

Bayesian Estimation Under Informative Sampling

Based on work of T. D. Savitsky, Daniell Toth, Michail Sverchkov

T. D. Savitsky

Office of Survey Methods Research

October 22, 2014



Table of contents

1 Goals

2 Motivating Application

3 Estimation Model

- Non sequitur Example
 - Estimating Population Density from Sampled Units
- Pseudo-population Estimation Statistics

4 Bayesian Accounting for Informative Sampling

- Empirical Generation of Pseudo-populations

5 A Simulation

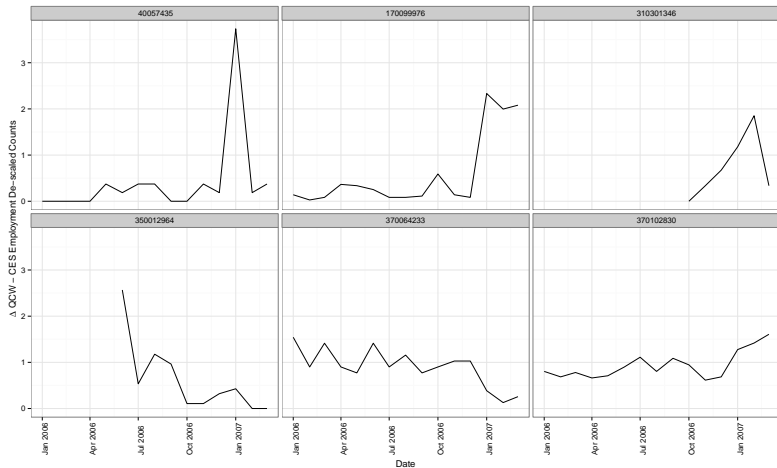
Employment Count Reporting Errors

- Make inference on the population distribution
- When the data are acquired under an informative sampling design
- Where $f(\mathbf{y}_o|\boldsymbol{\lambda}) \neq f(\mathbf{y}_p|\boldsymbol{\lambda})$
- Analyst may not know the sampling design to parameterize it

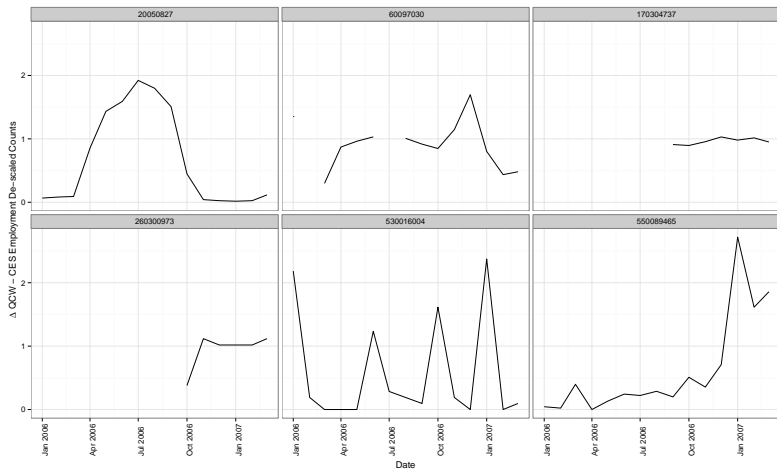
Employment Count Reporting Errors

- Establishments report number of employees
- On both Quarterly Census of Employment and Wages (QCEW)
- And Current Employment Survey (CES)
- Response Analysis Survey (RAS)
- Our data are $y = |y_{QCEW} - y_{CES}|$
- Collected on $N = 9777$ establishments, each for $T = 15$ months.
 - June, 2006– March, 2007
- Produces $N, T \times 1$ functions, $\{\mathbf{y}_i\}_{i=1, \dots, N}$
- Functions each standardized to variance 1
- Desired inference about set of latent functional bases

Randomly-selected observed y_i



More Randomly-selected observed y_i



Stratified Sampling Design

- Sampled 2917 establishments - 2013 responded from $N = 9777$
- Stratified design with $H = 6$ strata
- Assigned to strata based on *a priori* “interesting” features
- Focused on selected portions of functions
- e.g. spike at year-end indicates counting checks, rather than employees

