minimum
 - Numerous testing applications
 - Explosives
 - Metal fatigue
 - Development of Medical Procedures
 - Toxicology and Pharmacology
 - Median effective stress ==> median effective dose/dosage (ED₅₀)

Background--Conducting an UaD Bioassay Experiment

- Technique is simple to execute but has inspired extensive discussions on how to analyze the data
 - Trials are conducted one at a time until stopping criteria is met
 - For trial i , a subject is given some administered log dose (d_i) and the binary response is recorded (R_i)
 - The log dose for the next trial (d_{i+1}) is dependent on value for R_i
 - If a response occurs ($R_i = X$ or 1):
 $d_{i+1} = d_i - \Delta d$
 - If a response does not occur ($R_i = O$ or 0):
 $d_{i+1} = d_i + \Delta d$
 - Δd is the step size
 - In original method, Δd kept fixed in value throughout experiment
- Several modifications of the basic method have been developed
 - Example: various schemes for changing Δd during course of experiment

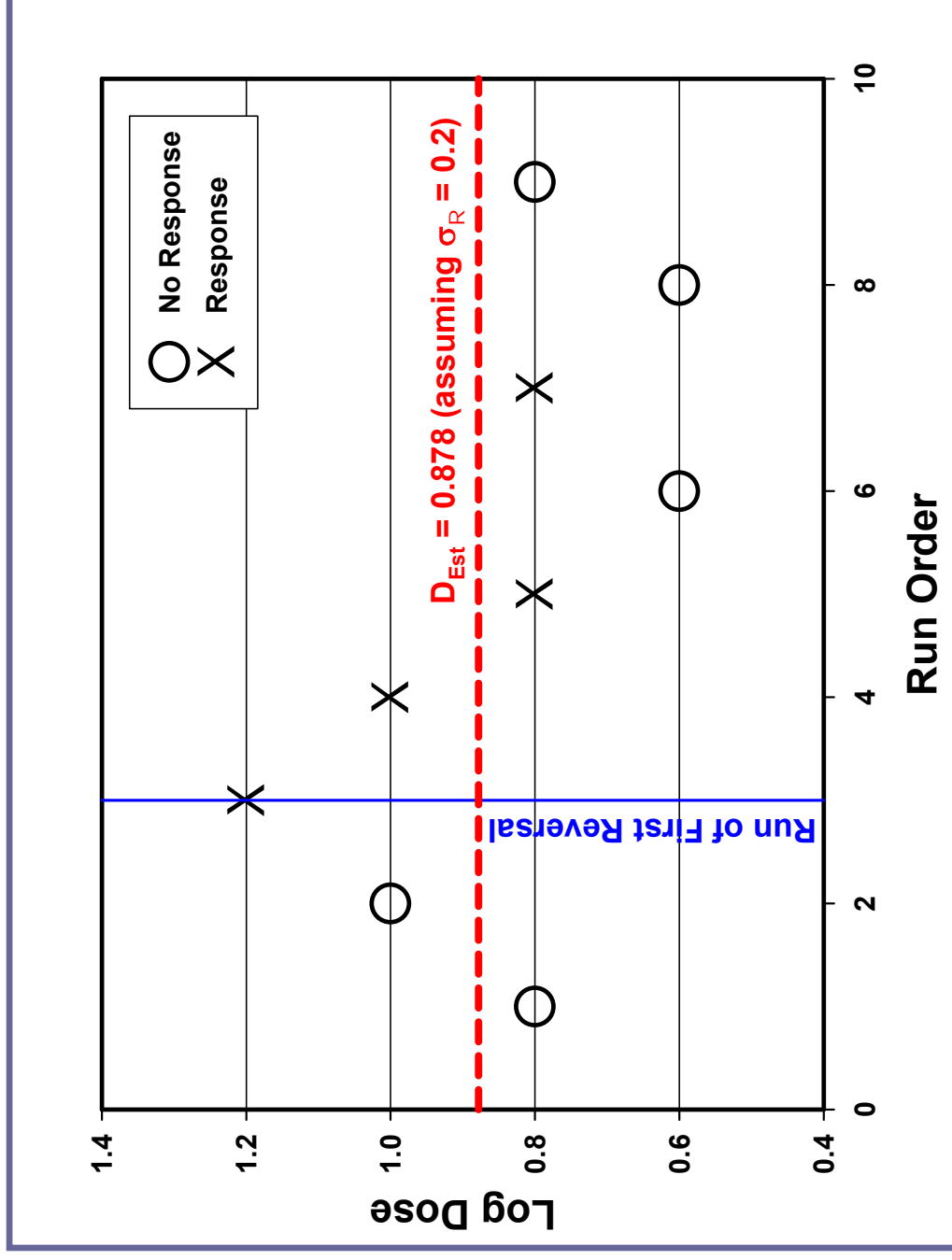
Obtaining an Estimate of the ED_{50}

- Maximum Likelihood Estimation (MLE) used to find estimate, D_{Est} , of the true $\log(ED_{50})$, D_{Tr}
- Normal distribution is favorite though others used (ex. Little (1974))
 - For bioassays, normal distribution obtained by working with the logarithms of the doses
- With small sample sizes, D_{Est} is solved for while σ_R is set equal to Δd (Dixon (1965))
- Precision of D_{Est} increases as Δd decreases
- Speed of convergence towards region of D_{Tr} decreases as Δd decreases

Characteristics of UaD Bioassay Experiment

- Optimum Δd range:
 - $(\sigma_R / 2) < \Delta d < (2\sigma_R)$
 - σ_R -- standard deviation of distribution of effective dosages, which needs to be estimated prior to start of experiment
 - Above range represents a trade-off between precision and efficiency
- Standard error (σ_{ED}) of D_{Est} from UaD experiment
 - σ_{ED} approximately equals $\{(\sigma_R)(2 / N)^{(1/2)}\}$ (from Dixon (1965))
 - It is assumed that σ_R equals Δd
 - N is the nominal sample size (versus N_T the total number of trials conducted)
 - Common practice is to base N on N_{FR} (run of first reversal {0 to X, or X to 0})
 - $N = N_T - N_{FR} + 2$

Example of an Up and Down Experiment



UaD Method and Errors in Dosage Administration and Measurement

- What happens when irregular spacing of dosages occurs--found with UaD in inhalation toxicology?
 - Dosage errors (assumed to be normally distributed)
 - Administration error: σ_A
 - Difficulty in precisely generating the target dosage
 - Measurement error: σ_M
 - Error involved in measuring the actual dosage produced by generation device
 - Total error: $\sigma_{\text{Total}} = (\sigma_A^2 + \sigma_M^2)^{(1/2)}$
 - Vast majority of work/theory on UaD method assumes that both types of errors essentially equal zero.
 - How are the accuracy and efficiency of UaD method affected by nonzero σ_A and/or σ_M ?

Experimental Method

- Monte Carlo simulations were performed using MINITAB® v. 13
- Three different dosage parameters were tracked
 - d_T -- The target value for the log dosage to be administered
 - d_A -- The actual log dosage administered
 - d_M -- The measured/observed value of the administered log dosage
- Binary responses generated from the following underlying normal distribution of effective dosages
 - D_{Tr} set equal to $\log_{10}(80)$
 - σ_R set equal to (1/15) (or 0.0667)
- Each UaD experiment consisted of 10 trials
 - Wish to examine situations where number of available runs at premium

Binary Response Simulation Procedure

- $d_{T,j}$ is calculated from result of previous trial or chosen as initial guess
 - Dosage step can be taken from either d_T or d_M of previous step
 - Method T (previous target): $d_{T,i+1} = d_{T,i} \pm \Delta d$
 - Method M (previous measured): $d_{T,i+1} = d_{M,i} \pm \Delta d$
- $d_{A,i+1}$ is randomly generated from $d_{T,i+1}$
 - Distribution: $N(d_{T,i+1}, \sigma_A^2)$
- $d_{M,i+1}$ is randomly generated from $d_{A,i+1}$
 - Distribution: $N(d_{A,i+1}, \sigma_M^2)$
- R_{i+1} generated from binomial distribution
 - Probability of response ($R_{i+1} = X$) equals the area under $N(D_{Tr}, \sigma_R^2)$ between the limits of $-\infty$ to $d_{A,i+1}$

Calculation of ED₅₀ for Individual Experiment

- A MLE solution for D_{Est} was obtained numerically using the Newton-Raphson method and a probit link function
 - A MINITAB® macro was written to perform the calculations
 - d_M values were used in the calculations for D_{Est}
 - As default for MLE calculations, σ_R was set equal to the Δd value used for the experiment
 - However, sometimes large values of σ_A and/or σ_M will produce situations where no solution was possible (using MINITAB®) unless σ_R was increased
 - Calculation artifacts would occur due to arbitrary cutoff of normit values by MINITAB®
 - Values below -7 were rounded up to -7, and values above 7 were rounded down to 7.
- Several other MLE related parameters were calculated
 - Wald Test statistic
 - Score Test statistic
 - Likelihood-Ratio Test statistic
 - Log Likelihood of final MLE fit



Factors Investigated

Factor	Description	Definition
A	Step Size	$(\Delta d / \sigma_R)$
B_A	Administration Error	(σ_A / σ_R)
B_M	Measurement Error	(σ_M / σ_R)
C	Location Initial Log Dosage	$([d_{T,1} \pm D_{Tr}] / \sigma_R)$
D	Basis for Calculation of Next Dosage	Method T: $d_{T,i+1} = d_{T,i} \pm \Delta d$ Method M: $d_{T,i+1} = d_{M,i} \pm \Delta d$



Factor Level Values

Factor	Level Values								
	1	2	3	4	5	6	7	8	9
A	0.500	0.707	1.000	1.414	2.000	2.828	4.000	5.657	
B _A	0.000	0.125	0.250	0.500	1.000	2.000			
B _M	0.000	0.125	0.250	0.500	1.000	2.000			
C	-3.000	-2.250	-1.500	-0.750	0.000	0.750	1.500	2.250	3.000
D	T	M							

Values in shaded boxes were not used with Method M

50 experiments per coordinate point

Method	T	M	Total
# Coordinates	2592	1440	4032
# Experiments	129600	72000	201600



- Efficiency
 - Number of experiments that “failed”
 - Experiments that consists of either all responses or all non-responses
 - Run of first reversal (RFR)
 - Indication of how quickly region of the median has been reached
 - Number of responses in an experiment
 - Indication of how well region of the median has been covered
- Precision
 - Standard error (SE) of D_{est}
 - Expressed as multiples of σ_R
 - Mean square error (MSE) of D_{est}
 - Expressed as multiples of σ_R^2

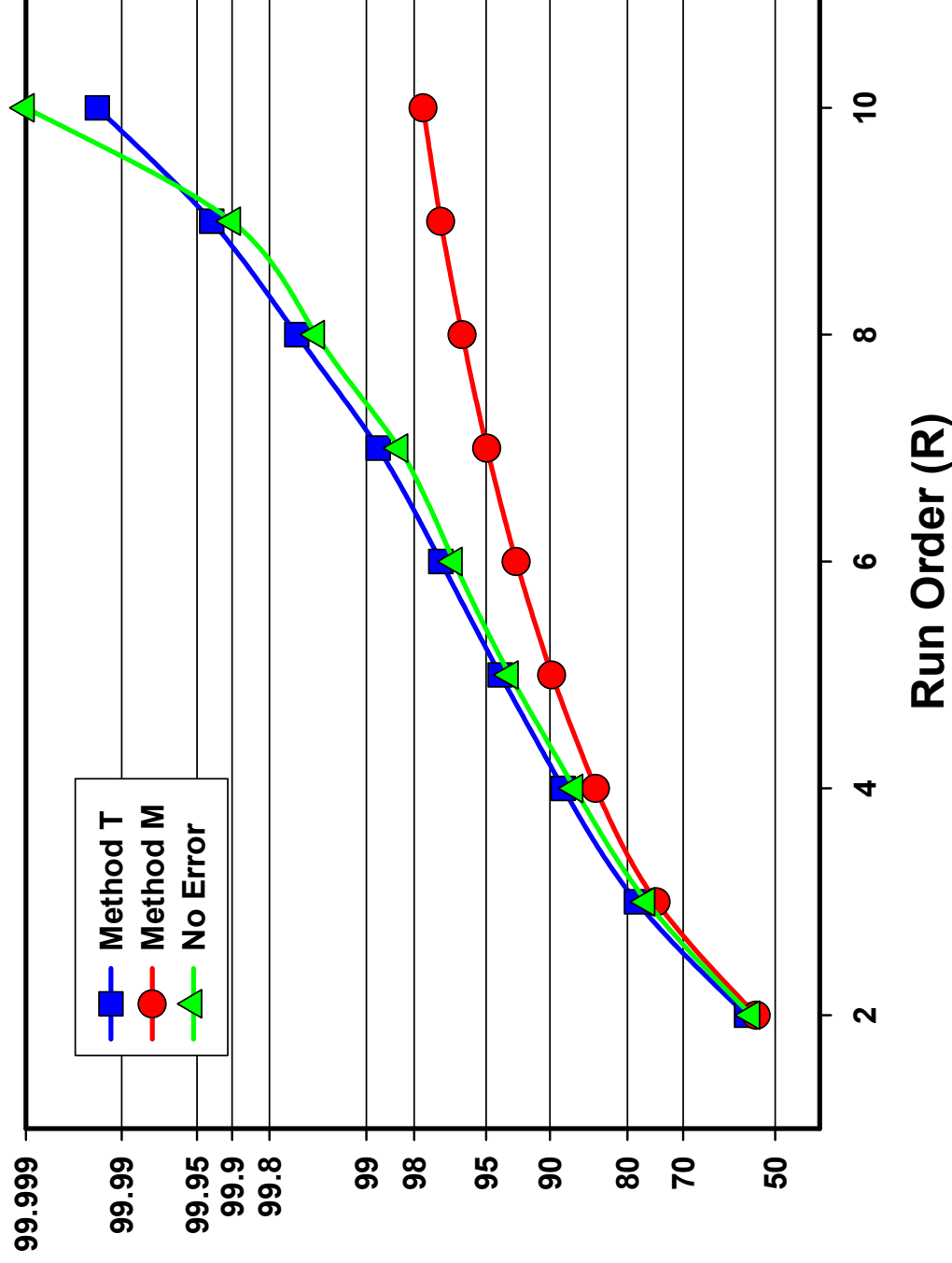


Efficiency Results

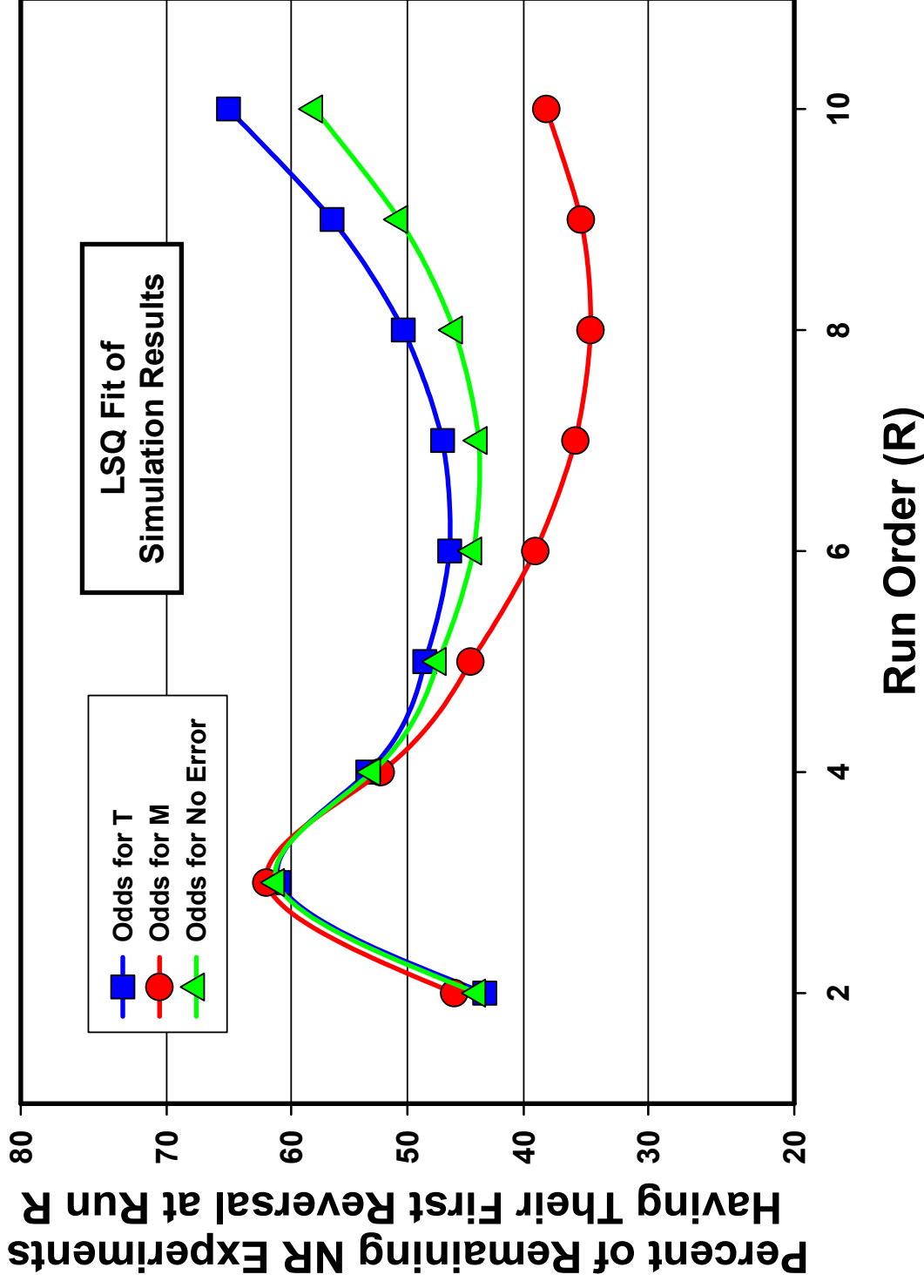
- Obtaining Early Run of First Reversal (RFR)
 - Factors in order of decreasing importance: A, D, CC, B_A , B_M and AA
 - Many statistically significant interactions also exist
 - Target dosage basis (D) more important than dosage errors (B_A and B_M)
- Success rate (T vs. M): Method T slightly better
 - Success rate of Method T closely parallels that of experiments executed with $B_A = B_M = 0$ as a function of Run of First Reversal (RFR)

