

# Plan de sondage informatif, distribution pondérée, et maximum de vraisemblance.

D. Bonnéry\*, F. Coquet\*, J. Breidt\*\*

\* Crest (Ensaï)

\*\* Colorado State University

5 novembre 2012

## ① Informative selection and asymptotic framework

- Informative selection mechanism
- Sample distribution

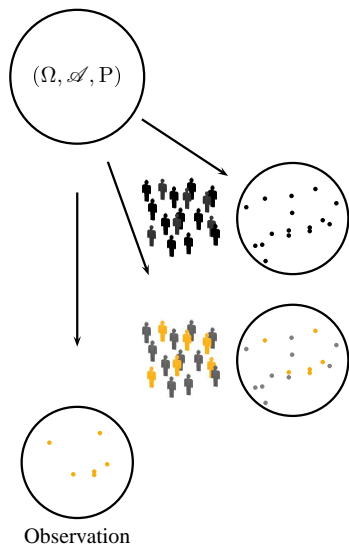
## ② Parametric estimation

- Maximum pseudo-likelihood estimator
- Convergence
- Simulations

An observation is the **outcome of two random processes**:

- The **population realization**: variables are generated for each element of a population  $U$  of size  $N$ .
- The **selection mechanism** : a **sample** is drawn from  $U$ .

The observation usually consists of the values of the study variables for each element in the sample.



Define:

- $(Y_1, \dots, Y_N)$ : **study variable**,
- $(Z_1, \dots, Z_N)$ : **design variables**,
- $\Pi$ : **design measure** and symmetric function of  $(Z_1, \dots, Z_N)$ ,
- $(J_1, \dots, J_N)$ : **sample** (vector of  $\mathbb{N}^N$ ),

$$P^{\Pi, \mathcal{Y}, \mathcal{Z}} \text{ -a.s. } (p, y, z), P^{\mathcal{J} | \Pi=p, \mathcal{Z}=z, \mathcal{Y}=y} = p.$$

- $\pi_k = \Pi(\{J_k \geq 1\})$ : **inclusion probability**,
- $n = \sum_{k \in U} J_k$ : **sample size**,

Assume:

- $\lim_{N \rightarrow \infty} N^{-1}n > 0$ .

## Definition

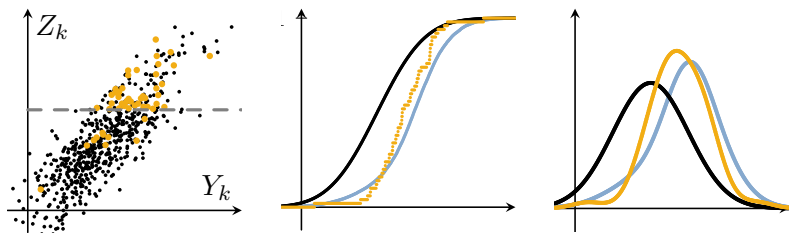
The weight function is:

$$\rho : y \mapsto \lim_{N \rightarrow \infty} \mathbb{E} [J_k | Y_k = y] / \mathbb{E} [J_k]$$

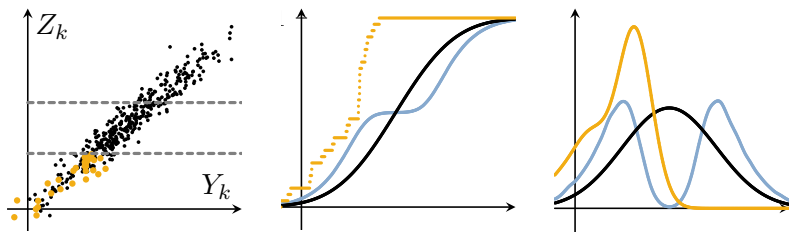
Denote  $(\ell)$  the  $\ell$ th drawn element.

$$\text{Does } P^{(Y_{(\ell)})_{\ell=1}^n} \sim (\rho \cdot P^{Y_1})^{\otimes n}?$$

- Observations behave like iid  $\rho.f$



- Observations do not behave like iid  $\rho.f$



Population cdf,  
empirical sample  
cdf and  
sample cdf

$f$ ,  $\rho.f$  and  
sample kernel  
density estimator

## Results:

Under asymptotic independence of draws in the one-dimensional case:

- the empirical sample cdf converges to  $\alpha \mapsto (\rho \cdot P^Y)((-\infty, \alpha])$  (Bonnéry, Breidt and Coquet, 2011)
- a kernel density estimator (kde) of the pdf converges to

$$\rho \cdot \frac{dP^Y}{d\lambda}$$

## Goal:

- Parametric estimation.