Bayesian Estimation Model: DP mixture over GP's

$$\mathbf{y}_i^{T \times 1} \sim \mathcal{N}_T(\mathbf{f}_i, \tau_\epsilon^{-1} \mathbb{I}_T), \quad i = 1, \dots, N$$

$$\mathbf{f}_i \sim \mathcal{N}_T\left(\mathbf{0}, \mathbf{C}^{T \times T}(\boldsymbol{\theta}_i)\right)$$

$$\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N | G \sim G$$

$$G \sim \text{DP}(\alpha_0, G_0)$$

\Updownarrow

$$\mathbf{f} | G \sim \int \mathcal{N}_T(\mathbf{0}, \mathbf{C}(\boldsymbol{\theta})) G(d\boldsymbol{\theta})$$

$$\mathbf{f}_i \sim \mathcal{N}_T \left(\mathbf{0}, \mathbf{C}(\boldsymbol{\theta}_i) \right)$$

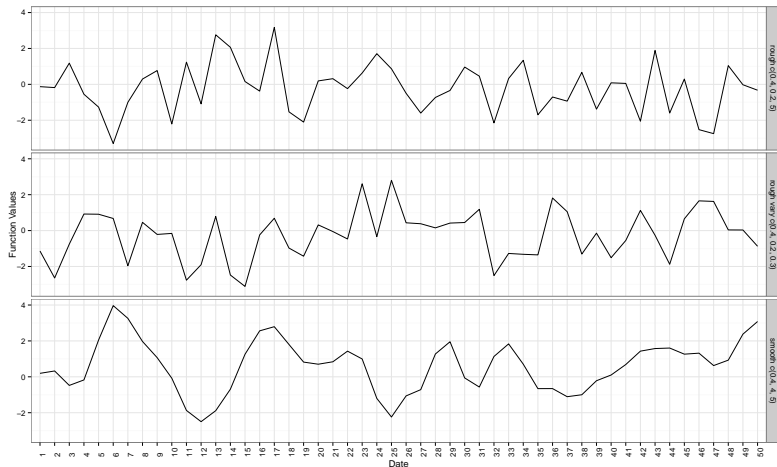
$$\mathbf{C}(\boldsymbol{\theta}_i) \equiv \mathbf{C}_i = (C_{f_{ij}, f_{il}})_{j, l \in (1, \dots, T)}$$

$$C_{f_{ij}, f_{il}} = \frac{1}{\theta_{i,1}} \left(1 + \frac{(t_{ij} - t_{il})^2}{\theta_{i,2}\theta_{i,3}} \right)^{-\theta_{i,3}}$$

- $\theta_{i,1}$ controls the vertical **magnitude**
- $\theta_{i,2}$ controls the base **length-scale** \leftrightarrow **wavelength**
- $\theta_{i,3}$ controls **smooth deviations** across scales
 - As $\theta_{i,3} \uparrow \infty$, converge to single length-scale, $\theta_{i,2}$

GP Draws

- Top - “Rough”: $(\theta_1, \theta_2, \theta_3) = (0.4, 0.2, 5.0)$
- Middle - “Rough with Variation”: $(\theta_1, \theta_2, \theta_3) = (0.4, 0.2, 0.3)$
- Bottom - “Smooth”: $(\theta_1, \theta_2, \theta_3) = (0.4, 4.0, 5.0)$



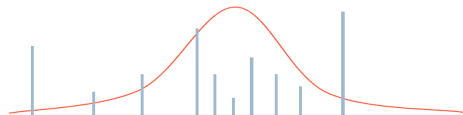
Draws from $G \sim \text{DP}$ are discrete, a.s.