Run of First Reversal Breakdown by Method and Run Order

Percent of Experiments with RFR At or Below R



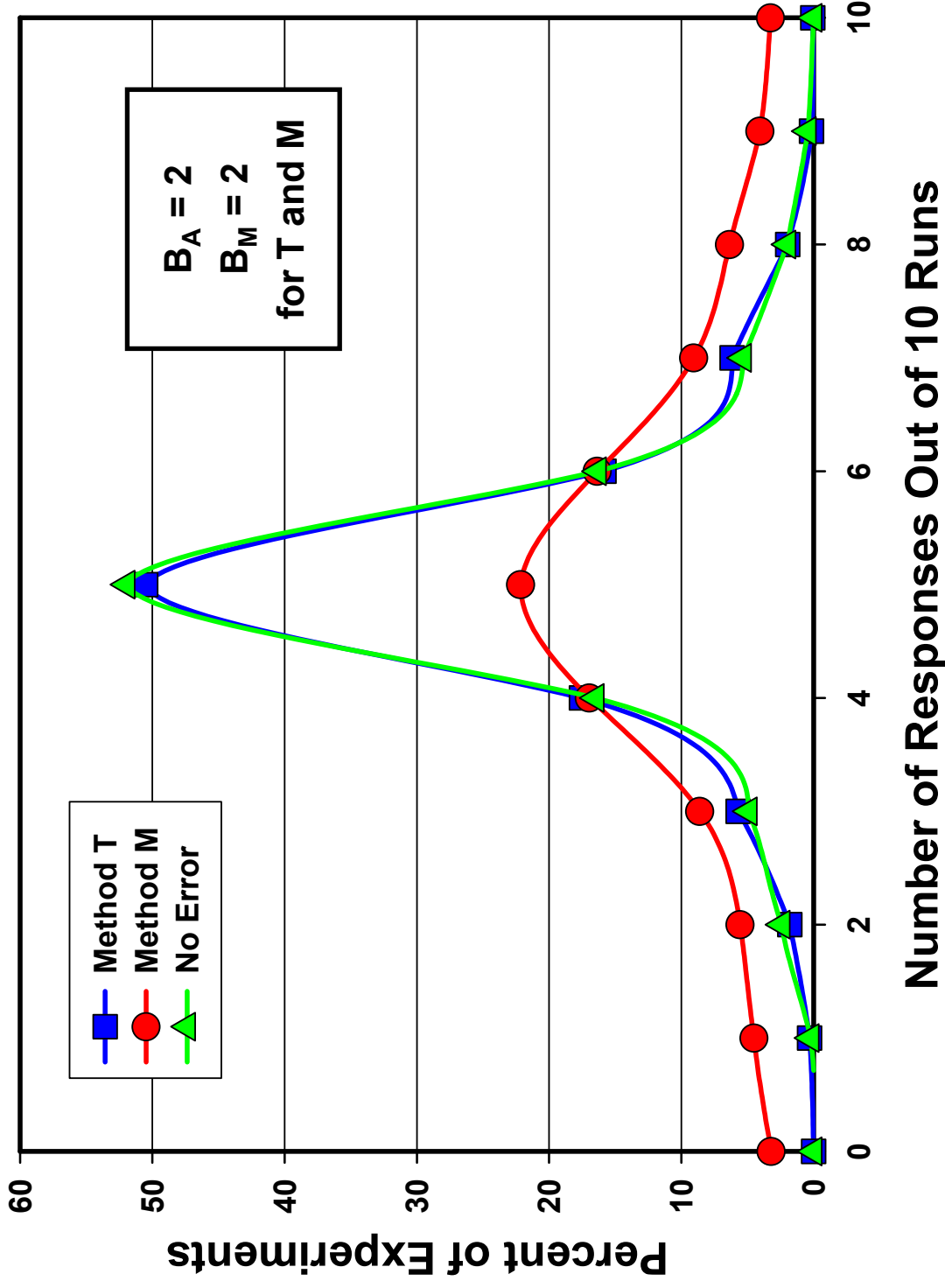
Chance of Non-Reversed (NR) Experiment having a Reversal



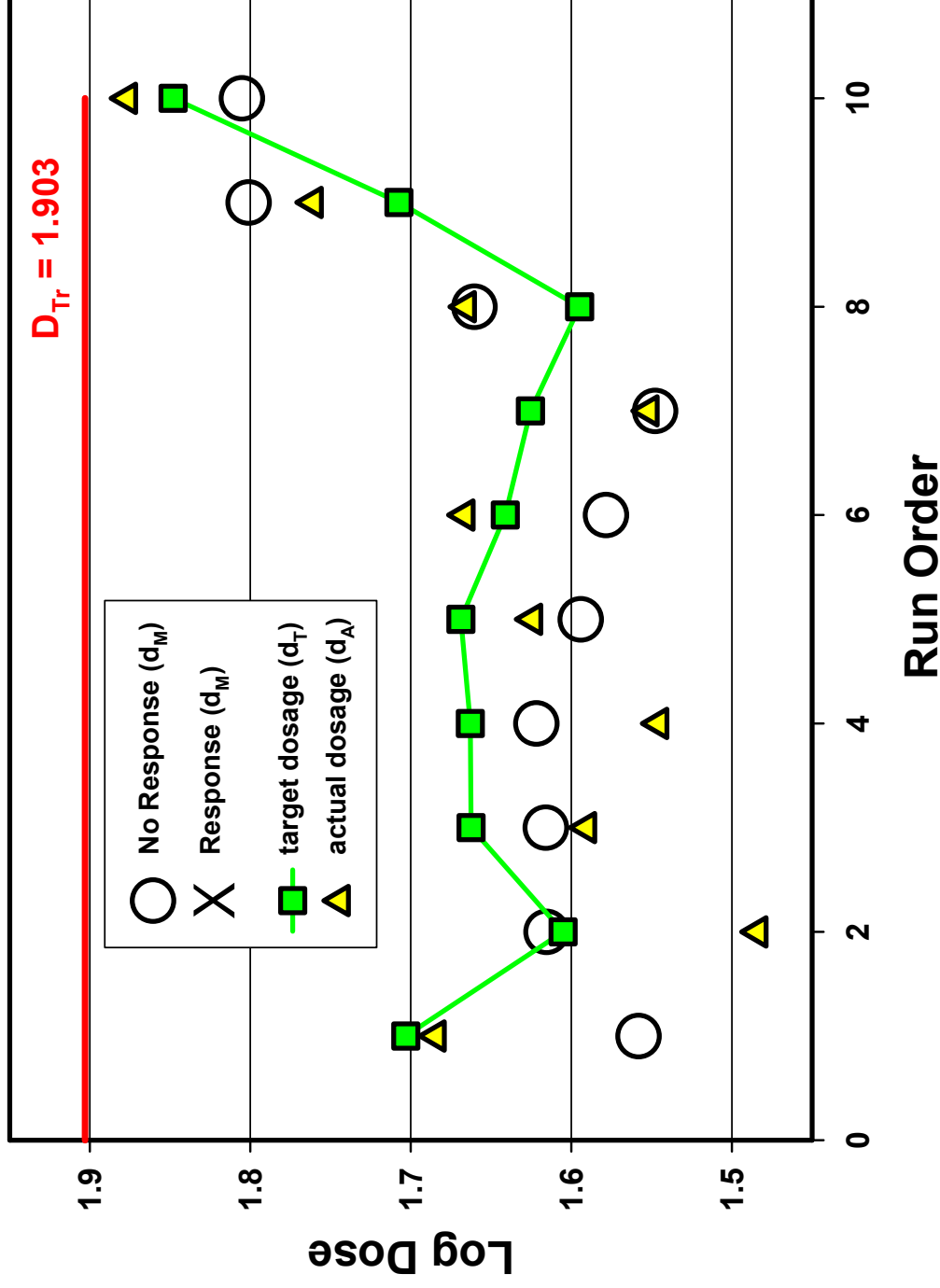
Efficiency Results (Cont.)

- Distribution of Number of Responses
 - Results for Method T (with any dosage errors present) mirrors that of UaD operated with no dosage errors
 - Distribution for Method M becomes shorter and broader as the dosage error increases
- Method T is more robust in recovering from “dead ends” than Method M
 - Target dosages in Method M are more readily influenced by fluctuations resulting from administration (σ_A) and measurement (σ_M) errors
 - Smaller step sizes can be easily overwhelmed
 - Advice on recommended step size ($\Delta d = \sigma_R$ or $A = 1$) may need to be revised