# Inference on population model



Consider the population model

- $(Y_k)_{k \in \{1, \dots, N\}} \sim (f_{\theta \cdot \lambda})^{\otimes N}$ ,
- $(Y_k, Z_k)_{k \in \{1, \dots, N\}}$  are iid realizations,
- $P^{Z_k|Y_k}$  is parametrized by  $\xi \in \Xi$

The target of the inference is  $\theta$ .

### Definition

$$\rho_{\theta, \xi} : y \mapsto \lim_{N \rightarrow \infty} \mathbb{E}_{\theta \xi} [J_k | Y_k = y] / \mathbb{E}_{\theta \xi} [J_k].$$

Following Pfeffermann and Krieger (1992) we define:

$$\hat{\theta}(\xi) = \arg \max_{\theta \in \Theta} \left\{ \sum_{\ell=1}^n \ln (\rho_{\theta \xi} f_{\theta} (Y_{(\ell)})) \right\}.$$

Assume A0. Let

- $\hat{\xi}$  be a consistent estimator of  $\xi$ ,
- $\bar{\mathcal{L}} \left( (Y_{(\ell)})_{\ell \in \{1, \dots, n\}}, \theta, \xi \right) = n^{-1} \sum_{k=1}^n \ln (\rho_{\theta \xi} f_{\theta} (Y_{(\ell)}, \theta, \xi))$

### Definition

The maximum pseudo-likelihood estimator associated to  $\hat{\xi}$  is:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \left\{ \bar{\mathcal{L}} \left( (Y_{(\ell)})_{\ell \in \{1, \dots, n\}}, \theta, \hat{\xi} \right) \right\},$$

## Definition

Define

$$m(y) = \mathbb{E} [J_1 | Y_1 = y]$$

$$m'(y_1, y_2) = \mathbb{E} [J_2 | Y_1 = y_1, Y_2 = y_2]$$

$$v(y) = \text{Var} [J_1 | Y_1 = y]$$

$$c(y_1, y_2) = \text{Cov} [J_1, J_2 | Y_1 = y_1, Y_2 = y_2]$$

$$\delta(y_1, y_2) = m'(y_1, y_2)m'(y_2, y_1) - m(y_1)m(y_2)$$

$$Y^* \sim \rho_{\theta\xi} f_{\theta} \cdot \lambda$$

$$\mathcal{I}_{11} = \mathbb{E}_{\theta, \xi} \left[ \left( \frac{\partial \ln(\rho_{\theta\xi} f_{\theta})}{\partial \theta} (Y^*, \theta, \xi) \right)^2 \right]$$

$$\mathcal{I}_{12} = \mathbb{E}_{\theta, \xi} \left[ \left| \left( \frac{\partial}{\partial \theta} \ln(\rho_{\theta\xi} f_{\theta}) \frac{\partial}{\partial \xi} \ln(\rho_{\theta\xi} f_{\theta}) \right) (Y^*, \theta, \xi) \right| \right]$$

## Assumption

- *Standard conditions on  $\rho_{\theta,\xi}f$ ,*
- *asymptotic independance of draws:*

$$\forall g \in \left\{ \mathbf{1}, \left( \frac{\partial \ln(\rho_{\theta,\xi}f_{\theta})}{\partial \theta}(\cdot, \theta, \xi) \right)^2, \left( \frac{\partial \ln(\rho_{\theta,\xi}f_{\theta})}{\partial \theta}(\cdot, \theta, \xi) \frac{\partial \ln(\rho_{\theta,\xi}f_{\theta})}{\partial \xi}(\cdot, \theta, \xi) \right) \right\},$$

- $E_{\theta\xi} [ |g(Y_1) g(Y_2)| c_{,\theta,\xi}(Y_1, Y_2) ] = o_{N \rightarrow \infty}(1),$
- $E_{\theta\xi} [ |g(Y_1) g(Y_2)| \delta_{,\theta,\xi}(Y_1, Y_2) ] = o_{N \rightarrow \infty}(1),$
- $E_{\theta\xi} [ (g^2(v_{\theta,\xi} + m_{\theta,\xi}^2))(Y_1) ] = o(N).$

## Theorem

*Suppose*

$$\sqrt{n} \begin{bmatrix} \left( \frac{\partial}{\partial \theta} \bar{\mathcal{L}} \right) \left( (Y_{(\ell)})_{\ell \in \{1, \dots, n\}}, \theta, \xi \right) \\ \hat{\xi} - \xi \end{bmatrix} \xrightarrow{\mathcal{L}} \mathcal{N} \left( 0, \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{22} \end{bmatrix} \right).$$