Dirichlet Process

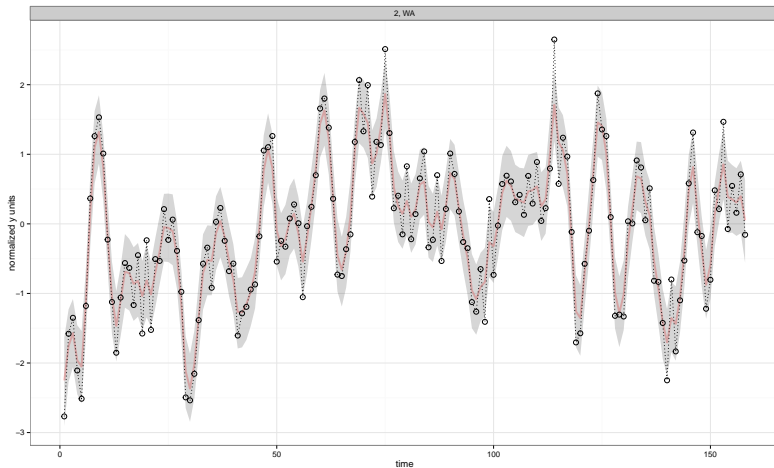
- ▶ Consider Gaussian G_0



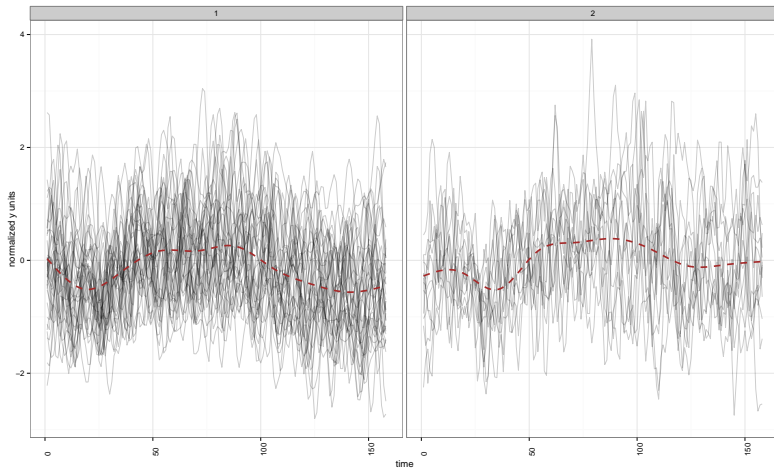
- ▶ $G \sim \text{DP}(\alpha, G_0)$



f_i : Monthly CPS data



f_* : Monthly CPS data



Estimating Population Density from Sampled Units

- Consider sequence of populations, $\{N_\gamma\}_{\gamma \in \mathbb{N}}$.
- Increasing in size as $\gamma \uparrow$.
- $m_\gamma(y) := \mathbb{E}(I_{\gamma i} | Y_i = y)$
- $c_\gamma(y_1, y_2) := \text{Cov}(I_{\gamma i}, I_{\gamma \ell} | Y_i = y_1, Y_\ell = y_2)$
- $w_\gamma(y) := \frac{1}{m_\gamma(y)}$
- Assume $m_\gamma(y) \xrightarrow{\text{a.s.}} m(y)$
- Assume $\forall y_1, y_2, c_\gamma(y_1, y_2) = o_\gamma(1)$

$$\frac{w_\gamma}{\mathbb{E}_o(w_\gamma(y))} f_o(y|\boldsymbol{\lambda}) \xrightarrow{\text{a.s.}} f_p(y|\boldsymbol{\lambda})$$

- Bonn ery et al. (2013) extends Glivenko-Cantelli

The Procedure

- Generate B pseudo-populations, each of size n , by **weighted sampling** of observed data **with replacement** using normalized inverse inclusion probabilities
- Estimate $(\lambda_b)_{b=1,\dots,B}$ under each pseudo-population
- Concatenate estimates, $\lambda = (\lambda_1, \dots, \lambda_B)$

Empirical Generation of Pseudo-populations

The Idea

- **Weighted** empirical predictive distribution, $\hat{f}(\mathbf{y}_p|\mathbf{y}_o)$, generates “pseudo-populations” (Pfeffermann et al., 1998, Pfeffermann and Sverchkov, 2003, Sverchkov and Pfeffermann, 2004), applies result of Bonn ery et al. (2013)
- \mathbf{y}_p are treated as parameters

$$\hat{f}(\boldsymbol{\lambda}|\mathbf{y}_o) = \int f(\mathbf{y}_p|\boldsymbol{\lambda}) \hat{f}(\mathbf{y}_p|\mathbf{y}_o) f(\boldsymbol{\lambda}) d\mathbf{y}_p = \int c(\mathbf{y}_p) f(\boldsymbol{\lambda}|\mathbf{y}_p) \hat{f}(\mathbf{y}_p|\mathbf{y}_o) d\mathbf{y}_p.$$

- Full uncertainty accounted for in estimation of $\boldsymbol{\lambda}$
 - population generation, $\hat{f}(\mathbf{y}_p|\mathbf{y}_o)$
 - $\boldsymbol{\lambda}$ estimation given the finite population, $\hat{f}(\boldsymbol{\lambda}|\mathbf{y}_p)$
 - conditioned on \mathbf{y}_o