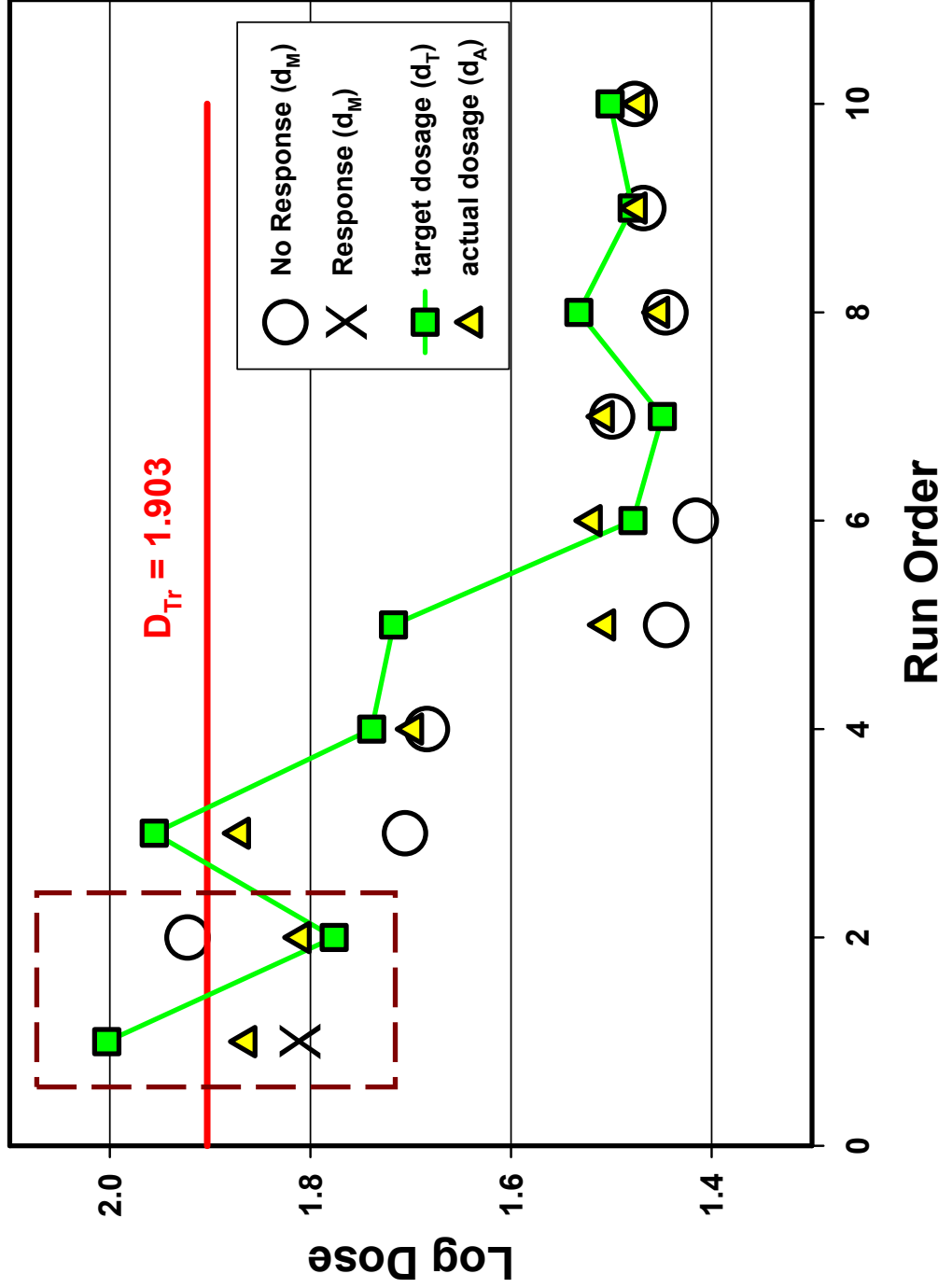
Distribution of Number of Response



Example of a Wandering Experiment with Method M



Example of a Negative Step Direction with Method M



Expt 28159M

Method M

$\Delta d = (1/30)$

$\sigma_R = (1/15)$

$A = 0.5$

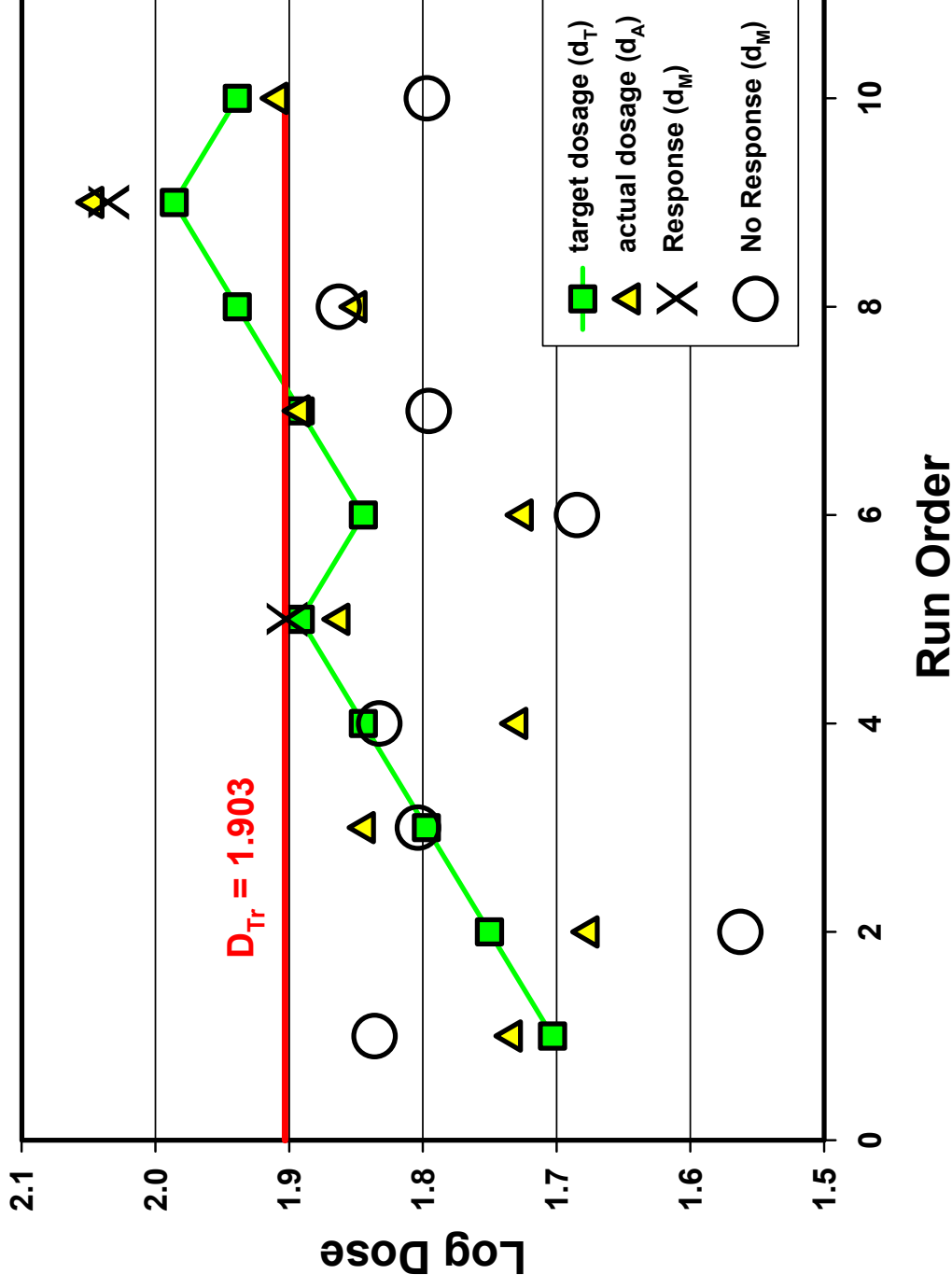
$B_A = 1$

$B_M = 1$

$C = 1$

$D_{Est} = 1.866$

Example of an Experiment with Method T



Expt 64013T

Method T

$\Delta d = (0.047)$

$\sigma_R = (1/15)$

$A = 0.707$

$B_A = 1$

$B_M = 1$

$C = -3$

$D_{Est} = 1.901$

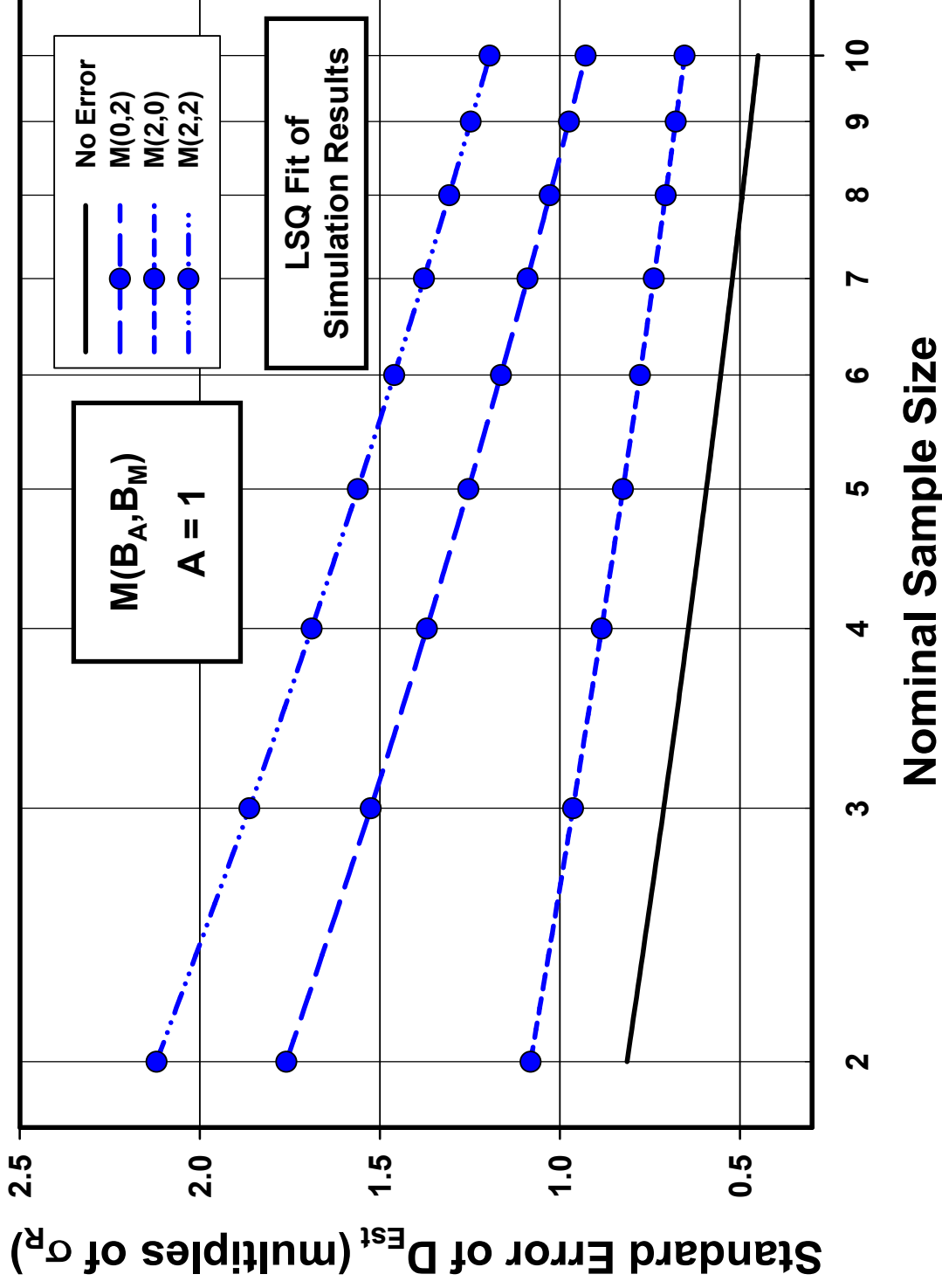
Run Order



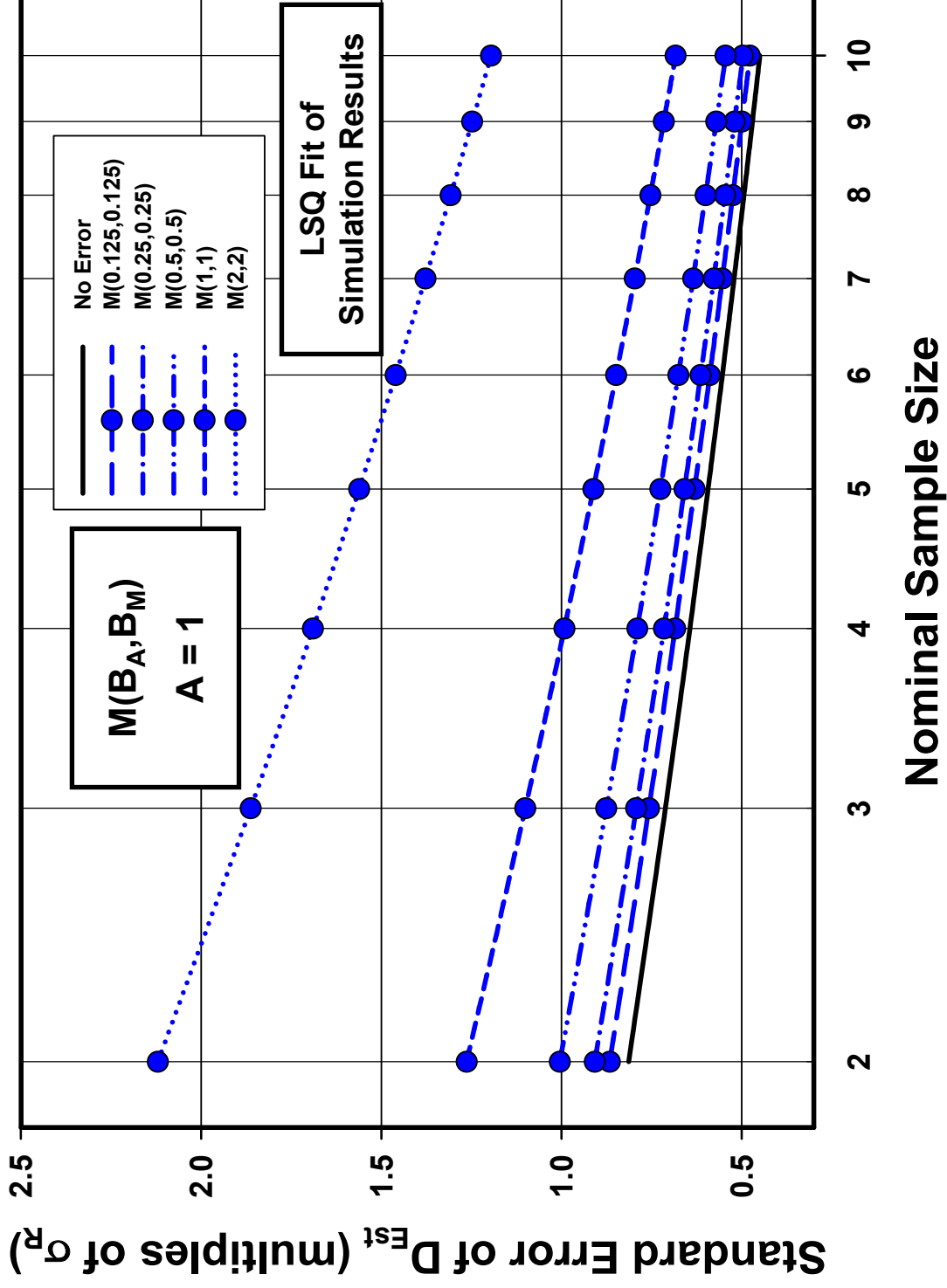
Precision Results

- Additional factor introduced--N (nominal sample size)
 - Accomplished by dividing up the 50 runs per coordinate point into five equal size groups having N_T values of 6, 7, 8, 9 and 10, respectively
 - Knowing RFR for each experiment, N was calculated from N_T and N_{FR}
 - Values of N from 2 to 10 obtained
- Factors in order of decreasing importance: $\log N$, B_M , A , B_A , $B_M B_M$ and $B_A B_A$
 - Many statistically significant interactions also exist
 - D is only significant when involved in an interaction (ex. AD and $B_M D$)
 - Dosage errors (B_A and B_M) more important than agent dosage basis (D)
 - This is in contrast to the reverse being true for experiment success rate
- Standard error for D_{Est} is smaller using Method T than with Method M

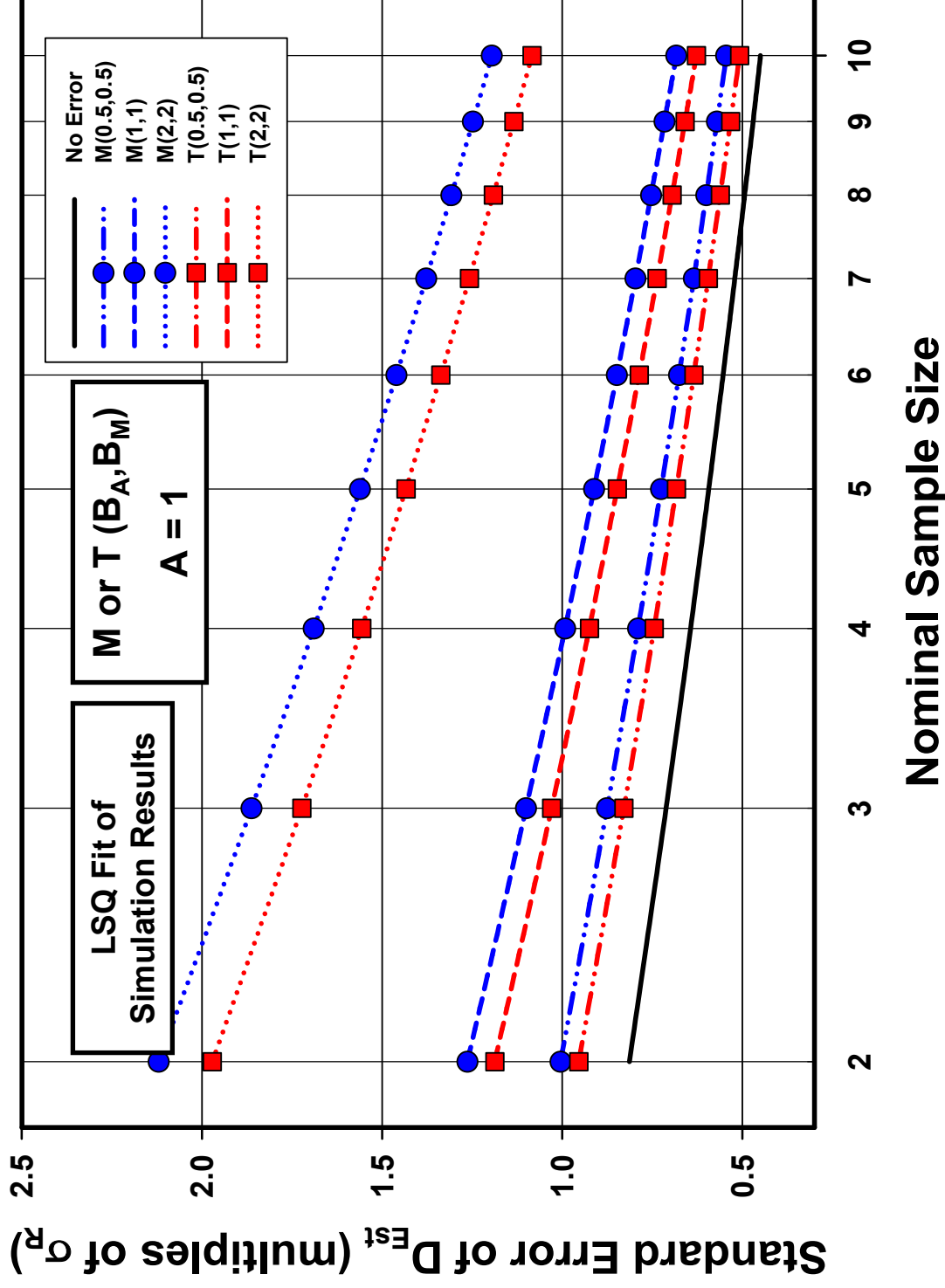
Standard Error of D_{Est} for Method M as a Function of N , B_A and B_M



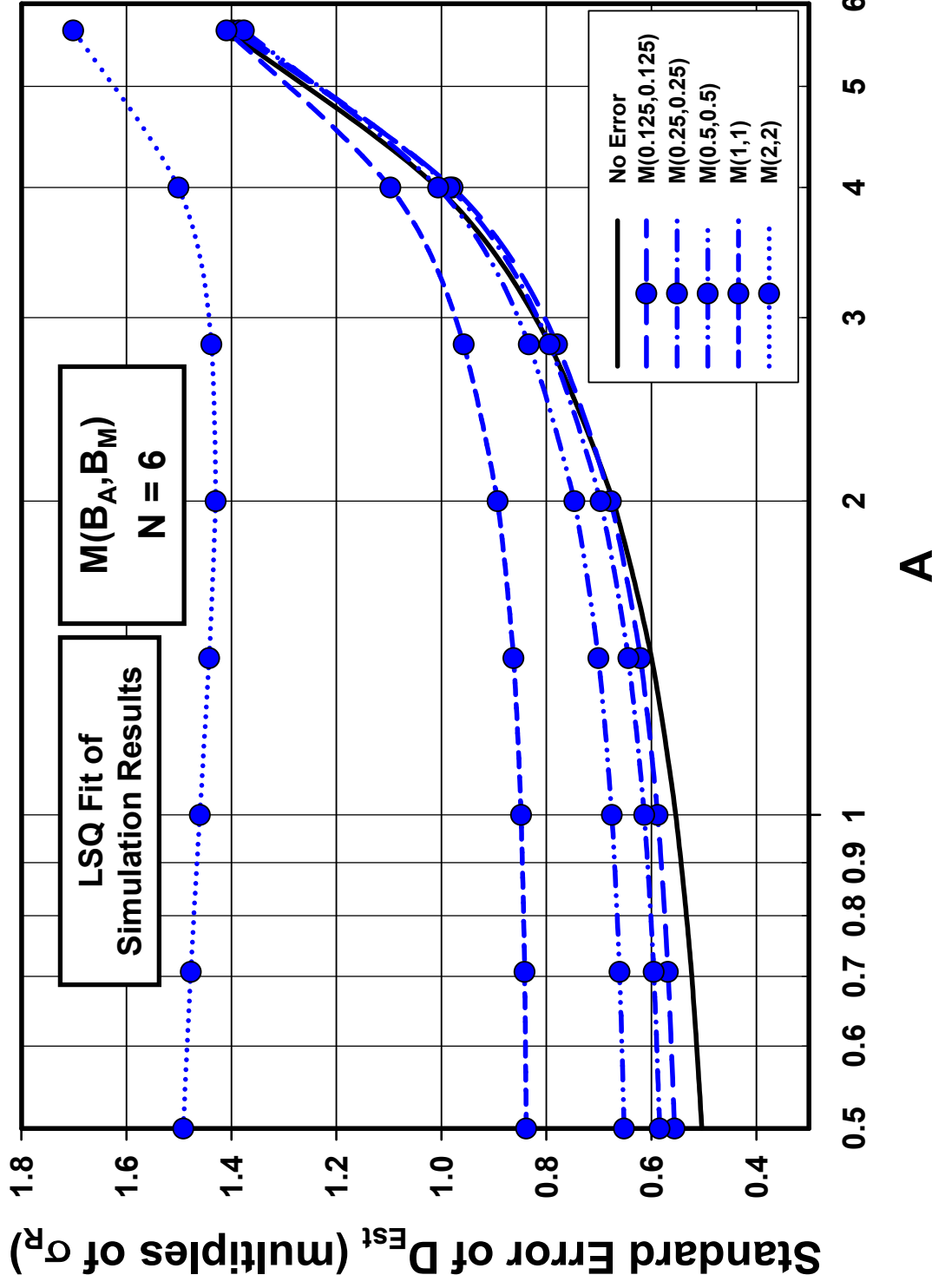
Standard Error of D_{Est} for Method M as a Function of N and $(B_A + B_M)$



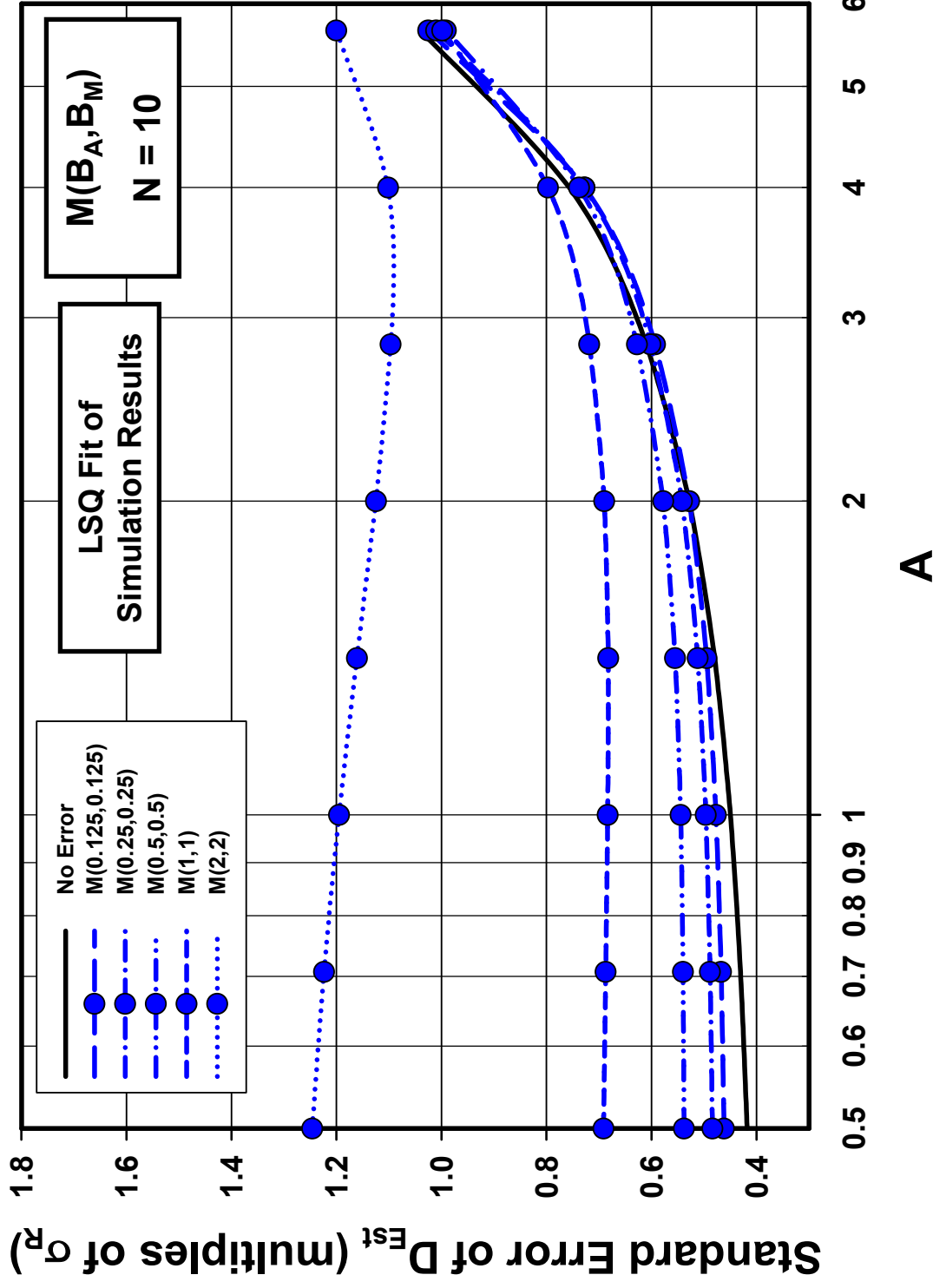
Standard Error of D_{Est} for Methods M and T as a Function of N and $(B_A + B_M)$



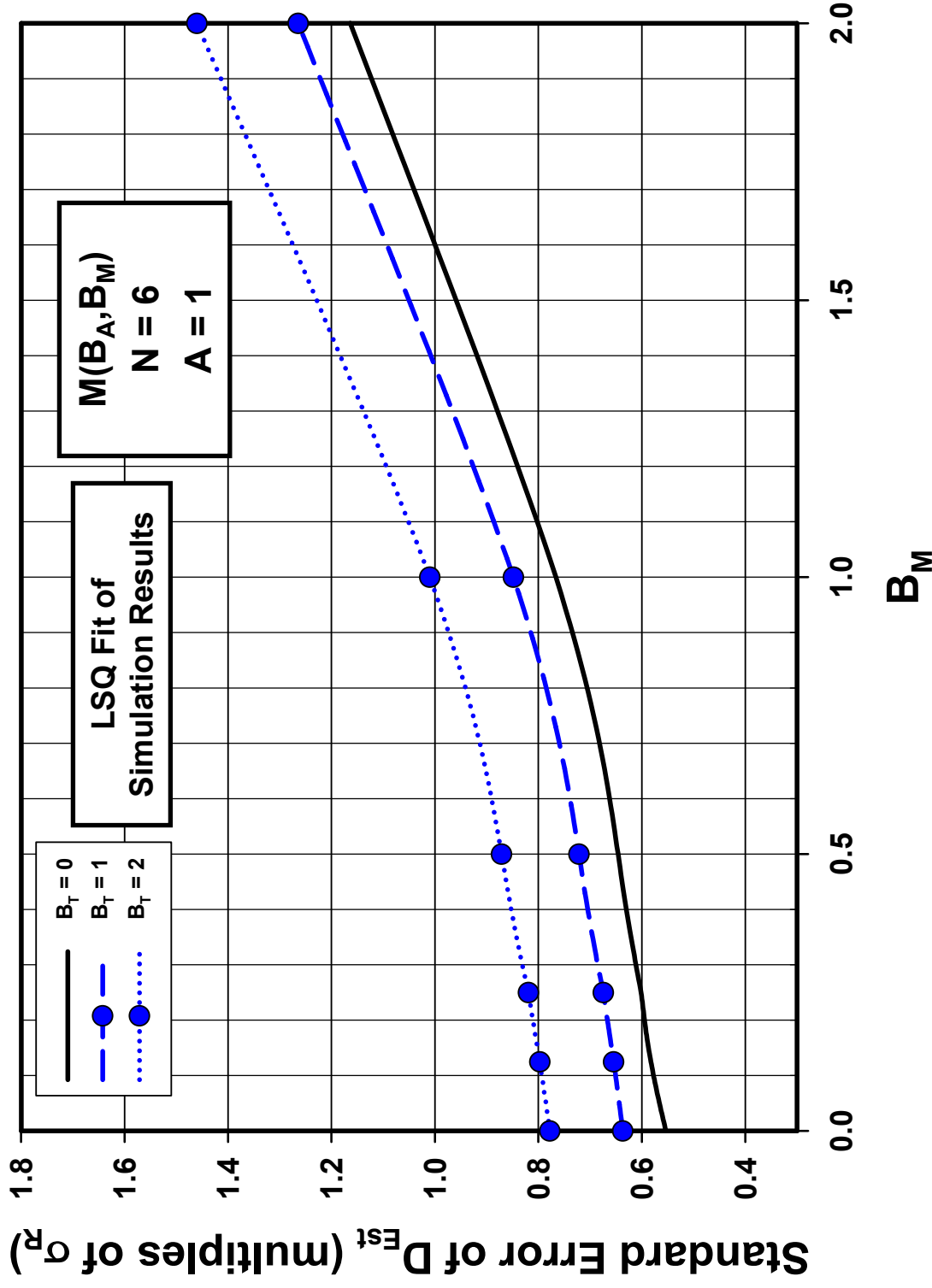
Standard Error of D_{Est} for Method M as a Function of A and $(B_A + B_M)$



Standard Error of D_{Est} for Method M as a Function of A and $(B_A + B_M)$



Standard Error of D_{Est} for Method M as a Function of B_A and B_M





Summary

- With non-zero dosage errors (B_A & B_M), UaD Method is robust with respect to efficiency & precision in the following parameter space
 - $-3 \leq C$ (initial dosage) ≤ 3
 - Neither B_A and/or B_M exceeds roughly 0.25 to 0.5
- When selecting next target dosage (Factor D), Method T is better than Method M overall
 - Difference between methods becomes more pronounced as:
 - A decreases
 - B_A and/or B_M increases
 - $|C|$ increases
 - Method T is more able to resist effects of large values of B_A & B_M



Summary (Cont.)