*Then*

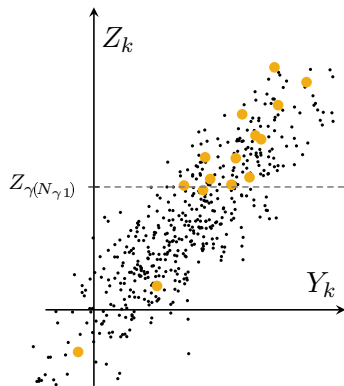
$$\sqrt{n} (\hat{\theta} - \theta) / \sigma \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1),$$

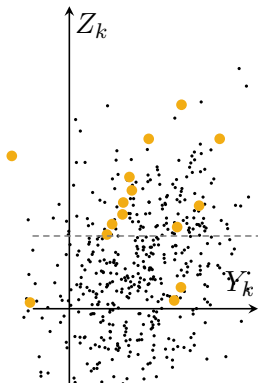
*with*

$$\sigma^2 = \frac{\Sigma_{11}}{\mathcal{J}_{11}^2} + \frac{\mathcal{J}_{12}}{\mathcal{J}_{11}^2} (\Sigma_{22} \mathcal{J}_{12} - 2\Sigma_{12}).$$

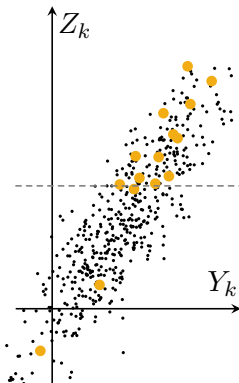
Consider:

- $Y \sim \mathcal{N}(\theta.1, Id_N)$ ,
- $\varepsilon \sim \mathcal{N}(0, Id_N)$ ,  $\varepsilon$  and  $Y$  are independent,
- $Z = \xi.Y + \eta\varepsilon$ ,  $\eta$  known,
- $N = 5000$ ,
- $\frac{N_1}{N} = 0.7$ ,  $\frac{N_2}{N} = 0.3$ ,
- $\frac{n_1}{N_1} = 1/70$ ,  $\frac{n_2}{N_2} = 4/30$ .

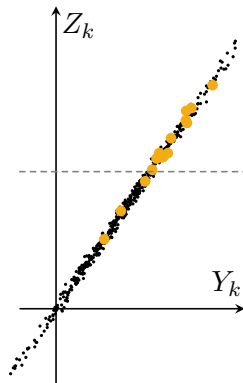




$\eta = 10$



$\eta = 1$



$\eta = 0.1$

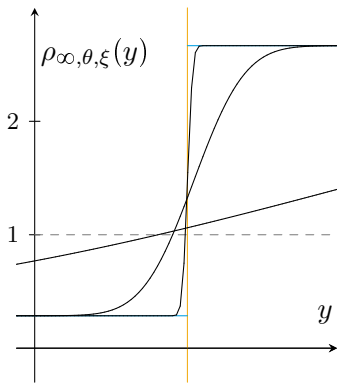
Let

- $t_1 = \lim_{N \rightarrow \infty} \frac{N_1}{N} = 0.7$
- $\zeta = \phi^{-1}(t_1)$
- $\tau_h = \lim_{N \rightarrow \infty} \left( \frac{n_h}{N_h} \right)$ ,  
 $\tau_1 = 1/70$ ,  $\tau_2 = 4/30$
- $p(y) =$   
$$P \left( \varepsilon < \frac{\zeta \sqrt{\xi^2 + \eta^2} + \xi(\theta - y)}{\eta} \right)$$

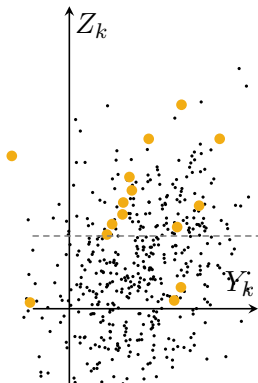
We get:

$$\rho_{\theta, \xi}(y) = \frac{\tau_1 p(y) + \tau_2 (1 - p(y))}{\tau_1 t_1 + \tau_2 (1 - t_1)}$$