Substitution “Plug-in” Method

- $f_o(\lambda|H(\mathbf{Y}_o)) \equiv$ sample model
- Use (\tilde{w}_i) to adjust $H(\cdot)$ from sample to population
- $\widehat{f}_p(\lambda|H(\mathbf{Y}_o, \tilde{\mathbf{w}})) \equiv$ population model

- Statistic, $\widehat{f}_p(\lambda|\mathbf{H}(\mathbf{y}_o, \tilde{\mathbf{w}})) = \prod_{j=1}^p \widehat{f}_p(\lambda_j|\mathbf{H}(\mathbf{y}_o, \tilde{\mathbf{w}}, \lambda_{-j}))$

$$\log \widehat{f}_p(\theta_{qm}^* | H(\{\mathbf{y}_i\}_{i \in s_m}, \{\tilde{w}_i\}_{i \in s_m}, \tau_\epsilon, \boldsymbol{\theta}_m^*))$$

$$\propto -\frac{1}{2} n_m \log(|\mathbf{C}^\tau(\boldsymbol{\theta}_m^*)|) - \frac{1}{2} \sum_{i \in s_m} \mathbf{y}'_{o,i} \mathbf{C}^\tau(\boldsymbol{\theta}_m^*) \mathbf{y}_{o,i} \tilde{w}_i + \log f(\theta_{qm}^*)$$

- $\mathbf{C}^\tau(\boldsymbol{\theta}_m^*) = \mathbf{C}(\boldsymbol{\theta}_m^*) + (1/\tau_\epsilon) \mathbb{I}_T$
- s_m collects the units assigned to cluster m
- $n_m = \sum_{i \in s_m} \tilde{w}_i$, $\sum_{m=1}^M n_m = n$

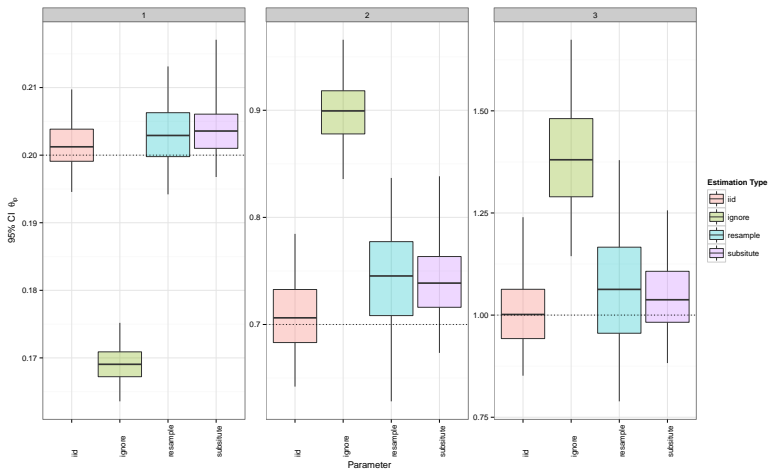
Comments on using Sampling Weights

- Both re-sampling and substitution don't require to re-parameterize
- This estimator employs harmonic averages for full conditionals
 - May break-down under high variation in the weights
- “Gold standard” to jointly model (\mathbf{y}_i, π_i)

A Warm Up

- Generate $N = 10000$ population units
- $\mathbf{f}_i \sim \mathcal{N}_t(\mathbf{0}, \mathbf{C}(\boldsymbol{\theta}))$
- $\boldsymbol{\theta} = (0.2, 0.7, 1.0)$ are global, not indexed by unit - No clustering
- Generate $y_{ij} = \mathcal{N}(f_{ij}, \tau_\epsilon^{-1})$
- τ_ϵ chosen to achieve noise-to-signal of 0.1
- Employ $H = 4$ strata
- Assign units based on quantiles of variances for $N, T \times 1$ functions, \mathbf{y}_i
- Idea is that more interested in high magnitude (error) functions
- (STSI) Sample $n = 750$ from $H = 4$ strata with,
 - $\mathbf{n} = (82, 113, 211, 344), \mathbf{p} = (0.03, 0.05, 0.08, 0.14)$
- Generate $B = 25$ pseudo-populations for re-weighted estimation