- Method T produces a distribution of number of responses and Run of First Reversal that is essentially identical to what is produced by traditional UaD (no dosage errors)
- Method T produces more precise D_{Est} values than Method M
 - However, dosage errors (B_A & B_M) have larger impact on precision than Factor D (Methods T and M) (the reverse is true for efficiency)



Conclusions

- Dosage administration and measurements errors do have an effect on the efficiency and precision of the Up-and-Down Method
- The effect of dosage errors on efficiency is significant overall
 - Can be practically eliminated by using Method T (basing next step on the previous target dosage) instead of Method M (basing next step on the previous measured dosage)
- The greater precision afforded by smaller step sizes is increasingly eroded by increasing dosage errors
 - Dosage measurement error has greater impact than both step size and dosage administration error on precision

Special Session
on
**Robust and Resilient Critical
Infrastructure Systems**

Overview, Problem Description, and Challenges

Jagdish Chandra
The Institute for Reliability and Risk Analysis
The George Washington University
October 30, 2003

Some Basic Definitions

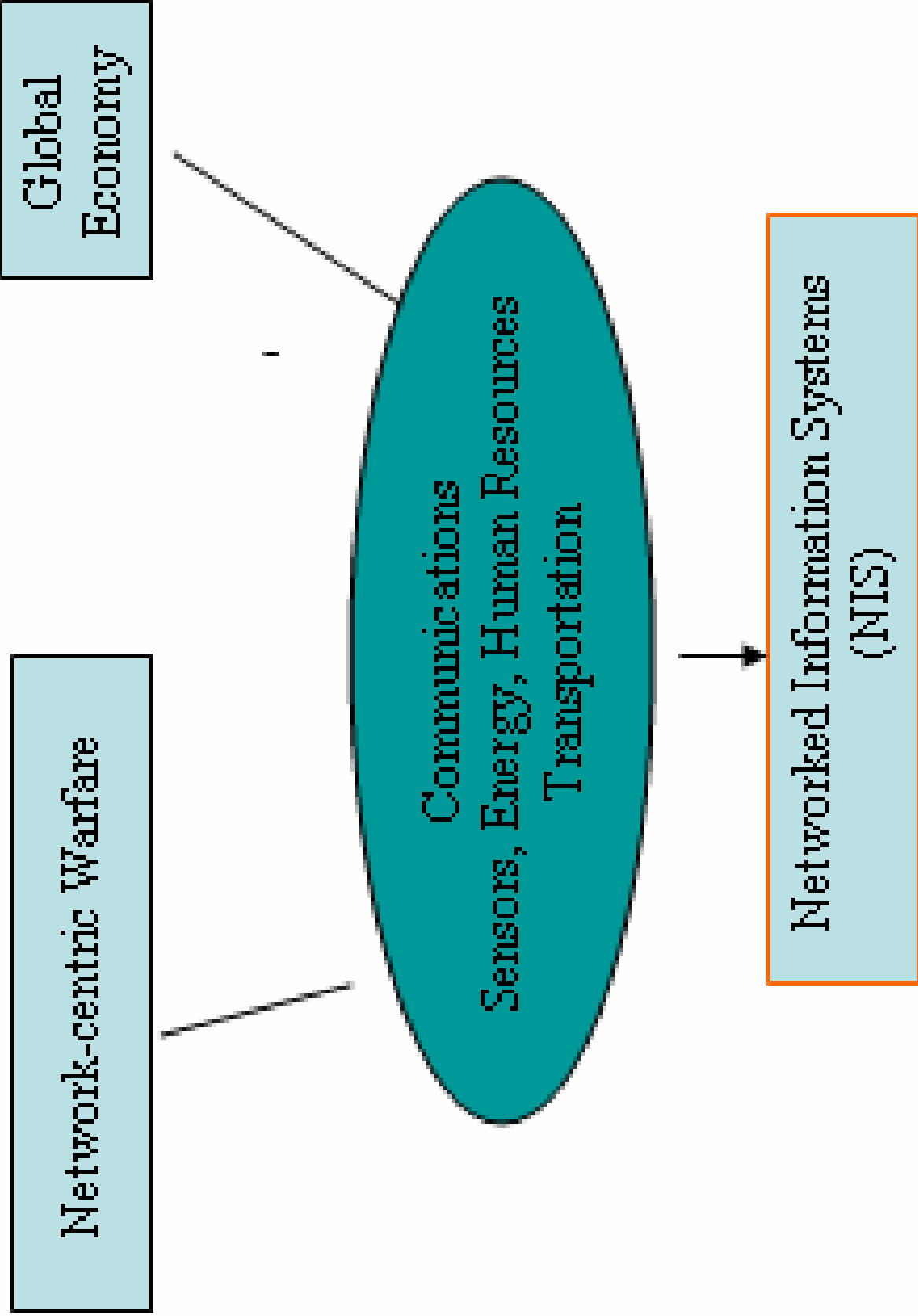
- **Infrastructures:** Linked system of facilities and activities that provide the range of essential services
- **Critical Infrastructures:** So vital that their incapacitation or destruction would have a debilitating impact on defense or national security (Clinton, 1997)
- **Robustness:** Failure- resistant through design and /or construction
- **Resilience:** Ability to recover quickly

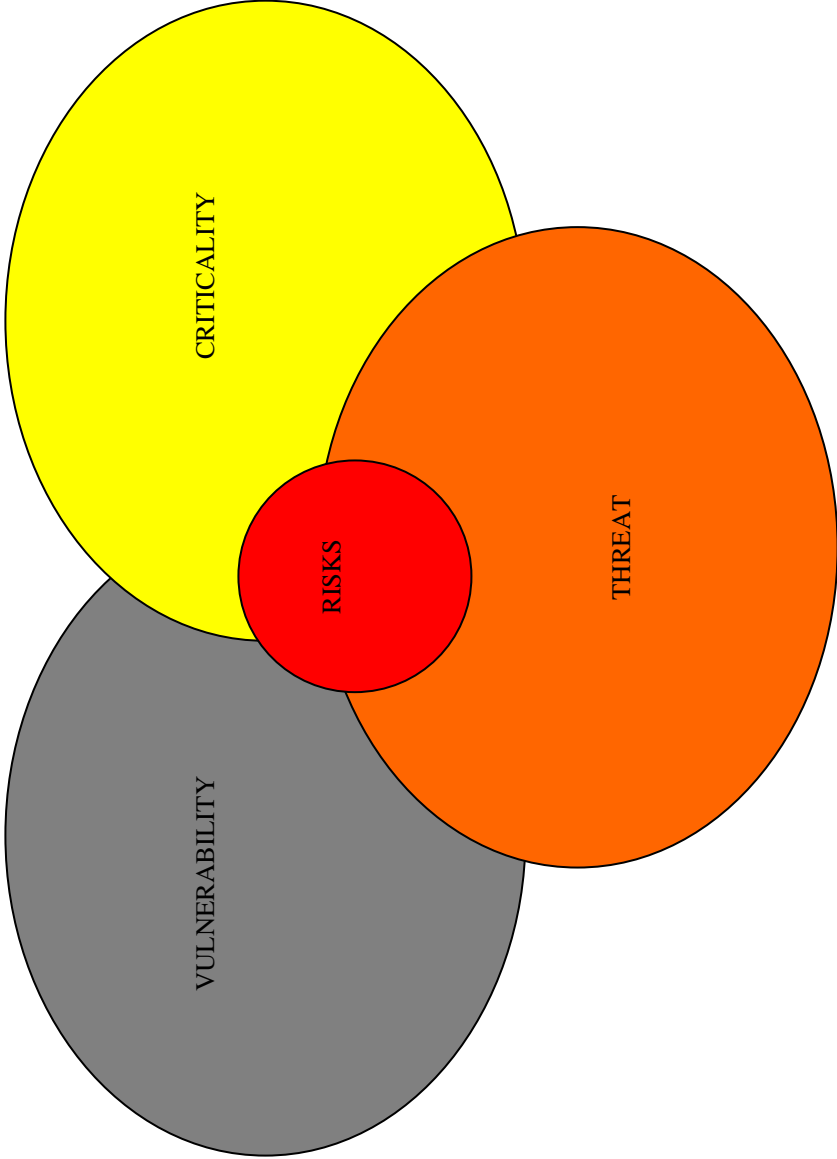
INTERDEPENDENCIES

- Physical Interdependency
 - Cyber Interdependency
 - Geographic Interdependency
 - Logical Interdependency
- ***Modeling and Simulation of Interdependent Infrastructures is a complex, multifaceted, and multi-disciplinary problem***

Factors: Analyses of Interdependencies

- Time Scales
- Geographic/Spatial Scales
- Higher Order Effects/ Cascading
- Human/Social/Psychological
- Operational Procedures
- Business Policies/Government Regulations
- Restoration/Recovery Procedures
- Stakeholders Concern





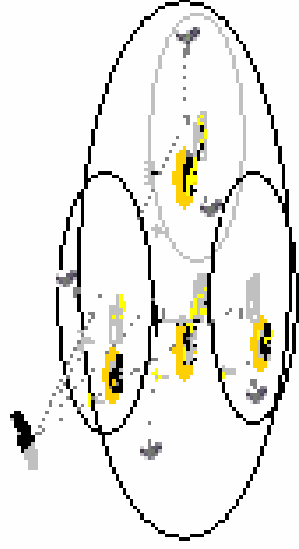
INFRASTRUCTURE VULNERABILITIES

- Natural Hazards
 - Degradation of Material and Components
 - Complexity/ Interconnections
 - Malevolent Acts
- ***Characterize Vulnerabilities, Role of Dependencies, and Propagation of Failures***

THREAT ANALYSIS

- Analyze data, patterns of threat scenarios and intrusion (some data-centric tools combining data-mining and adaptive case-based reasoning have been developed at FSU)
- Information fusion and management (Bayesian techniques for both cooperative and adversarial/compromised sensors/sources-GWU)
- A stochastic framework for intrusion detection (optimal filtering techniques for IDS- UW)

Fault Tolerance and Recovery in Mobile Wireless Networks



- Hybrid totally wireless networks
- Standby (backup) mobile routers are used to provide recovery and enhance reliability
- Developed distributed recovery protocols: the backup routers are scattered among the primary mobile routers (UCF)
- Using a flock-like dynamics, studying the arrangement of backups to maximize reliability (GWD)

RESILIENT INFRASTRUCTURES

- Reliable communication in a dynamic battlefield (developed fault-tolerant and distributed recovery protocols for hybrid totally mobile wireless networks-UCF)
- Robustness and Resiliency (analysis and design of strategies for optimal deployment of back-up mobile routers-GWU)

Risk Assessment and Management

- What can go wrong? What is the likelihood that it will? And, what are the consequences
- What can be done and what options we have? What are the trade-offs in terms of costs, benefits, and risks? And, how these decisions impact the future?
- ***Characterize optimal defensive strategies for sabotage risks (e.g., game theory as a paradigm for critical infrastructure protection-UW)***
- ***Risk management strategies for high consequence/low probability events***

Information Assurance of NIS

- Human introduced errors
- System probing (malicious, non-malicious)
- System penetration
- Subversion of networks
- Devise security and control mechanisms
- Misuse of policy, authority, power
- ***The interface between technology and human behavior; human factor is the Achilles heel of information security***

Networked Systems Simulation

- Modeling and simulation for resiliency designs (developed distributed modular intrusion detection system for ad hoc and hybrid totally-mobile wireless networks- UCF)
- Complex systems simulation (optimizing performance in networked systems- UW)

Risk Assessment: A Game Theoretic Approach

**Vicki Bier, Aniruddha Nagaraj,
Vinod Abhichandani**

University of Wisconsin-Madison

Background



- ◆ **Game theory is a useful model for security risk assessment:**
 - ◆ Appropriate when protecting against intelligent and adaptable adversaries
 - ◆ Recognizes that defensive strategies must take attacker behavior into account
 - ◆ Can identify qualitative properties of optimal solutions (e.g., randomization)

Background...

- ◆ **Game theory is only beginning to be used in security risk assessment**
- ◆ **Military analogies (Schneier):**
 - ◆ **“The defender has to defend against every possible attack”**
 - ◆ **“The attacker...only has to choose one attack, and he can concentrate his forces on that one attack”**

Background...



- ◆ **Most applications are still exploratory:**
 - ◆ Illustrative applications to the choice of attack and defense strategies (Cohen)
 - ◆ Experiments demonstrating relevance of game theory to information warfare (Burke)
 - ◆ Application of game theory to financial institution risks (Chaturvedi et al., Gupta)
 - ◆ Importance of perverse incentives (Anderson)

Outline of this work

- ◆ **Games between attackers and defenders:**
 - ◆ Simple series/parallel systems
 - ◆ Components with inherent values, and also a value to system function

Overall goal

- ◆ Study optimal allocation of resources for protection of series and parallel systems against intentional attacks
- ◆ Protective investment c_i reduces the probability of successful attack against component i to $p_i(c_i)$:
 - ◆ $p_i(c_i)$ convex, decreasing, twice differentiable and invertible

Cases being considered

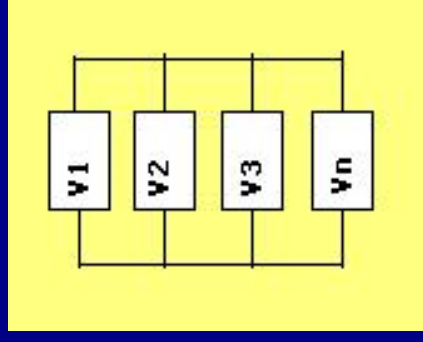
- ◆ **Results to date:**
 - ◆ Components in parallel
 - ◆ Components in series
 - ◆ “Additive” models
(Components have different “values” v_i)
- ◆ **In process:**
 - ◆ Arbitrary series/parallel structures
(NP-hard, may use heuristic approaches)
 - ◆ Other configurations
(Explore merits of perimeter defense, etc.)

Assumptions

- ◆ Realistic levels of defensive investment will not deter attacks:
 - ◆ Models applicable to determined attackers
- ◆ Attacks against different components succeed or fail independently:
 - ◆ Models applicable to functionally diverse and spatially separated defenses
- ◆ Likely to apply to most serious threats against security-critical systems

Components in parallel

- ◆ Defender wishes to maximize (expected value of system) – (defense cost), or equivalently:
 - ◆ Choose c_i to minimize $\alpha [\prod p_i(c_i)v + \sum p_i(c_i)v_i] + \sum c_i$
where α is probability of an attack on the system,
 v is the value of the system functionality, and
 v_i is the inherent value of component i
- ◆ Optimum occurs when
 - ◆ $p_i'(c_i) \geq -1/\alpha[v + v \prod_{j \neq i} p_j(c_j)]$, and
 - ◆ $c_i [\alpha v_i p_i'(c_i) + \alpha v p_i''(c_i) \prod_{j \neq i} p_j(c_j) + 1] = 0$
- ◆ Multiple local optima are possible



Components in parallel...

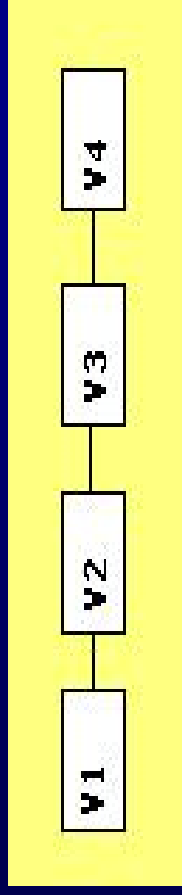
- ◆ Local optimum is unique when $p_i(c_i)$ are **log convex**:
 - ◆ Success probability decreases “faster than exponentially” in c_i
- ◆ This seems unlikely to be the case

Components in parallel...

- ◆ **General insights:**
 - ◆ **Optimal defense strategy depends on the cost-effectiveness with which components can be improved**
 - ◆ As measured by the $p_i'(c_i)$ and on the **values of the components**
 - ◆ As measured by the v_i
 - ◆ **Components that are too costly to defend (relative to their value) will not be hardened**

Components in series

- ◆ Can occur for many reasons:
- ◆ Physically in series (e.g., pipelines)
- ◆ Multiple failure modes
- ◆ Attacker can afford only one target
- ◆ First successful attack is much more serious (e.g., for symbolic reasons)



Components in series...

- ◆ **Attacker has a choice of targets**
- ◆ **Two bounding cases:**
 - ◆ **Attacker has no knowledge of defensive investments**
 - ◆ **Attacker can obtain perfect knowledge about defensive investments at no cost**

Components in series...

Perfect knowledge

- ◆ Assumption of perfect knowledge may not always be unrealistic:
 - ◆ Due to the openness of our society
- ◆ Public demands knowledge of defense
 - ◆ Even when this weakens its effectiveness!
- ◆ This increases the difficulty of defense:
 - ◆ E.g., anthrax protection

Components in series...

Perfect knowledge

- ◆ Assume attacker has only one attempt (multiple attacks are considered later)
- ◆ Attacker objective is to:
 - ◆ Choose i to maximize $[p_i(c_i) v_i]$
- ◆ For optimal allocation of defensive resources:
 - ◆ Defense must equalize the expected values of attacks against all targets
 - ◆ “Each of the defended targets [must] yield the same payoff to the attacker” (Dresher)

Components in series...

Perfect knowledge

- ◆ Unlike in defending against accidents or acts of nature:
 - ◆ **Optimal allocation does not depend on cost-effectiveness of investments!**
- ◆ **Defender is deprived of flexibility:**
 - ◆ **Must defend all targets of comparable expected value equally (regardless of cost)**

Insight

- ◆ **“Investment in defensive measures,**
- ◆ **unlike investment in safety measures,**
- saves a lower number of lives.” (Ravid)**

Components in series...

Perfect knowledge

- ◆ Now, assume that the attacker can attack **each component** once (multiple attacks)
- ◆ Attacker objective is to:
 - ◆ Choose i to maximize $\alpha [\sum p_i(c_i) v_i + v \{1 - \prod [1 - p_i(c_i)]\}] + \sum c_i$
- ◆ For optimal allocation of defensive resources:
 - ◆ Defense need not focus exclusively on components that cause highest expected damage
 - ◆ Investment in other components may pay off, if attacks against such “first-choice” targets fail
 - ◆ Optimal defense strategy again depends on the cost-effectiveness with which components can be improved

Components in series...

Perfect knowledge

- ◆ **Insights:**
 - ◆ Properties of the optimal solution for series systems with multiple attacks are similar to those for parallel systems (e.g., multiple optima)
 - ◆ If one component dominates the risk, then the optimal solution with multiple attacks will be similar to that with a single attack

Components in series...

No knowledge

- ◆ Assume:
 - ◆ Attacker targets component i with constant probability q_i (regardless of defense c_i)
 - ◆ Attacker has only one attempt
- ◆ Defender objective similar to previous:
 - ◆ Choose $\{c_i\}$ to minimize $\sum q_i v_i p_i(c_i) + \sum c_i$
- ◆ Optimum occurs when $p_i'(c_i) \geq -1/(q_i v_i)$
 - ◆ and $c_i [q_i v_i p_i'(c_i) + 1] = 0$
- ◆ Expenditure c_i is increasing in $q_i v_i$

Arbitrary Structures

- ◆ Find defensive strategy when optimal attack strategy is NP-hard (joint work with Cox, Azaiez):
 - ◆ Cox's work on least cost diagnosis (1989, 1996) suggests near-optimal heuristic attack strategies
 - ◆ Identify optimal (or near-optimal) defenses against near-optimal attacks
 - ◆ Determine when heuristic attack strategies are in fact optimal

Conclusions

- ◆ Protection of series systems from knowledgeable adversaries is a **fundamentally different** challenge:
- ◆ Investments less cost-effective (since attacks can be deflected to other targets)
- ◆ Defender loses flexibility to allocate resources cost-effectively
- ◆ Importance of redundancy, secrecy (and deception) as defensive strategies

Conclusions...