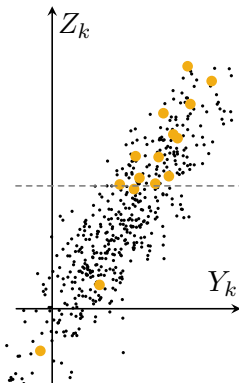
Figure: Plot of  $\rho$  for  $\theta = 1.5$ ,  
 $\xi = 2$ ,  $\eta \in \{.1, 1, 10\}$



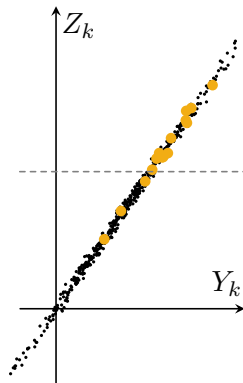




$\eta = 10$



$\eta = 1$



$\eta = 0.1$

To estimate  $\xi$ , we use

$$\hat{\xi} = \frac{\sum_{\ell=1}^n Z_{(\ell)} Y_{(\ell)} / \pi_{(\ell)}}{\sum_{\ell=1}^n Y_{(\ell)}^2 / \pi_{(\ell)}}.$$

Compare  $\hat{\theta}$  to

- $\tilde{\theta} = \sum_{\ell=1}^n \frac{Y_{(\ell)}}{\pi_{(\ell)}} = \arg \max_{\theta \in \Theta} \left\{ \sum_{k=1}^N \frac{\ln(f_{\theta}(Y_k)) J_k}{\pi_k} \right\}$ .,
- $\bar{\theta} = n^{-1} \sum_{\ell \in \{1, \dots, n\}} Y_{(\ell)}$ .

Table: Calculus of mean and mean square error on 1000 simulations

$\theta$	$\xi$	$\sigma$		Mean[.]	MSE[.]	$\sqrt{\frac{\text{MSE}}{\text{MSE}(\hat{\theta})}}$	$\frac{1}{n_\gamma} \lim_{\gamma \rightarrow \infty} n_\gamma \text{Var} [.]$
1.5	2	0.1	$\hat{\theta}$	1.502	$7.643 \cdot 10^{-4}$	1	$6.962 \cdot 10^{-4}$
			$\tilde{\theta}$	1.5	$4.811 \cdot 10^{-3}$	2.509	$4.523 \cdot 10^{-3}$
			$\bar{\theta}$	2.329	$6.887 \cdot 10^{-1}$	30.02	$3.979 \cdot 10^{-3}$
1.5	2	1	$\hat{\theta}$	1.5	$1.975 \cdot 10^{-3}$	1	$2.975 \cdot 10^{-3}$
			$\tilde{\theta}$	1.501	$5.583 \cdot 10^{-3}$	1.681	$6.024 \cdot 10^{-3}$
			$\bar{\theta}$	2.241	$5.509 \cdot 10^{-1}$	16.7	$3.971 \cdot 10^{-3}$
1.5	2	10	$\hat{\theta}$	1.497	$5.501 \cdot 10^{-3}$	1	$2.943 \cdot 10^{-3}$
			$\tilde{\theta}$	1.5	$1.030 \cdot 10^{-2}$	1.368	$1.030 \cdot 10^{-2}$
			$\bar{\theta}$	1.662	$2.999 \cdot 10^{-2}$	2.335	$4.027 \cdot 10^{-3}$

## Summary:

- Definition of sample pdf and limit sample pdf, applicable to with or without replacement and fixed or random size samples,
- Simple and verifiable conditions on the sequence of sample schemes,
- Pertinence: The sample behaves like an independent sample (Uniform cdf convergence),
- The limit sample pdf can be used for inference (Convergence of the maximum pseudo likelihood estimator),
  - Asymptotic normality under stratified sampling and fixed number of strata.

- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3).
- Sugden, R. A. and Smith, T. M. F. (1984). Ignorable and informative designs in survey sampling inference. *Biometrika*, 71(3).

- Bonnéry, D., Breidt, F. J., and Coquet, F. (2011). Uniform convergence of the empirical cumulative distribution function under informative selection from a finite population. *To appear in Bernoulli*.
- Pfeffermann, D. and Krieger, A. M. (1992). Maximum likelihood estimation for complex sample surveys. *Survey Methodology*, 18(1):225–255.
- Gong, G. and Samaniego, F. J. (1981). Pseudomaximum likelihood estimation: theory and applications. *The Annals of Statistics*, 9(4).
- Pfeffermann, D., Krieger, A. M., and Rinott, Y. (1998a). Parametric distributions of complex survey data under informative probability sampling. *Statistica Sinica*, 8(4).

Thank you for your attention.