$\theta = (0.2, 0.7, 1.0)$ and $n = (82, 113, 211, 344)$



2— stage Sampling Design

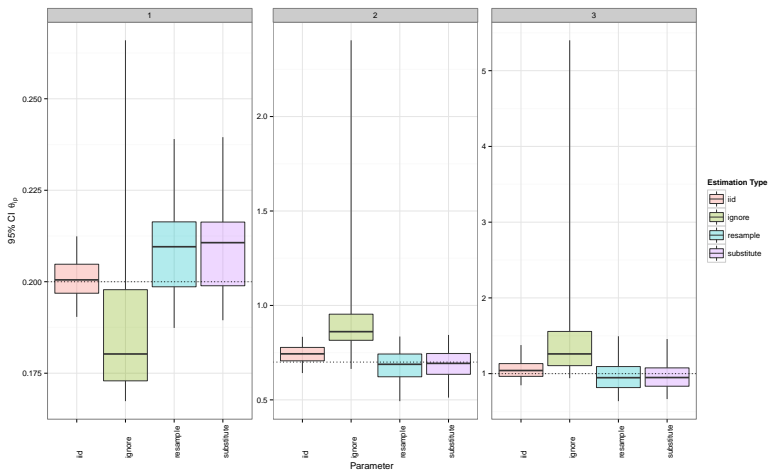
- Form $n_I = 10$ homogenous blocks by variance quantiles (of y_i)
- 1st stage: Sample 7 blocks proportional to variance
 - $p = (0.04, 0.06, 0.07, 0.08, 0.09, 0.10, 0.11, 0.12, 0.14, 0.19)$
- 2nd stage: $H = 4$ strata based on variance quantiles within blocks

	block/stratum	$H=1$	2	3	4
$p_I(\mathbf{s}_I) * p_i(\mathbf{s}_{Ii})' =$	1	0.001	0.002	0.004	0.006
	2	0.002	0.003	0.006	0.009
	3	0.002	0.004	0.007	0.010
	4	0.003	0.004	0.008	0.012
	5	0.003	0.005	0.009	0.013
	6	0.004	0.005	0.010	0.015
	7	0.004	0.006	0.011	0.017
	8	0.004	0.006	0.012	0.019
	9	0.005	0.007	0.014	0.022
	10	0.007	0.010	0.019	0.029

- $n = 750$ units from $N = 10000$ population

2— stage Posterior Distributions for $\theta = (\theta_1, \theta_2, \theta_3)$

- Monte Carlo draws of 10 samples, each selecting 7 of 10 blocks
- “Most” population elements have non-zero selection probabilities
- Conditional inclusion dependence is low



Generate Functions Under $M = 3$ Clusters

- Generate $N = 10000$ population units
- Here, $\{ \theta_i \}_{i=1, \dots, N}^{P \times 1}$
- 3 clusters with **location** values, $\Theta^* = (\theta_1^*, \dots, \theta_M^*)$:

parameter/cluster	$m=1$	2	3
$\Theta^* =$			
$p=1$	0.30	0.30	0.30
2	0.30	2.50	0.70
3	0.80	1.50	2.00
	rough-varied	smooth	rough

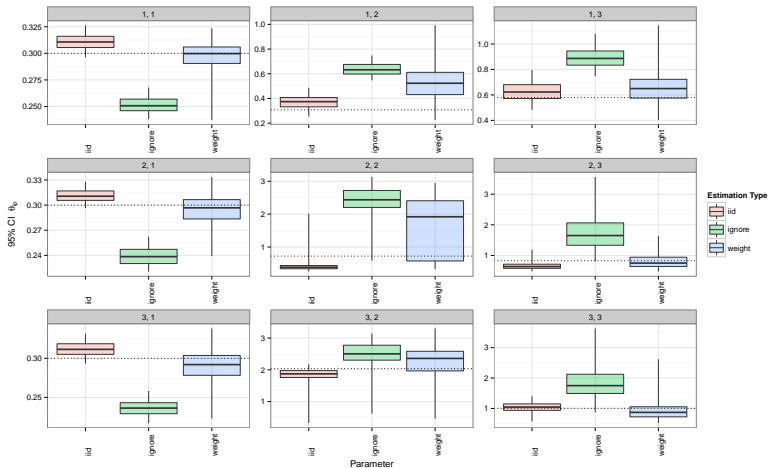
- $\mathbf{s} = s_1, \dots, s_N$ indexes randomly-assigned cluster memberships
- $s_i \in (1, \dots, M = 3), i = 1, \dots, N$
- $\theta_i \equiv \theta_{s_i}^*$