- ◆ Defender should consider the **success probabilities** of attacks against various components:
 - ◆ Not only their inherent values
- ◆ Some high-value targets with a low probability of being successfully attacked may not merit any investment:
 - ◆ Lower-value, more vulnerable targets may merit defense
- ◆ Contrast this to the heuristic proposed by Brookings (2002):
 - ◆ Protecting only the most valuable assets



COLLEGE OF ENGINEERING
UNIVERSITY OF WISCONSIN-MADISON

Optimizing Performance in Networked Systems

Julien Granger (UW-Madison)
Ananth Krishnamurthy (RPI)
Stephen M. Robinson (UW-Madison)

Work Sponsored by ARO and AFOSR

U.S. Army Conference on Applied Statistics
Napa, CA, 30 October 2003

Outline

- Problem of improving performance in a network with uncertainty
- Difficulty with repeated simulation
- Using approximation methods to avoid repeated simulation
- Example: Analytical Airfield Model (AAM)
 - Small numerical example that nonetheless gives significant managerial insight
- Conclusion

The Problem

- Given a network with uncertainty, decide how to invest resources to improve its performance
 - These often arise in logistics, but occur also in many other areas
- Performance measure: often throughput
- Investment: change system parameters to
 - Decrease times for critical operations such as transshipment, refueling, maintenance
 - Increase capacity of intermediate stations that perform these operations

Measuring Performance

- We usually have to measure performance with simulation models
 - Significant uncertainties in these systems
 - Deterministic models are easy to use but they ignore the uncertainties and may give seriously wrong answers
- Simulation models give “snapshots” of performance of a given system, but are slow
 - More so, if system is large and complex
 - And we need to simulate *every time we change parameter values*, to check improvement

Problem: Slowness

- These methods simulate the system at each step
 - On complex systems they can be very slow
- How can we get around this?
 - One way: speed up the simulations
 - » Network of workstations (Condor Project)
 - Another: maybe avoid simulating each time
 - » Use approximation methods instead

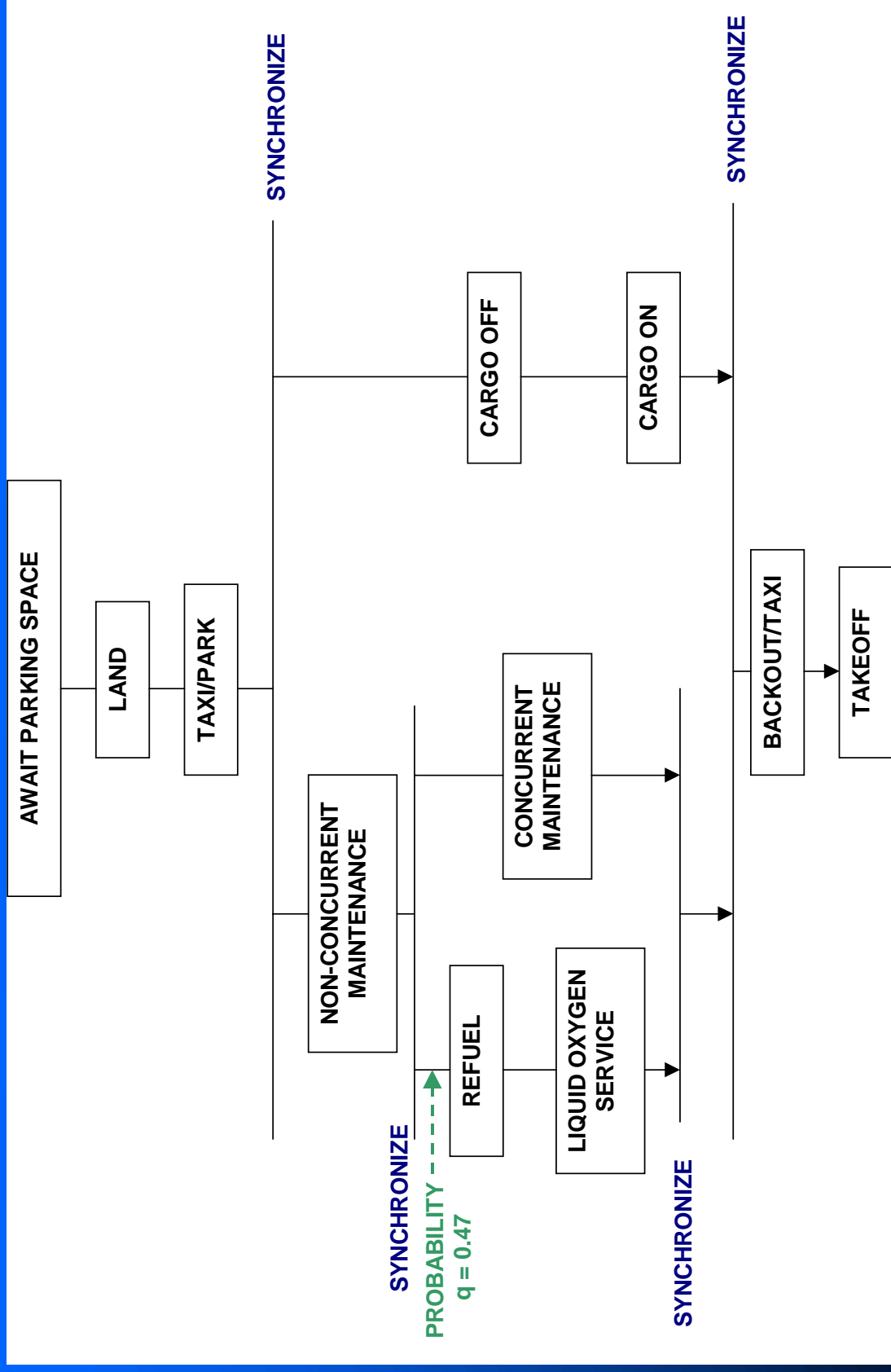
Network Approximations

- Idea: don't simulate the system exactly
 - Approximation requires *only solving a system of nonlinear equations*
 - » Much faster than simulating
 - On various problems we've tried, ratio of simulation time to approximation time has been in range 16:1 to 80:1 or more
- So: use approximations to improve the system, then simulate once to confirm the results

Example

- Analytical airfield model (AAM)
 - Ref: Dietz, D. C., *J. Aircraft* (1999) *et seq.*
- Relatively complex network of immediate military interest
 - We'll illustrate how to use optimization to answer managerial questions
 - What's new here is better control methods through use of new 2-moment approximations
 - Tools apply just as well to many other networks

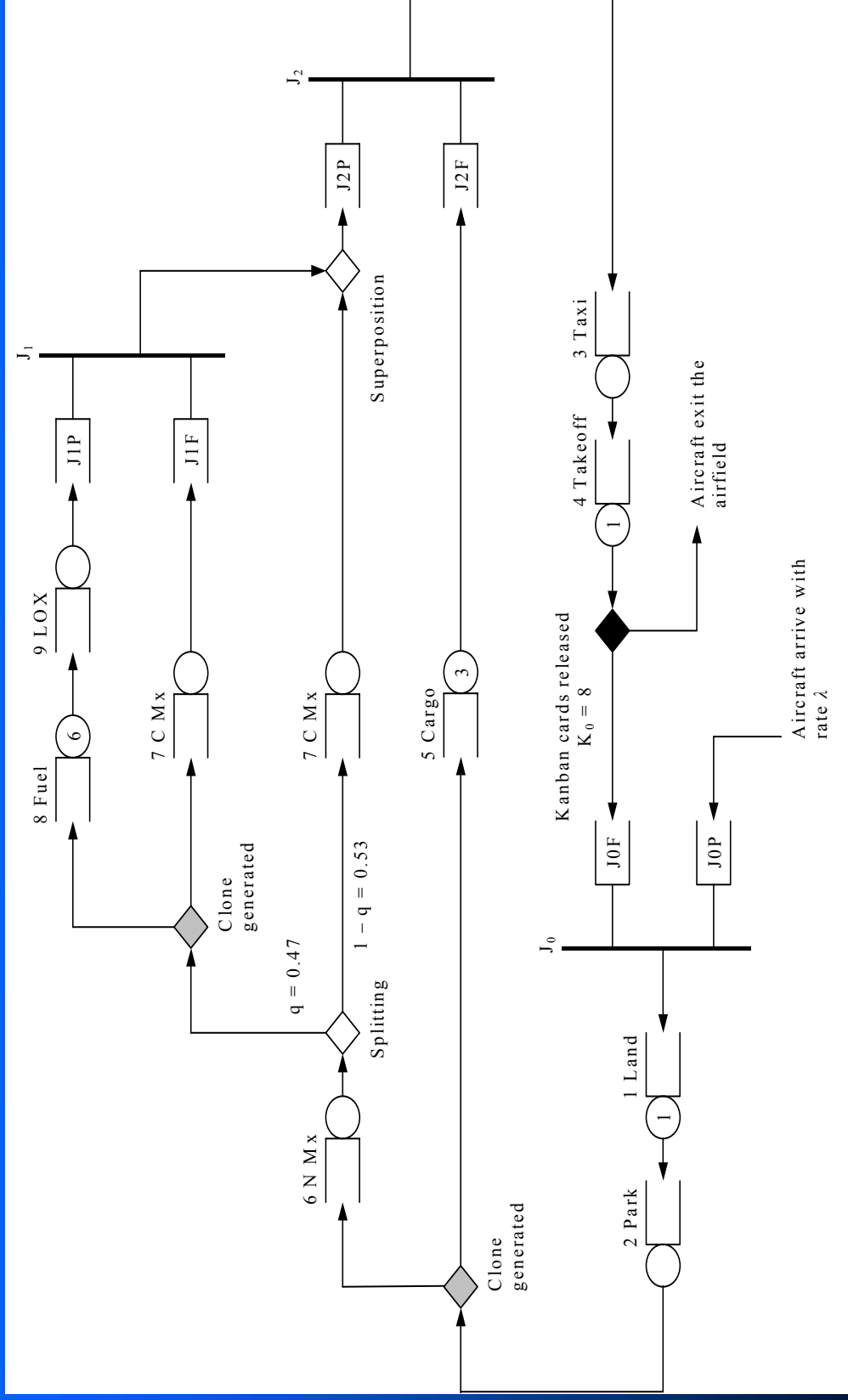
Task Precedence Graph for the AAM Problem



Modeling Essential Features

- Concurrent use of resources can be modeled as a fork/join mechanism:
 - Aircraft generate temporary clones with identical attributes
 - Overall waiting time is the maximum over all clone paths
- Aircraft waiting authorization for landing obey a join mechanism governed by Kanban cards
 - Number of Kanban cards equals maximum number of aircraft allowed on ground simultaneously

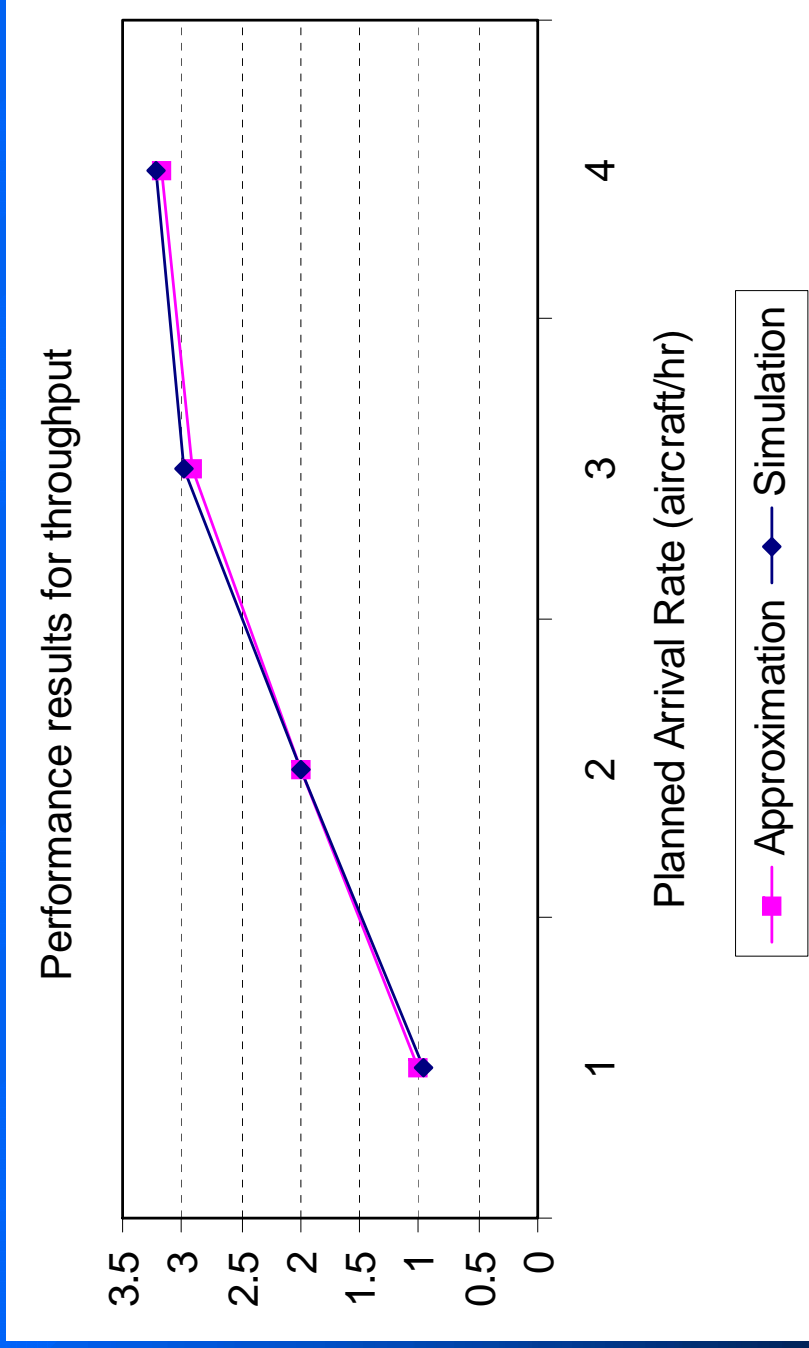
Layout of Queuing Model



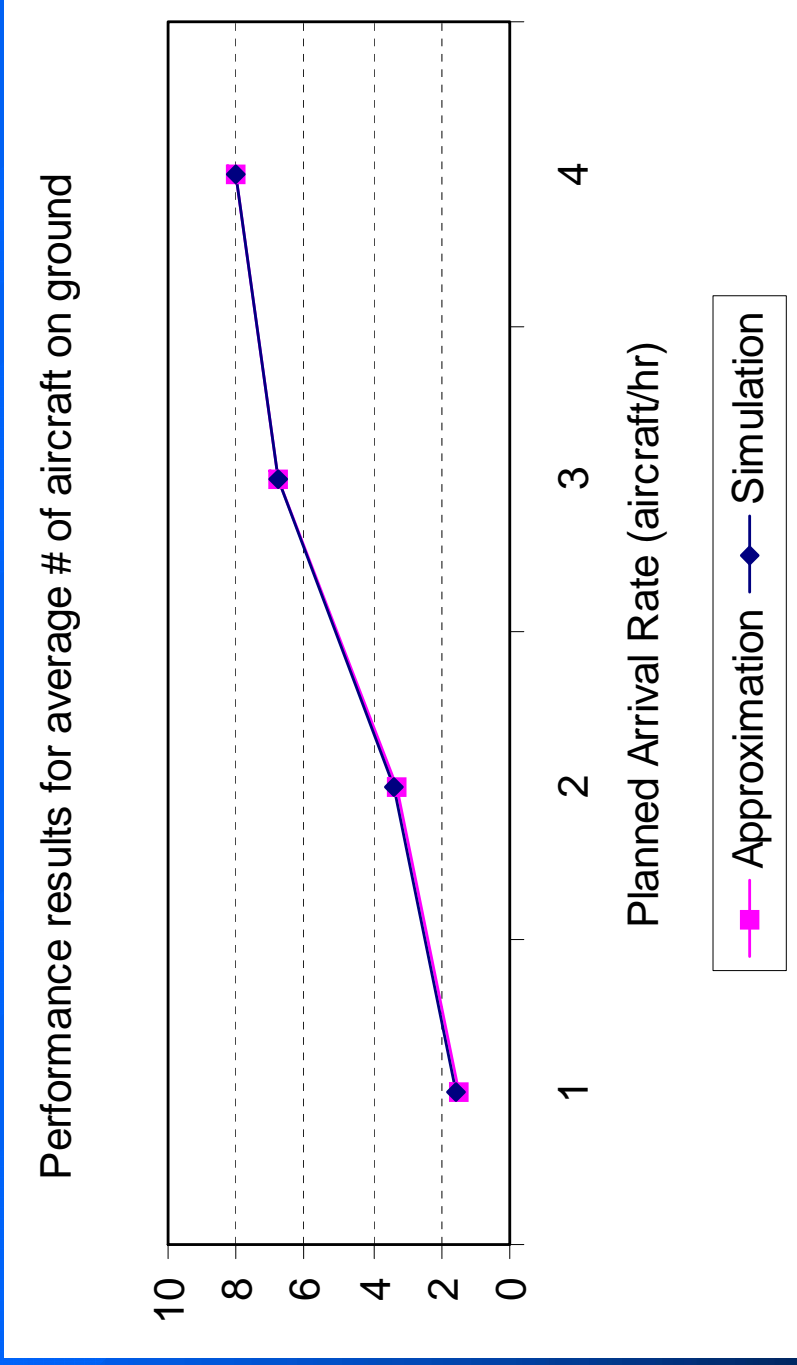
Approximating the AAM

- Capacity planners need fast answers to changes in the system
 - Performance measures are usually estimated with simulation, which is accurate but costly
- Instead, we do most of the analysis calculations with 2-moment approximations
 - Very fast: solution in < 5 seconds on a PC
 - Relatively accurate (relative error $< 10\%$)
- Kanban modeling better captures real constraint for landing authorization than current modeling

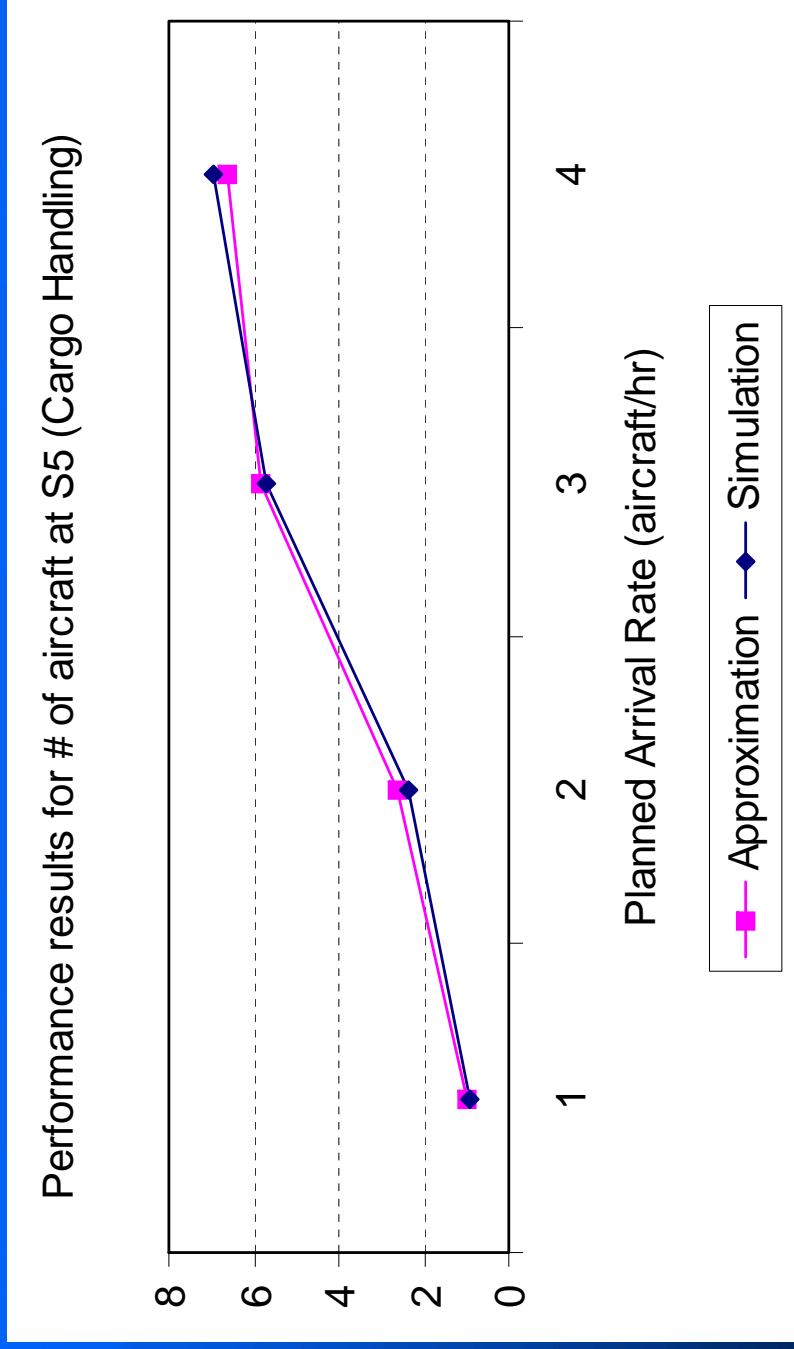
Example: Throughput



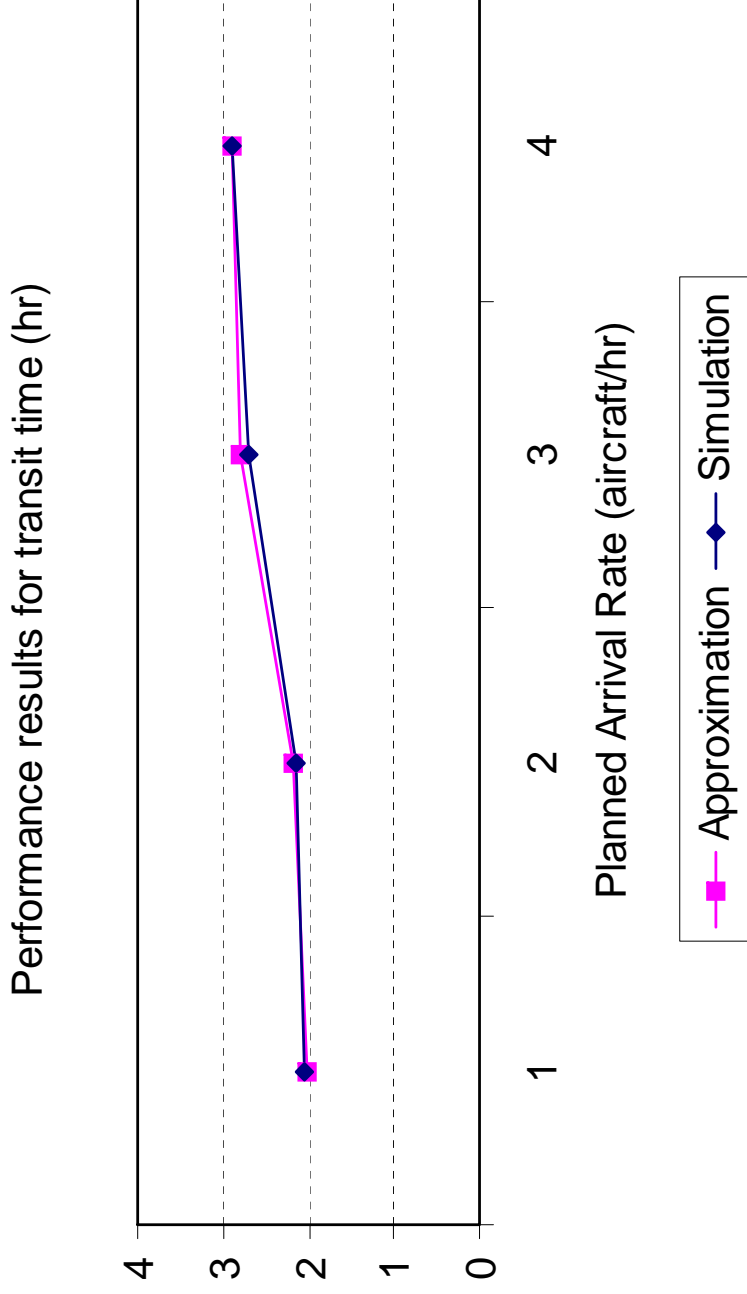
Example: Average on Ground



Cargo Handling is a Bottleneck



Transit Time Up By > 40 %



Network Improvement


- Initial situation:
 - $\lambda = 3.0$ aircraft/hr
 - $K_0 = 6$ (# allowed on ground)
 - Transit time (Approximation) = 2.482 hr.
- Constraint: transit time should be decreased by 20 %.
- Question: What is the minimum # of capacity units to invest at S5 (Cargo Handling) to have transit time decreased by 20 %?

Network Improvement

- Answer: **2**, so that number of servers at Cargo Handling increases from 3 to 5.

Transit time (hr.)

Number of servers at Cargo Handling	Approximation	Simulation
3	2.48	1.84
5	1.89	1.49




Network Improvement

- Additional insight: Number of aircraft at Cargo Handling decreases by 31 %.

Number of aircraft at Cargo Handling

Number of servers at Cargo Handling	Approximation	Simulation
3	4.67	4.25
5	3.22	2.92

-31%



Conclusion

- Fast numerical procedure that gives *quick*, reasonably accurate results
 - Can validate results by using just one simulation run instead of many
- New technical results allow modeling of more complex systems
 - Better modeling of join stations
- Results
 - We can ask and answer managerial questions quickly and easily
 - Improve operations by investing optimally

(Backup slides follow)

Main Assumptions for 2-Moment Approximations

- Traffic processes are approximated as renewal processes
- Traffic process rate and SCV (squared coefficient of variation) are enough to completely characterize the performance measures of the network
- Example: for a GI/G/1 queue

$$W_q = \left(\frac{c_a^2 + c_s^2}{2} \right) \left(\frac{\lambda \tau^2}{1 - \lambda \tau} \right)$$

W_q : delay in queue

c_a^2 : SCV of arrival process

c_s^2 : SCV of service process

λ : arrival rate

τ : mean service time

New Result for a Join Station

- Developed in (Krishnamurthy, 2002)
- Rate and SCV of departure process of a Join station as a function of the 6-tuple:

$$\left(\lambda_1, c_1^2, K_1, \lambda_2, c_2^2, K_2 \right)$$

- Example:

$$\text{If } \rho \equiv \frac{\lambda_1}{\lambda_2} \neq 1:$$

$$\lambda_D = \lambda_1 \left[\frac{1 - \rho^{K_1 + K_2}}{1 - \rho^{K_1 + K_2 + 1}} \right] \left[1 - 0.5 \left(\frac{c_1^2 + c_2^2}{2} \right) \left(\frac{1 - \rho}{1 - \rho^{2(K_1 + K_2) + 1}} \right) \rho^{2(K_1 + K_2)} \right]$$

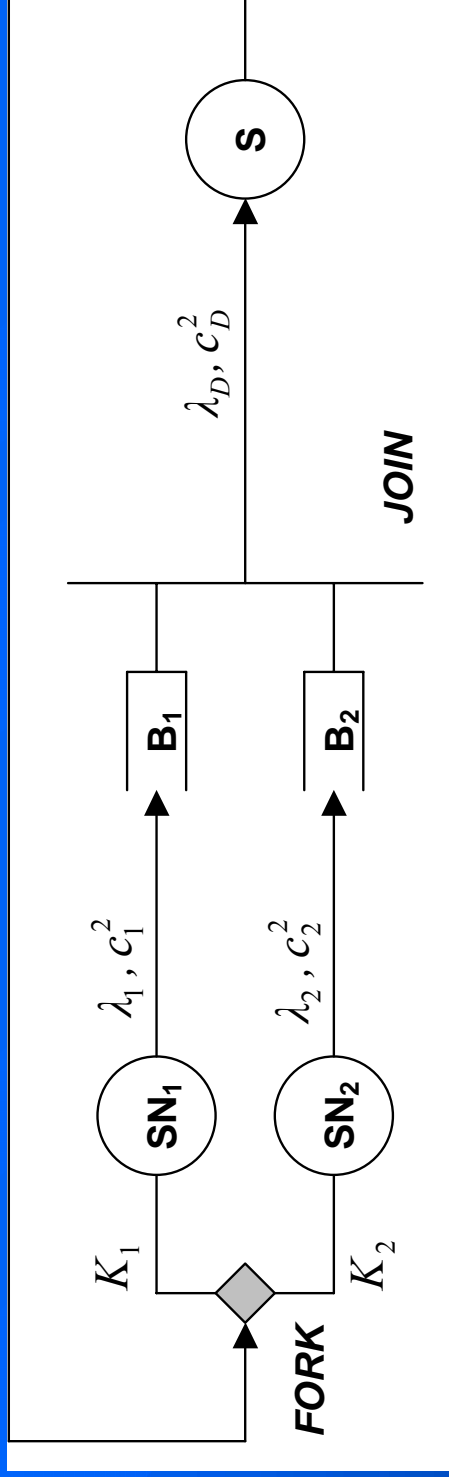
References

- New results for a Join station:
 - A. Krishnamurthy, Analytical performance models for material control strategies for manufacturing systems, Ph.D. Dissertation, University of Wisconsin-Madison, 2002
- Known results for other types of queues and overall performance evaluation of the network:
 - R. Suri, J.L. Sanders and M. Kamath, Performance evaluation of production networks, in: S.C. Graves et al. (Eds.), *Logistics of Production and Inventory* (Elsevier, Amsterdam, 1993)

Reference for the AAM Problem

- Dietz, D. C., Mean value analysis of military airfield operations at an individual airfield, *Journal of Aircraft* 36, No. 5, 1999.
- Start with:
 - Total fleet size equals 20.
 - Maximum number of aircraft allowed on ground $K_0 = 8$.
- Vary planned arrival rate from 1.0 to 4.0 aircraft/hr.
- Compute performance measures of interest:
 - Throughput.
 - Transit time on airfield.
 - Average number of aircraft on ground.
 - Queue length at bottlenecks.

Fork-join Dynamics



SN_1 : Sub - network where entity 1 travels independently

SN_2 : Sub - network where entity 2 travels independently

B_1 : Buffer where entity 1 waits for an arrival of entity 2

B_2 : Buffer where entity 2 waits for an arrival of entity 1

S : Sub - network where joined entity travels

λ_1, c_1^2, K_1 : Rate, SCV and shut - down level of arrival process of entity 1

λ_2, c_2^2, K_2 : Rate, SCV and shut - down level of arrival process of entity 2

λ_D, c_D^2 : Rate, SCV of departure process of joined entity

Generalized Inference: Applications to Mixed Linear Models

Hari Iyer

Department of Statistics
Colorado State University
Fort Collins, COLORADO
<http://www.stat.colostate.edu/~hari>

TERMINOLOGY

- Tsui & Weerahandi (1989) introduced the terms
 - Generalized P-value (GPV)
 - Generalized Test Variable (GTV).
- Weerahandi (1993) introduced the terms
 - Generalized Pivotal Quantities (GPQ)
 - Generalized Confidence Intervals (GCI).
- The term Generalized Inference (GI) refers to inference procedures that are based on the above concepts.

TERMINOLOGY

- Using this approach one can develop hypothesis tests and confidence interval procedures for certain classes of parametric models when exact pivotal quantities are not available.
- In particular, the approach leads to "good" inference procedures in (balanced) normal mixed linear models. But the method is more generally applicable.

OUTLINE

1. What is a GPQ? How it leads to a GCI?
2. What is a GTV? How it leads to a GPV?
3. Simple Examples of GPQs, GCIs, GTVs, and GPVs.
4. Recipe for Constructing GPQs
5. Examples (a) Exact Methods (b) Approximate Methods
6. Generalized Inference in Balanced Mixed Linear Models
7. Some nonstandard applications
8. References
9. Remarks: Historical connections
Issues in unbalanced situations
Extensions

Notation

D = observable data vector

d = observed value of D

ξ = vector of parameters

$\tau = h(\xi)$, a scalar function of ξ about which inference is to be made (test or confidence interval)

WLOG we can assume that $\xi = (\tau, \zeta)$ where τ is the scalar parameter of interest and ζ is a vector of nuisance parameters

Generalized Pivotal Quantity

$R = R(D; \mathbf{d}, \xi)$, a function of D , \mathbf{d} , and ξ , is called a **Generalized Pivotal Quantity** if it satisfies the following two properties (Weerahandi, 93):

1. Distribution of R is free of unknown parameters.
2. The **observed pivotal quantity** $r_{obs} = r = R(\mathbf{d}; \mathbf{d}, \xi)$ depends on ξ only through τ .

Generalized Confidence Interval

An equal tailed $1 - \alpha$ GCI (confidence set) for τ is obtained as the set

$$\{\tau \mid R_{\alpha/2} \leq r \leq R_{1-\alpha/2}\}$$

Often the confidence set reduces to an interval $[L, U]$.

Example: One-Sample Problem

$Y_1, \dots, Y_n \sim \text{iid } N(\mu, \sigma^2)$

y_1, \dots, y_n are observed values

\bar{Y} = sample mean, S = sample standard deviation.

\bar{y} , s the corresponding realized values (known constants)

Usual pivotal quantity: $T = \frac{\bar{Y} - \mu}{S/\sqrt{n}}$

Realized value: $t = \frac{\bar{y} - \mu}{s/\sqrt{n}}$.

Classical t -interval is $\{\mu | t_{\alpha/2} \leq t \leq t_{1-\alpha/2}\}$

i.e., $\left\{ \mu \mid \bar{y} - t_{1-\alpha/2:n-1} \left(\frac{s}{\sqrt{n}} \right) \leq \mu \leq \bar{y} - t_{\alpha/2:n-1} \left(\frac{s}{\sqrt{n}} \right) \right\}$

Example: One-Sample Problem

Define $R = \bar{y} - \left(\frac{s}{S}\right) (\bar{Y} - \mu) = \bar{y} - \left(\frac{s}{\sqrt{n}}\right) \frac{\bar{Y} - \mu}{S/\sqrt{n}}$.

R is said to be a **Generalized Pivotal Quantity** for μ .

The **observed pivotal** is $r = \mu$.

A GCI for μ is $\{\mu | R_{\alpha/2} \leq \mu \leq R_{1-\alpha/2}\}$

Note $R = \bar{y} - \left(\frac{s}{\sqrt{n}}\right) T$ **so,** $R_{\gamma} = \bar{y} - \left(\frac{s}{\sqrt{n}}\right) T_{1-\gamma}$.

$\{\mu | R_{\alpha/2} \leq r \leq R_{1-\alpha/2}\} = \{\mu | t_{\alpha/2} \leq \frac{\bar{y} - \mu}{s/\sqrt{n}} \leq t_{1-\alpha/2}\}$.

Thus, the GCI is the same as the classical t -interval.

Generalized Test Variable and GPV

We wish to test

$$H_0 : \theta \leq \theta_0 \quad \text{versus} \quad H_a : \theta > \theta_0$$

$T = T(\mathbf{D}; \mathbf{d}, \boldsymbol{\xi})$ is called a **GTV** if it satisfies:

1. The distribution of T depends on $\boldsymbol{\xi}$ only through θ . In particular, it is completely determined when θ is specified.

2. The **observed value** of the test variable

$$t = t_{obs} = T(\mathbf{d}; \mathbf{d}, \boldsymbol{\xi}) \text{ is free of unknown parameters.}$$

3. For fixed \mathbf{d} , $\boldsymbol{\xi}$, and t^* , $Pr[T(\mathbf{D}; \mathbf{d}, \boldsymbol{\xi}) > t^*]$ is a **nondecreasing** function of θ .

$$4. \text{ GPV} = Pr[T(\mathbf{D}; \mathbf{d}, \boldsymbol{\xi}) > t_{obs} \mid \theta = \theta_0]$$

Example: One-Sample Problem - GTV

Define $V = \mu - \bar{y} + \left(\frac{s}{S}\right) (\bar{Y} - \mu)$

V is said to be a **Generalized Test Variable (GTV)** for testing

$H_0 : \mu \leq \mu_0$ versus $H_0 : \mu > \mu_0$.

The **observed test variable is $v = 0$** .

$GPV \stackrel{def}{=} P [V \geq v | \mu = \mu_0]$

Here $GPV = P \left[\frac{\bar{Y} - \mu}{S/\sqrt{n}} \geq \frac{\bar{y} - \mu_0}{s/\sqrt{n}} \right] = P [T \geq t_0]$

where t_0 = the usual computed t -statistic.

So, in this example, $GPV =$ Ordinary P-value

A Recipe for Constructing GPQs

Joint distribution of the data vector D is indexed by a

k -dimensional parameter $\xi = (\xi_1, \dots, \xi_k) \in \Omega \subseteq R^k$.

$\tau = h(\xi)$, a scalar function for which inference is required.

Assume

(a) There exist a mapping $f : R^k \times R^k \rightarrow R^k$, such that,
 $(U_1, \dots, U_k) = U = f(D; \xi)$ has a joint distribution free
of ξ .

(b) For each D , there exists a mapping $g(D; \cdot) : R^k \rightarrow R^k$
such that

$$g(D; U) = g(D; f(D; \xi)) = (g_1(D; U), \dots, g_k(D; U)) = \xi.$$

Recipe-continued

Define

$$R = R(\mathbf{D}; \mathbf{d}, \boldsymbol{\xi}) = h(\mathbf{g}(\mathbf{d}; \mathbf{f}(\mathbf{D}; \boldsymbol{\xi}))) = h(\mathbf{g}(\mathbf{d}; \mathbf{U}))$$

1. R is a GPQ for $\tau = h(\boldsymbol{\xi})$.
2. $R_{\alpha/2} \leq \tau \leq R_{1-\alpha/2}$ is an equal-tailed 2-sided GCI for τ .
(One-sided Generalized Bounds obtained in an obvious manner).
3. $T = T(\mathbf{D}; \mathbf{d}, \boldsymbol{\xi}) = h(\boldsymbol{\xi}) - R = \tau - R$, is a Generalized Test Variable for testing $H_0 : \tau \leq \tau_0$ versus $H_a : \tau > \tau_0$.
4. $GPV = Pr [T(\mathbf{D}; \mathbf{d}, \boldsymbol{\xi}) \geq 0 \mid \tau = \tau_0]$

Balanced Mixed Linear Models

In many Balanced Mixed Linear Models the ANOVA sums of squares and sample cell means form a set of complete sufficient statistics. For instance, this is the case when the model is saturated.

The joint distribution of this set has a simple structural representation.

The recipe may be applied to produce tests and confidence intervals for functions of the model parameters.

Simulation studies show that these work well.

Some general theoretical results exist that provide insight into why these methods perform well.

General Setting

- Suppose SS_1, \dots, SS_q are the sums of squares in the ANOVA table corresponding to the random effects. Also suppose $\hat{\beta}$ is the vector of estimates of the cell means generated by the fixed factors.
- Denote the cell means by β_1, \dots, β_p .
- Denote the expected Mean Squares (EMS) by $\theta_1, \dots, \theta_q$
- $\hat{\beta} \sim N(\beta, \Sigma)$ where

$$\Sigma = \theta_1 V_1 + \dots + \theta_q V_q$$

and V_i are matrices of known constants.

GPQs

- The θ_i admit GPQs of the following form:

$$R_{\theta_i} = \frac{ss_i}{U_i} = \frac{(ss_i)(\theta_i)}{SS_i}, \quad i = 1, \dots, q$$

where $U_i \sim \chi_{\nu_i}^2$ (jointly independent).

- β admits a GPQ given by $R_{\beta} = b - R_C Z$ where b is the observed value of $\hat{\beta}$ and R_C is the Cholesky factor (lower triangular) of the matrix $R_{\Sigma} = R_{\theta_1} V_1 + \dots + R_{\theta_q} V_q$.
- Let τ be any function, say $f(\theta, \beta)$ of the model parameters for which a confidence interval is sought. Then $R_{\tau} = f(R_{\theta}, R_{\beta})$ is a GPQ for τ .

One-way Nested Random Model

$$X_{ij} = \mu + A_i + e_{ij} \quad i = 1, \dots, a; j = 1, \dots, n.$$

$$A_i \sim N(0, \sigma_A^2) \text{ and } e_{ij} \sim N(0, \sigma_e^2).$$

All random variables jointly independent.