Generate Functions Under $M = 3$ Clusters

- $\mathbf{f}_i \sim \mathcal{N}_t(\mathbf{0}, \mathbf{C}(\boldsymbol{\theta}_{s_i}^*))$
- Generate $y_{ij} = \mathcal{N}(f_{ij}, \tau_\epsilon^{-1})$
- τ_ϵ chosen to achieve noise-to-signal of 0.1
- Employ $H = 4$ strata
- Assign units based on quantiles of variances for $N, T \times 1$ functions, \mathbf{y}_i
- (STSI) Sample $n = 800$ from $H = 4$ strata with unequal probabilities, $\mathbf{p} = (0.03, 0.06, 0.07, 0.16)$, $\mathbf{n} = (75, 155, 178, 392)$
- Generate $B = 25$ pseudo-populations for re-weighted estimation

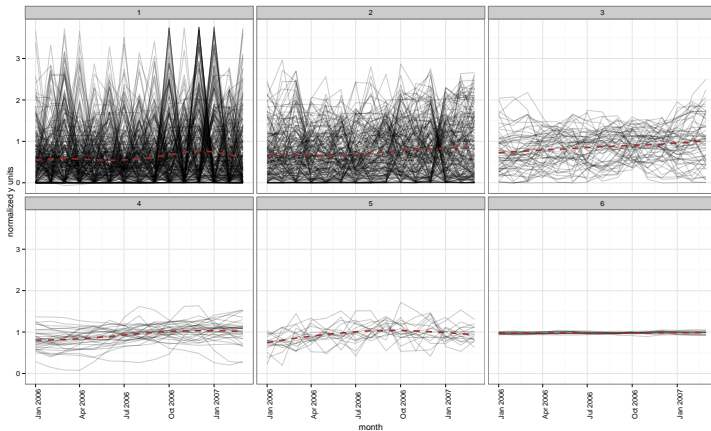
Posterior Distributions for $\theta_i = (\theta_{i,1}, \theta_{i,2}, \theta_{i,3})$

- One unit randomly chosen within each cluster
- Non-informative performs worst on members of smooth cluster
- Does best on rough cluster - over-samples \uparrow variance functions



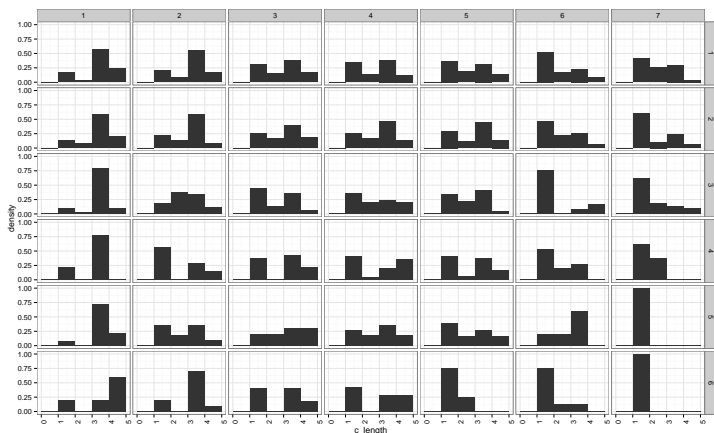
Clustering of $\{f_i\}$ for RAS Data

- 6 clusters contain (950, 648, 204, 106, 57, 48) establishments
- Based on relative magnitude of errors committed than on shape



Closing Period by Size \times Cluster

- 7 employee size categories,
(1 – 9, 10 – 19, 20 – 49, 50 – 99, 100 – 249, 250 – 499, > 499)
- BLS defines four time-indexed closing periods
- larger sized establishments in lower error clusters report sooner



References

- Bonnéry, D., Breidt, F. J., and Coquet, F. (2013). Uniform convergence of the empirical cumulative distribution under informative selection from a finite population. Technical report, Submitted to Bernoulli.
- Pfeffermann, D., Krieger, A., and Rinott, Y. (1998). Parametric distributions of complex survey data under informative probability sampling. *Statistica Sinica* 8, 1087-1114 (1998).
- Pfeffermann, D. and Sverchkov, M. (2003). *Fitting Generalized Linear Models under Informative Sampling*, pages 175–195. John Wiley & Sons, Ltd.
- Sverchkov, M. and Pfeffermann, D. (2004). Prediction of finite population totals based on the sample distribution. *Survey Methodology*, 30(1):79–82.