We want a confidence interval for σ_A^2 .
Methods that have appeared in the literature:

- Tukey-Williams
- Moriguti-Bulmer
- Howe
- Graybill-Wang

Example

$$Z = \frac{\bar{X} - \mu}{\sqrt{\frac{n\sigma_A^2 + \sigma_e^2}{na}}},$$

$$U_2 = \frac{SSb}{\sigma_e^2 + n\sigma_A^2}$$

$$U_1 = \frac{SSw}{\sigma_e^2},$$

$$Z \sim N(0, 1), \quad U_1 \sim \chi_{a(n-1)}^2 \quad U_2 \sim \chi_{a-1}^2$$

$$\mu = \bar{X} - Z \sqrt{\frac{SSb}{naU_2}},$$

$$\sigma_e = \sqrt{\frac{SSw}{U_1}} \quad \sigma_A = \sqrt{\max \left\{ 0, \frac{SSb}{nU_2} - \frac{SSw}{nU_1} \right\}}$$

GPO

A GPO for σ_A^2 is given by

$$\begin{aligned} R &= \max \left\{ 0, \frac{ssb}{nU_2} - \frac{ssw}{nU_1} \right\} \\ &= \max \left\{ 0, (n\sigma_A^2 + \sigma_e^2) \frac{ssb}{nSSb} - \sigma_e^2 \frac{ssw}{nSSw} \right\} \end{aligned}$$

EXAMPLE

Weights of bottles selected from filling machines

Machines			
1	2	3	4
14.23	16.46	14.98	15.94
14.96	16.74	14.88	16.07
14.85	15.94	14.87	14.91

Type 3 Analysis of Variance

	Sum of		Mean	
Source	DF	Squares	Square	EMS
machine	3	5.329425	1.776475	$\sigma_E^2 + 3\sigma_A^2$
Residual	8	1.454600	0.181825	σ_E^2

Confidence Interval for σ_A^2

The GPQ for σ_A^2 is

$$\begin{aligned} R &= \max \left\{ 0, \frac{ssb}{nU_2} - \frac{ssw}{nU_1} \right\} \\ &= \max \left\{ 0, \frac{5.329425}{3U_2} - \frac{1.4546}{3U_1} \right\} \end{aligned}$$

Graybill-Wang interval for σ_A^2 is [0.107, 8.16]

GCI for σ_A^2 is [0.09624, 8.19219] by simulation

GCI for σ_A^2 is [0.09605, 8.152] by numerical evaluation

Other methods work well also

$$\sigma_A^2 + \sigma_E^2$$

$$\sigma_A^2 + \sigma_E^2 = \frac{1}{n}(n\sigma_A^2 + \sigma_E^2) + \frac{n-1}{n}\sigma_E^2.$$

$$R = \max \left\{ 0, \frac{ssb}{nU_2} + \frac{(n-1)ssw}{nU_1} \right\} = \max \left\{ 0, \frac{5.329425}{3U_2} + \frac{2(1.4546)}{3U_1} \right\}.$$

Welch-Satterthwaite, Graybill-Wang work well. GCI is competitive.

	<u>GCI</u>	<u>Graybill-Wang</u>
Lower bound	0.313	0.306
Upper bound	8.403	8.360

Exact calculation for GCI gives [0.31241, 8.398] (maple)

GCI for σ_A^2/σ_E^2 coincides with the usual exact interval based on the ratio MS_B/MS_E .

Brand 1	Machine 1	15.66	15.66	15.70	15.70	15.68	15.70
	2	15.69	15.71	15.68	15.72	15.71	15.72
	3	15.73	15.68	15.73	15.71	15.67	15.72
	4	15.72	15.73	15.74	15.74	15.73	15.75
Brand 2	Machine 1	15.78	15.80	15.78	15.79	15.78	15.79
	2	15.78	15.76	15.76	15.77	15.76	15.77
	3	15.76	15.80	15.78	15.78	15.79	15.78
	4	15.77	15.80	15.78	15.78	15.77	15.78

ANOVA for EXAMPLE

BRAND - 1

Source	DF	Sum of Squares	Mean Square
machine	3	0.008083	0.002694
Residual	20	0.007167	0.000358

BRAND - 2

Source	DF	Sum of Squares	Mean Square
machine	3	0.001312	0.000437
Residual	20	0.002150	0.000108

$$(\sigma_{A1}^2 + \sigma_{E1}^2) / (\sigma_{A2}^2 + \sigma_{E2}^2)$$

$$R = \frac{\frac{ssb1}{n_1 U_{21}} + \frac{(n_1 - 1)ssw1}{n_1 U_{11}}}{\frac{ssb2}{n_2 U_{22}} + \frac{(n_2 - 1)ssw2}{n_2 U_{12}}} = \frac{0.008083}{6U_{21}} + \frac{5(0.007167)}{6U_{11}}}{\frac{0.001312}{6U_{22}} + \frac{5(0.002150)}{6U_{12}}}$$

90% Lower bound GCI Burdick-Graybill 1.068 1.03

90% Upper bound 24.412 ***

$$\sigma_A^2 + \sigma_B^2 + \sigma_E^2$$

Cage	1	2	3	4
Mosquito	1	2	3	3
	58.5	59.5	77.8	77.8
	50.7	49.3	63.8	63.8
	1	2	3	3
	80.9	84.0	83.6	83.6
	65.8	56.6	57.5	57.5
	1	2	3	3
	70.1	68.3	69.8	69.8
	77.8	79.2	69.9	69.9
	69.8	56.0	54.5	54.5
	69.2	62.1	64.5	64.5

ANOVA

Source	DF	Sum of Squares	Mean Square
cage	2	665.675833	332.837917
mosquito(cage)	9	1720.677500	191.186389
Residual	12	15.620000	1.301667

GCI

Burdick-Graybill

95% One sided Lower bound 63.3212

65.5

95% One sided Upper bound 1763.68

1724

Two way Crossed Random Model

$$Y_{ijk} = \mu + A_i + B_j + (AB)_{ij} + E_{ijk}$$

$i = \text{Mice Strain (5)}; j = \text{Day (6)}; k = \text{Mice (5)}$

Source	Df	SS	MS
Strain	4	0.3680	0.0920
Days	5	0.0505	0.0101
Interaction	20	0.1040	0.0052
Error	120	0.4080	0.0034

(Weir, 1949)

Need CI for $\sigma_A^2 / (\sigma_A^2 + \sigma_B^2 + \sigma_{AB}^2 + \sigma_E^2)$.

	<u>GCI</u>	<u>Leiva-Graybill</u>
90% Lower bound	0.157	0.196
90% Upper bound	0.770	0.814

A Mixed Model Example

Source	Df	SS	MS
Temperature (A)	2	616.78	308.39
Speed (B)	3	175.56	58.52
Pressure (C)	2	5.04	2.52
AB	6	809.46	134.91
AC	4	179.08	44.77
BC	6	115.56	19.26
ABC	12	231.12	19.26
Error	36	1248.12	34.67

(Montgomery, 1984)

Temp FIXED, Speed RANDOM

Pressure RANDOM

Need CI for $\mu_1 - \mu_2$

$$V(\widehat{\mu_1 - \mu_2}) = (\theta_{AB} + \theta_{AC} - \theta_{ABC})/24$$

$$= (6\sigma_{AB}^2 + 8\sigma_{AC}^2 + 2\sigma_{ABC}^2 + \sigma_e^2)/24$$

No exact interval available

	<u>GCI</u>	<u>Banerjee (1960)</u>
95% Lower bound	-1.146	-3.28
95% Upper bound	17.742	19.9

An Unbalanced Example

- An artifact measured by each of k labs (or, k methods)
- Lab i makes n_i measurements
- Data are $Y_{ij}, j = 1, \dots, n_i, i = 1, \dots, k$.
- Model is: $Y_{ij} = \mu_i + \epsilon_{ij}$
- μ = the true value.
- $\mu_i - \mu = b_i$ is the “bias” of lab i .
- Parameter of interest is μ .
- Estimate μ using combined information from all labs.
- $\epsilon_{ij} \sim N(0, \sigma_i^2)$.

Three Models

- **Model-1:** One-way random effects model with unequal sample sizes and heterogeneous variances. (See Rukhin and Vangel (1998), Vangel and Rukhin (1999), Rukhin, Biggerstaff, and Vangel (2000), Paule and Mandel (1971, 1982) – Large sample methods).
- **Model-2: (Bounded Bias Model)** b_i are in the (known) interval $[m_i, M_i]$. (Eberhardt, Reeve, and Spiegelman (1989).) Eberhardt et al. derived a minimax MSE linear estimator for μ . Proposed approximate CIs.
- **Model-3: (GUM model)** b_i have known distribution, say F_i . (see, Expression of Uncertainty in Measurement (ISO GUM) (1995); the distributions are referred to as *type-B* distributions.)

GPQ in Model 3

$$R^*(\mathbf{D}; \mathbf{d}, \boldsymbol{\theta}) = \bar{y}_{\mathbf{W}} - \bar{b}_{\mathbf{W}} - Z^* \left(\sum_{i=1}^k n_i Q_i / SS_i \right)^{-1/2}$$

where,

$$\bar{y}_{\mathbf{W}} = \sum_{i=1}^k W_i \bar{y}_i / \sum_{i=1}^k W_i, \quad \bar{b}_{\mathbf{W}} = \sum_{i=1}^k W_i b_i / \sum_{i=1}^k W_i,$$

$$U_i = \bar{Y}_i - b_i, \quad W_i = \frac{n_i SS_i}{\sigma_i^2 SS_i}$$

$$\tau_0^2 = \frac{1}{w_1 + \dots + w_k}$$

$Z^* = (\bar{U}_{\mathbf{w}} - \mu) / \tau_0 \sim N(0, 1)$, and

$Q_i = SS_i / \sigma_i^2 \sim \chi_{n_i-1}^2, i = 1, \dots, k$.

Example

Zinc ($\mu\text{g/g}$) in non-fat milk powder

Method	n_i	\bar{y}_i	s_i	M_i
1	8	45.21	1.68	5.880
2	12	46.63	0.47	0.466
3	22	46.26	0.82	0.927
4	8	47.05	1.44	0.230

Distribution	Lower Bound	Upper bound
Uniform $[-M_i, M_i]$	45.85	47.05
$N(0, M_i/3)$	46.03	46.86

Example

Tolerance bounds for the distribution of true values when there are measurement errors

$$X_{ij} = \mu + A_i + e_{ij} \quad i = 1, \dots, a; j = 1, \dots, n.$$

$$A_i \sim N(0, \sigma_A^2) \text{ and } e_{ij} \sim N(0, \sigma_e^2).$$

All random variables jointly independent.

Need a γ -content, $1 - \alpha$ confidence, upper tolerance-bound for the distribution of $\mu + A_i$, i.e., for $N(\mu, \sigma_A^2)$.

This is equivalent to an upper confidence bound for

$$\mu + z_\gamma \sigma_A.$$

Wang and Iyer (1994, Technometrics) have discussed this problem.

Example-continued

$$Z = \frac{\bar{X} - \mu}{\sqrt{\frac{n\sigma_A^2 + \sigma_e^2}{na}}}, \quad U_1 = \frac{SSw}{\sigma_e^2}, \quad U_2 = \frac{SSb}{\sigma_e^2 + n\sigma_A^2}$$

$$Z \sim N(0, 1), \quad U_1 \sim \chi_{a(n-1)}^2, \quad U_2 \sim \chi_{a-1}^2$$

$$\mu = \bar{X} - Z \sqrt{\frac{SSb}{naU_2}},$$

$$\sigma_e = \sqrt{\frac{SSw}{U_1}} \quad \sigma_A = \sqrt{\max \left\{ 0, \frac{SSb}{nU_2} - \frac{SSw}{nU_1} \right\}}$$

GPO

A GPO for $\theta = \mu + z_\gamma \sigma_A$ is given by

$$\begin{aligned} R &= \bar{x} - Z \sqrt{\frac{ssb}{naU_2}} + z_\gamma \sqrt{\max \left\{ 0, \frac{ssb}{nU_2} - \frac{ssw}{nU_1} \right\}} \\ &= \bar{x} - (\bar{X} - \mu) \sqrt{\frac{ssb}{SSb}} \\ &\quad + z_\gamma \sqrt{\max \left\{ 0, (n\sigma_A^2 + \sigma_e^2) \frac{ssb}{nSSb} - \sigma_e^2 \frac{ssw}{nSSw} \right\}} \end{aligned}$$

Historical Connections

- GPV and GCI are intimately related to Fisher's Fiducial Inference and Fraser's Structural Inference.
- Fiducial Inference and Structural Inference allow one to make probability statements about model parameters somewhat akin to Bayesian Posterior Distributions for parameters, but do not rely on any prior distributions for the parameters – (Savage: “..eat the Bayesian omelette without breaking the Bayesian Egg”).
- The probability statements are EXACT in their own setting but do not seem to have satisfactory frequentist interpretations. This led to quite a controversy over the use of Fiducial/Structural methods during the mid and latter part of the 20th century.

Historical Connections

- GPs and GCs have a fiducial/structural flavor to them but Tsui and Weerahandi have put forward these ideas in a frequentist context. Although the inference is only APPROXIMATE in all but the simplest problems, Weerahandi refers to them as EXACT methods. The exactness properties of the procedures refers to their own setting and not to the usual frequentist setting.
- Methods for developing GPs and GTVs were discussed by Iyer and Patterson (2002) where they used Fraser's structural distributions for parameters to construct GPs and GTVs. No other general methods appear to be available.

Historical Connections

- Andy Chang (2001) proposed the method of SURROGATE VARIABLES and derived some confidence interval procedures for a class of mixed models. His approach is essentially an application of Fraser's structural inference to construct GPQs for this class of problems.
- It is not clear whether generalized inference is EQUIVALENT to structural inference.
- IGNORING philosophical issues related to the meanings of fiducial or structural probability statements, if one examines frequentist properties of these methods, more often than not, they lead to competing procedures and often methods better than what is currently available. In many cases, the methods can be shown to also be obtainable using Bayesian arguments.

Structural Distributions – Basic Idea

Let $Y \sim N(\mu, 1)$. Then Y has the structural representation

$$Y = \mu + Z$$

where $Z \sim N(0, 1)$.

Suppose an observed value of Y is $y = 2$. We infer that a value z for Z has been realized such that

$$2 = \mu + z$$

If we want to know how plausible it is that $\mu = 10$, this is equivalent to asking how plausible it is that $z = -8$. The known distribution of Z helps us assess this. Thus, the distribution of Z induces a distribution on μ (called ‘Structural Distribution’ by Fraser). In this example the induced distribution on μ is $N(2, 1)$. We may write

Structural Distributions - 2

Y_1, \dots, Y_n iid sample from $N(\mu, \sigma^2)$.

Sufficient statistics: \bar{Y}, S^2 .

Structural representation:

$$\bar{Y} = \mu + \frac{\sigma}{\sqrt{n}}Z \quad \frac{(n-1)S^2}{\sigma^2} = U$$

Substitute observed values (\bar{y}, s^2) for \bar{Y}, S^2 and get:

$$\mu = \bar{y} - \left(\frac{s}{\sqrt{n}}\right) \left(\frac{Z}{\sqrt{U/(n-1)}}\right) = \bar{y} - \left(\frac{s}{\sqrt{n}}\right) T_{n-1}$$

Thus, μ and σ may be thought to have a joint structural distribution induced by the joint distribution of Z and U .

Structural Distributions - continued

Let $\tau_\gamma = \mu + Z_\gamma\sigma$, the γ^{th} percentile of the $N(\mu, \sigma^2)$ distribution. Suppose we want an upper confidence bound for τ_γ with confidence coefficient $1 - \alpha$.

Derive the structural distribution of τ_γ and use the γ^{th} percentile of this distribution as an upper confidence bound for τ_γ . Such a bound is often referred to as a γ -content, $1 - \alpha$ coverage, one-sided tolerance bound for the distribution $N(\mu, \sigma^2)$.

FACT: The structural approach results in exactly the same tolerance bound as does the classical method based on a noncentral t -distribution.

Detecting Meaningful Changes in Short-Term Military Attrition: Application of the Random Effect Model and Agreement Testing

Yuanzhang Li, Timothy Powers and Margot Krauss

Introduction

About one third of the first-term enlistees in each of the military services fail to complete their enlistment terms. The period of highest attrition occurs during the first six months of service. The cost of recruiting, processing and providing basic training to an individual is estimated to be as high as \$30,000. Roughly 15% of the 120,000 recruits who begin basic training are discharged prior to completion, resulting in a replacement cost of over \$500M per year. Accordingly, attrition reduction targets are frequently discussed as a sensible means of cost savings. These targets are often discussed in the context of short time spans, such as reducing monthly attrition by a specified amount.

When a specific attrition reduction goal is established, it is naturally desired that subsequent data be examined to determine whether the goal is being met. Similarly, it is always of interest to know if there is a short-term, unexpected upward spike in attrition, either overall or within a particular service or training site. Such deviations in attrition, when not in conjunction with a change in demographic features of the recruit population or other known covariates of attrition may be related to policy or other factors within the control of military managers. Knowing what these factors are can help to minimize attrition by eliminating unnecessary risk factors. Unfortunately, such a determination is generally difficult to make, as relatively large fluctuations in short-term attrition rates may be caused by seasonal patterns, time trends, and differences in the demographic profile of recruits over time, or simply random fluctuations.

The aim of this study is to develop attrition modeling that will account for these factors in order to better detect changes in core attrition rates at short-term intervals. We will use monthly accession and attrition data over 1995-1999 to develop short-term attrition prediction models. These models will then be used to predict monthly attrition for CY 2000, and these predictions will then be compared to actual monthly attrition for this period. Particular attention will be given to variance estimates of the predictions, as this will play an important role in determining whether an attrition rate is significantly different from that which was predicted.

Subjects and Methods

All first-time enlistees beginning active duty military service during January 1995 - December 2000 were included in the analyses. The enlistees were grouped according to the month and year of beginning military service (accession). Accession records on these individuals were linked with military personnel records to determine whether or not a subsequent early attrition occurred. For each month/year accession group, attrition percentages during the first 1, 2 and 3 months of service were then determined. In addition, a "demographic profile" was developed for each group, including the

distribution of Armed Forces Qualification Test scores (AFQT), gender distribution, race distribution, etc. These factors have been found in previous studies to be strongly related to likelihood of attrition.¹

Service-specific attrition rates by month/year group over the 60-month period 1995-1999 were first examined, and then adjusted for both seasonal and long-term trends by differencing. Remainder attrition after this differencing for the sequence of month/year groups was then examined for homogeneity. Finally, regression models were then developed to regress the remainder attrition rates against the demographic profiles. Both fixed and random effect models were used in the analysis.

A dynamic regression model was then used to predict attrition rates for the month-year groups of CY 2000. For example, in predicting attrition for enlistees beginning duty in January 2000, all 60 months of historical data from January, 1995 – December 1999 were used in the regression. Similarly, when predicting the attrition rate for enlistees beginning duty in February 2000, data from the previous 61 months (including January 2000) were used in the regression. A measure of the agreement between the predicted attrition rate and actual attrition rate a future month/year group will experience, proposed in the appendix, is used to detect significant differences in actual attrition from expected levels.

Results

First time enlistees were grouped by active duty month and year as well as by service, and raw attrition rates at 1, 2, and 3 months from beginning duty were computed. Figures 1-4 show the attrition rate within three months service for the Army, Navy, Marines, and Air Force from January 95 to June 2000. It can be seen that there is a common seasonal pattern to these attrition rates, with recruits beginning duty in the early part of the calendar year generally having somewhat higher attrition than those entering in the summer months.

¹ AMSARA Annual Report, 1999.

Figure 1. The Attrtiron RateWithin 3 Months by AD Months of service in the Army

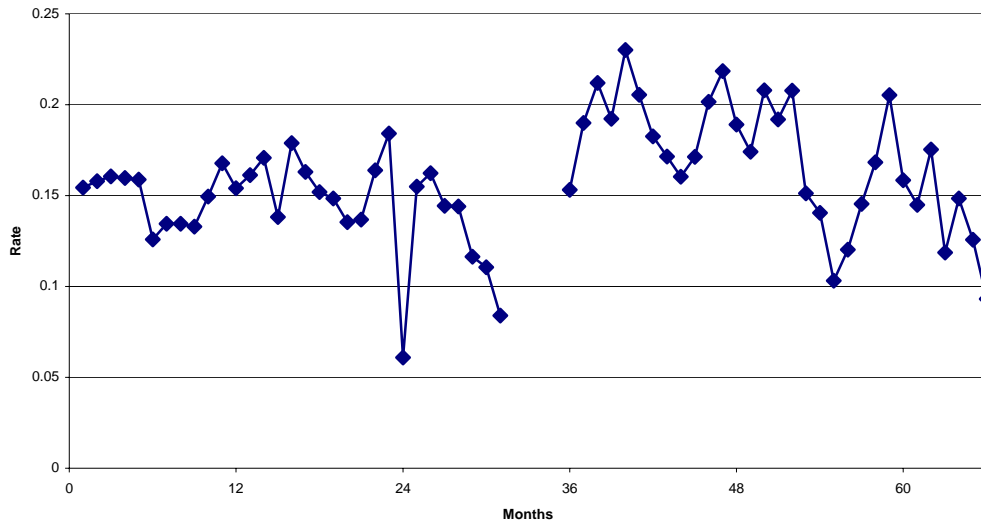


Figure 2. The Attrition Rate 3 Months by AD Months of service in the Navy

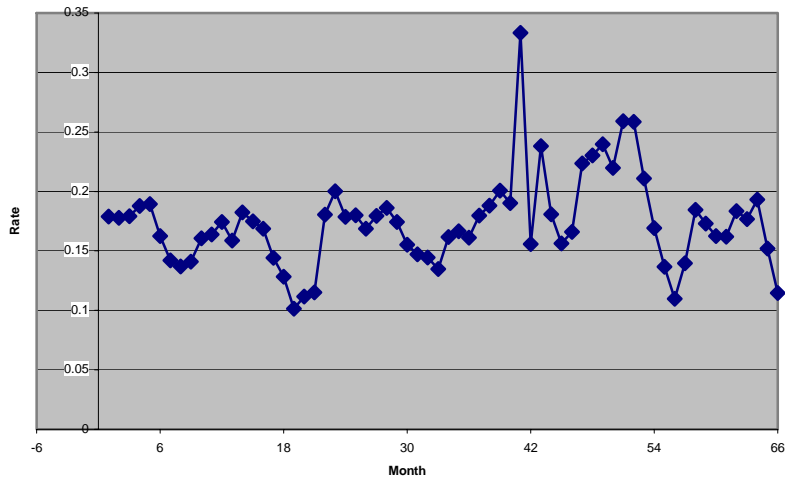


Figure 3. Attrition Rate within 3 months of service by AD months in the Marines

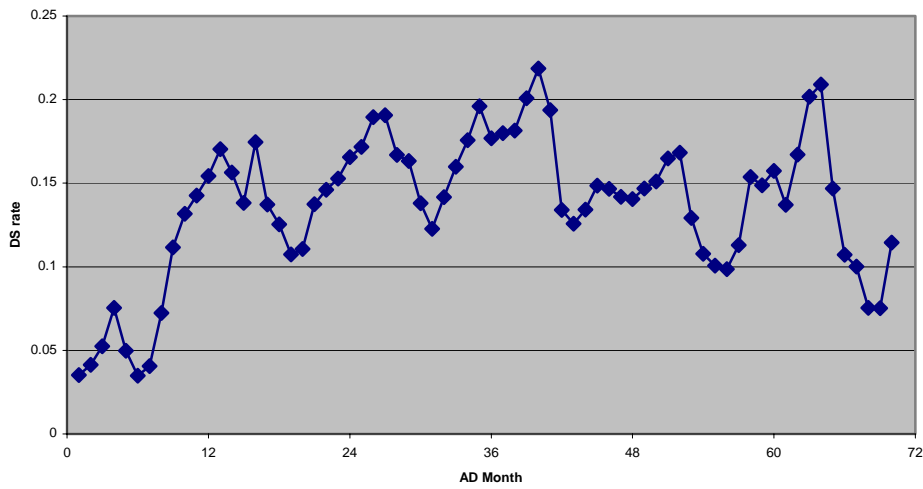
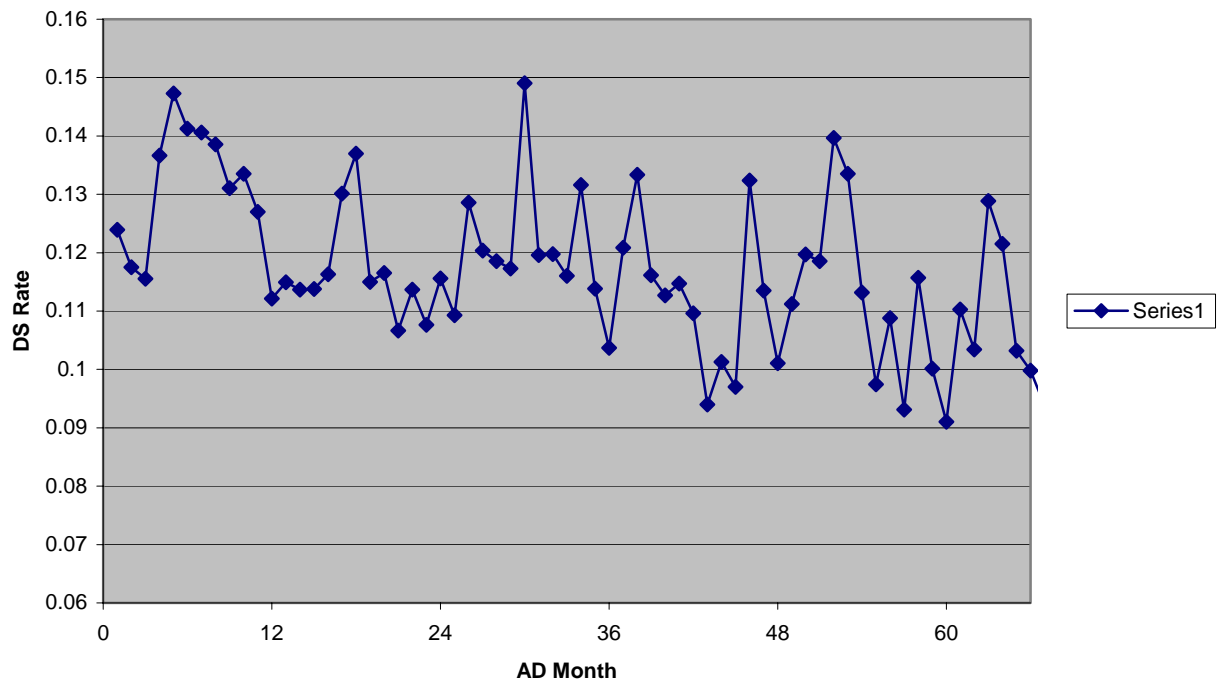


Figure 4. Attrition Rate within 3 months of service by AD Months in the Air Force



Anecdotal information and prior AMSARA analyses suggest that this phenomenon may be related to differences in the types of individuals beginning service at various times of the year. For example, those coming during the summer months have been shown to be a quite homogeneous group, consisting mostly of young applicants who have just recently graduated from high school. Those entering at other times of the year are somewhat more

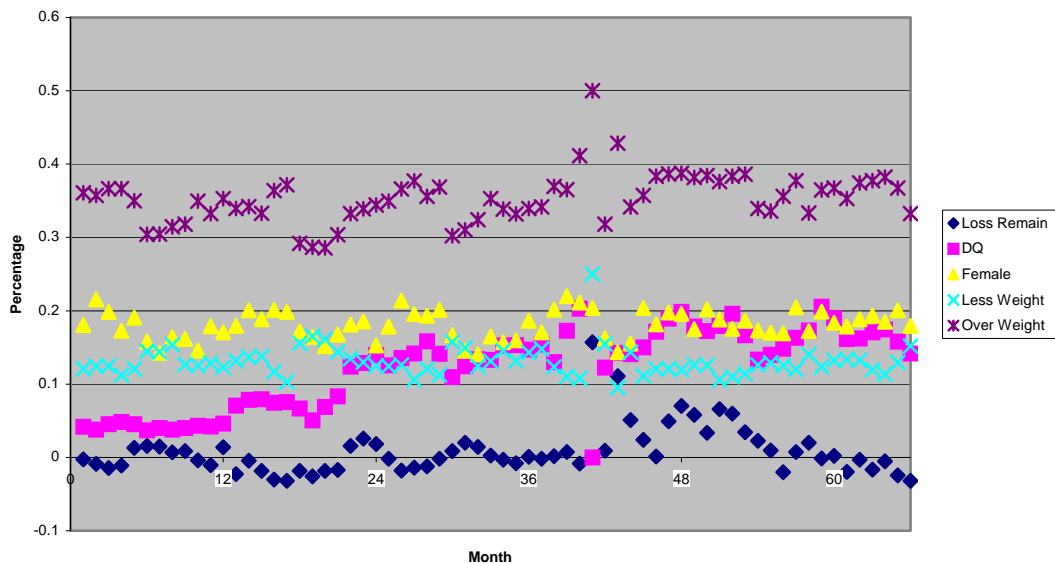
heterogeneous, consisting of some recent high school graduates and a mix of older individuals with various reasons for joining the military. Also by anecdotal information, the fitness level of this latter group is generally lower than that of recent high school graduates.

We account for this seasonal attrition pattern by subtracting from each monthly attrition rate the average level of attrition for that month over the time period studied (1995-2000). For example, to calculate the remainder attrition for January 1995, the average attrition rate for each January over 1995-2000 is subtracted from the raw rate in January 1995.

Figure 5 shows the remainder attrition rates after the seasonal differencing, along with demographic features of the month/year groups. It can be seen that the monthly remainder loss rates are still not pure random noise, but instead are related to the demographic factors which show non-random patterns over time. It can also be seen that the remainder attrition rates exhibit long-range changes over the time period examined. For example, the levels in 1995 and 1996 are considerably lower than those in 1999, and this relation is not linear.

We regress this remainder attrition against several demographic factors which were found in previous AMSARA studies to be strongly related to attrition, such as sex, race and age distributions, and body mass index (BMI). We will also consider variables of time and the square of time to account for long-term differences in attrition rates over the modeled time period.

Figure 5. The Monthly Remain Attrition Rate and Demographic Profiles



In order to avoid including too many variables in the model simultaneously, all candidate variables are examined in a forward stepwise manner in the order of their significance. The final model includes only significant variables ($p < 0.10$).

We used both a random effect regression model and a fixed regression model to fit the historical data, with attrition rates after 1, 2 and 3 months of service as the outcome. It was found that the variance between months was virtually zero for each service except the Air Force. In other words, after controlling for the demographic factors, the residual attrition rates were subject to the same distribution and thus homogeneous, with the exception of those for the Air Force. The remaining analyses will therefore employ only a fixed effects model for the Army, Navy and Marines, while considering both a fixed and random effects model for the Air Force.

Comparison of Predicted to Actual Attrition

Table 1 shows predicted 3-month attrition rates and associated standard errors, by service, for an example month/year group (those beginning service in March 2000). Also shown are the corresponding actual rates and associated standard errors, and the measure of agreement between the predicted and actual for each group. It is seen that the actual 3-month attrition percentage in the Army was 7.8%, which is quite close to the predicted level of 7.7%. Accordingly, the z-score for this difference was not statistically significant (i.e. $|z| < 1.96$), indicating that the 3-month attrition levels observed for the Army was in accordance with what would normally be expected after accounting for time trends, seasonal trends, and features of the recruits who began duty at this time.

For the Navy, the actual attrition was 14.9%, whereas the predicted level was 16.2%. This difference is larger than that seen between the actual and predicted levels for the Army, but not enough to achieve statistical significance. The z-score comparing actual to predicted attrition for the Navy was -1.90 .

Conversely, the z-score for the Marines indicates a high level of significance. The actual 3-month attrition was 14.5%, whereas the predicted level was only 8.2%. Accordingly, the high z-score for this difference indicates that the observed attrition was much higher than was predicted from the modeling.

Finally, it is seen that the Air Force loss rate, when examined by the fixed effect model, is significantly higher than expected. However, when the random effects model is used to account for variability that was observed across months for the Air Force, the result is no longer statistically significant. It is this latter result that would be used in practice, as the random effects model was determined to be a better choice for the Air Force.

**Table 1. Example of Actual versus Predicted Attrition and Agreement Testing:
Subjects Beginning Active Duty in March, 2000**

Service (Model)	Actual		Predicted			Agreement Z-score
	Loss rate	Std Err	Loss Rate	Std Err	Param Err	
Army (Fixed)	0.077	0.004	0.078	0.003	---	-0.28
Navy (Fixed)	0.149	0.007	0.162	0.002	---	-1.90
Marines (Fixed)	0.145	0.008	0.082	0.003	---	7.12
AF (Fixed)	0.097	0.006	0.082	0.001	---	2.43
AF (Random)	0.097	0.006	0.083	0.002	0.009	1.51

Table 2 summarizes the agreement results of modeled versus actual attrition at 1, 2 and 3 months of service for recruits beginning duty January – September 2000. It can be seen from the January results that attrition among recruits beginning duty in that month was significantly lower than expected in both the Army and Navy at all lengths of follow-up (1, 2 and 3 months). Army attrition was then lower than expected in June and July, whereas Navy was lower than expected in July. Attrition among Marines recruits was higher than expected at virtually all follow-up times over March – July.

Of the nine month/year groups examined, the Army had six months with at least one significant attrition deviation from the predicted level, the Navy had five, the Marines seven, and the Air Force four. These results indicate that attrition among recruits beginning military service during CY 2000 was not completely explainable on the basis of long-term trends, seasonal trends, and demographic makeup of the recruit populations.

It is difficult to say whether this large number of significant results is indicative of features of the particular recruit populations that were not included in the modeling, or of the ever-changing military training environment. The modeling results do not appear to have systematic bias, as the actual attrition is roughly evenly distributed above and below predicted levels.

Month (2000)	Month	Army	Navy	Marines	Air Force
Enter AD	Since AD	Fixed	Fixed	Fixed	Random Effect
January	1	-3.66	-4.29	0.12	-0.88
	2	-5.24	-4.18	-0.09	-0.92
	3	-4.73	-3.78	0.65	-0.42
February	1	-0.33	0.88	-0.93	-1.87
	2	0.36	0.16	1.90	-1.38
	3	0.70	-0.52	2.87	-0.82
March	1	1.88	-1.33	3.03	0.51
	2	0.44	-2.26	4.86	1.61
	3	-0.28	-1.90	7.12	1.51
April	1	0.44	-0.73	4.71	-0.53
	2	-1.97	-0.35	6.70	0.33
	3	-1.65	-0.27	7.12	0.98
May	1	1.22	-0.71	3.11	-3.40
	2	1.31	-3.69	4.33	-2.22
	3	2.00	-4.15	3.95	-2.31
June	1	4.11	-5.66	0.57	-2.22
	2	2.96	-7.53	2.14	-0.72
	3	3.73	-6.29	2.08	-1.01
July	1	4.42	-1.35	2.00	-0.72
	2	7.55	-0.80	3.15	0.14
	3	6.69	-0.66	2.17	0.15
August	1	4.28	-0.30	-1.15	-3.42
	2	-1.53	-1.23	0.47	-2.65
	3	-4.21	-1.53	0.72	-2.98
September	1	0.47	-4.40	-3.24	-4.47
	2	-0.78	-0.41	-0.45	-4.63
	3	-1.16	-1.11	0.57	-4.36

Discussion

Determining a reason (or set of reasons) for particular spikes in short-term recruit attrition for a particular service branch will require deeper focus on that branch, and perhaps on particular basic training sites within that branch. For example, this might include examining the two Marines basic training sites separately to see if the increase seen during March-June was a local phenomenon, or whether it was observed at both sites.

The coded reasons for Marines discharges during this period might also be compared to see if there was a spike in a certain category of discharges that might indicate group dynamics or other such effects. For example, there have been occasional episodes of

“contagious” psychological problems within groups of recruits, such as an outbreak of suicide ideation episodes within a recruit class at one training site a few years ago.

Policy changes, traumatic or other unusual events, and contagious motivational or attitude problems are other possibilities that might be investigated. For example, the past few years have seen a considerable increase in the number of programs designed to keep recruits in basic training who would have been discharged in past years. For example, an injury rehabilitation program at the Army’s Fort Jackson is now a mandatory stop for recruits with injuries that previously would have led to a discharge. This program has recently been extended to the other four Army basic training sites. Even if such a program served only to delay attrition, this would result in a downward spike in short-term attrition rates.

Future study of these short-term attrition rates should therefore involve closer collaboration with the services, and the individual training sites. Accounting for local phenomena may be the key to fully and successfully modeling and monitoring short-term attrition rates.

Appendix: Measure of agreement between a pooled estimation and a new estimation

Let r_1, \dots, r_K be attrition rates of first K months, r_W be the weighted average of these, and r_{K+1} be the attrition rate for month K+1. The agreement between r_W and r_{K+1} is usually defined by $(r_{K+1} - r_W) / \sqrt{V(r_{K+1}) + V(r_W)}$. In general, $V(r_{K+1}) = v_{K+1}$, where v_i is the sampling variance of the estimation of r_i , $i=1, 2, \dots, K+1$. However, if a random effects model is used, variance between months may exist, represented by τ^2 . Moses (2002) suggests add the variance between months to variance of r_{K+1} , i.e. $V(r_{K+1}) = v_{K+1} + \tau^2$. This suggestion is natural. The estimation in the K+1 month should be treated the same as those, which derived from the previous K months. We will follow his suggestions to measure the agreement between the prediction from the regression model and the observed attrition rate for the month K+1.

Measure of agreement between a predicted value and the observed value in the K+1 month:

Considering the attrition rate r_i in the month i , $i=1, 2, \dots, K$, which may depends on demographic profile by months and other factors, such as season, time trend etc. The fixed regression model

$$F(r_i) = \beta X_i + \varepsilon_i,$$

X_i is the demographic profile by months, ε_i is the residual error term, which is from the regression. However, since $F(r_i)$ is a function of r_i , which is an estimation based on the monthly base, with the sampling error e_i , hence a random effect model should be used,

$$F(r_i) = \beta X_i + \varepsilon_i + e_i,$$

We assume the residuals are subject to a normal distribution with mean of zero and the variances are

$$V(e_i) = \sigma^2_i \text{ (variance within month)}$$

$$V(\varepsilon_i) = \tau^2 \text{ (variance between months)}$$

$$\text{then for random effect model: } V(F(r_i)) = \sigma^2_i + \tau^2$$

Using the above model, based on the demographic profile of the month K+j, $j=1, 2, \dots$, the predictor of the attrition rate in the month K+j is \hat{r}_{K+j} . The agreement between the predictor and the observed attrition rate could defined by

$$z_1 = (F(\hat{r}_{K+j}) - F(r_{K+j})) / \sqrt{v(F(\hat{r}_{K+j})) + v_{K+j}} \text{ or}$$

$$z_2 = (F(\hat{r}_{K+j}) - F(r_{K+j})) / \sqrt{v(F(\hat{r}_{K+j})) + v_{K+j} + \tau^2}$$

where z_1 is the classic measure for the agreement, z_2 is the new way to measure the agreement, which add the variance between months to the variance of the month $K+j$.

In the dynamic regression process, K is increasing when more data from the future months are available. The prediction is always for the month $K+1$.

References

Li Y, Powers T, Roth, N, 1994. Random effect Linear Regression Meta-Analysis Models with Application to the Nitrogen Dioxide Health Effects Studies, Air & Waste, Vol 44, 1994

Loannidis, et al. Issues in comparisons between Meta-analyses and Large trials. JAMA April 1998

Mosses et. al. Comparing results of large clinical trials to those of meta-analyses. Statistics in Medicine 2002; 21