

# ROBUSTNESS AND HIGH DIMENSIONAL DATA

Peter J. Bickel

UC Berkeley

MSU 9/2012

(Joint with Boaz Nadler, Bin Yu, N. el Karoui, Derek Bean and  
Chingway Lim)

# Outline

- 1 Robust  $M$  estimation in Linear Regression for fixed number of covariates  $p$
- 2 What is known if,  $\frac{p}{n} \rightarrow 0$ ,  $p \rightarrow \infty$ ?
- 3 Least squares and Lasso: Some current results
- 4 Some curious simulations
- 5 Heuristics
- 6 Projection Pursuit
- 7 Discussion

# The Regression Model and Least Squares Data

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$$X_i = (\mathbf{Z}_i, Y_i) \quad i = 1, \dots, n \quad \mathbf{Z}_i \quad p \times 1 \text{ iid}$$

**Assumed model:**  $(n = 1)$

$$Y = \mathbf{Z}^T \beta_0 + \mathbf{e}$$

$$\mathbf{e} \perp \mathbf{Z}$$

Used as an approximation to general model

$$Y = \mu(\mathbf{Z}) + e, \quad E(e|\mathbf{Z}) \equiv 0$$

# Basic Theorem

If  $\hat{\beta} = \arg \min \{ \|Y - \mathbf{Z}^T \beta\|_n^2 \}$  where  $\|f(X)\|_n \equiv \frac{1}{n} \sum_{i=1}^n f^2(X_i)$

a)  $\hat{\beta} = \left[ \frac{1}{n} (\mathcal{Z}^{(n)} [\mathcal{Z}^{(n)}]^T) \right]^{-1} (\mathbf{Z}^{(n)}, Y)_{(n)}$

where  $\mathbf{Y} \equiv (Y_1, \dots, Y_n)^T$ ,  $(\mathbf{Z}^{(n)}, \mathbf{Y})_{(n)} \equiv \frac{1}{n} \sum_{i=1}^n Y_i \mathbf{Z}_i$

$$\mathcal{Z}_{p \times n}^{(n)} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)$$

$$\mathcal{Z}^{(n)} [\mathcal{Z}^{(n)}]^T = \sum_{i=1}^n \mathbf{Z}_i \mathbf{Z}_i^T$$

b) If  $\Sigma \equiv E(\mathbf{Z}\mathbf{Z}^T)$  is nonsingular,  $\beta_0$  is TRUE

$$\sqrt{n}(\hat{\beta} - \beta_0) \implies N(\mathbf{0}, \sigma^2 \Sigma^{-1})$$

$$\beta_0 = \Sigma^{-1} [E(\mathbf{Y}\mathbf{Z})]$$

# Robust $M$ Estimation in Regression

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$$\beta_0 = \arg \min E_0 \rho(Y - \mathbf{Z}^T \beta)$$

$\rho$  convex, symmetric about 0.

$p$  fixed

$$\hat{\beta}_\rho \equiv \arg \min \frac{1}{n} \sum_{i=1}^n \rho(Y_i - \mathbf{Z}_i^T \beta)$$

**Thm (Huber)** If  $\psi \equiv \rho'$  is smooth,  $p$  is fixed,  $n \rightarrow \infty$ ,  
 $E_0 \psi^2(e) < \infty$ ,  $E_0 \psi'(e) \neq 0$  and  $\mathbf{Z}$  is full dimensional,  
 $E \mathbf{Z} \mathbf{Z}^T$  non singular

# Robust $M$ Estimation in Regression

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$$\hat{\beta}_\rho = \beta_0 - \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \frac{\psi(e)}{E_0 \psi'(e)} + o_P(n^{-\frac{1}{2}})$$

$$\sqrt{n}(\hat{\beta} - \beta_0) \implies N_p(\mathbf{0}, [E_0(\mathbf{Z}\mathbf{Z}^T)]^{-1} \sigma^2(\rho))$$

$$\sigma^2(\rho) = \frac{E_0 \psi^2(e)}{[E_0 \psi'(e)]^2}$$

E.g.:  $\psi(t) = t$  LSE not robust against heavy tails

$$\begin{aligned} \psi(t) &= h_k(t) = t, \quad |t| \leq k \text{ (Huber)} \\ &= k \operatorname{sgn}(t), \quad |t| > k \end{aligned}$$

$$\psi(t) = \operatorname{sgn}(t) \quad (L1)$$

# The current focus of interest: $p, n$ both large

What if  $p \rightarrow \infty$ ?

**Theorem** (Huber) (1973) (Negative)

If  $\frac{p}{n} \rightarrow c > 0$

$\exists$  contrast  $\mathbf{t}^T \beta_0$

$\mathbf{t}^T (\hat{\beta}_{\text{LSE}} - \beta_0)$  is not asymptotically Gaussian

**Note:**  $E[X^T (\hat{\beta}_{\text{LSE}} - \hat{\beta}_0)]^2 = \sigma^2 \frac{p}{n}$

$\implies$  Data picked contrast is inconsistent

# (Huber, Portnoy) (Positive)

(Huber) If the projection matrix  $[\mathcal{Z}^{(n)}]^T \{ \mathcal{Z}^{(n)} [\mathcal{Z}^{(n)}]^T \}^{-1} \mathcal{Z}^{(n)}$  has diagonal  $\pi_{ii} \equiv \frac{p}{n}$  and

$$\boxed{\frac{p^3}{n} \rightarrow 0}$$

$$a^T (\hat{\beta} - \beta) \sim N(0, \sigma^2(a, \psi))$$

$$\sigma^2(a, \psi) = \frac{E\psi^2(e)}{(E\psi'(e))^2} a^T [(\mathcal{Z}^{(n)} [\mathcal{Z}^{(n)}]^T)^{-1}] a$$

**Improved Conditions: Portnoy (1985) AS**



What if  $\frac{p}{n} \rightarrow 0$  more slowly or  $\frac{p}{n} \rightarrow c$ ,  $0 < c \leq \infty$

## Gaussian Linear Regression Model

$$\mathbf{Y}_{n \times 1} = \mathbf{Z}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \mathbf{e}_{n \times 1}$$

$$\mathbf{e} = (e_1, \dots, e_n)^T \text{ iid } N(0, \sigma^2)$$

$$\mathbf{z}^{(j)} \equiv \begin{pmatrix} Z_{1j} \\ \vdots \\ Z_{nj} \end{pmatrix}, \quad j = 1, \dots, p$$

$$\mathbf{Z} \equiv (\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(p)})_{n \times p} = [\mathcal{Z}^{(n)}]^T$$

Suppose  $|\mathbf{z}^{(j)}|^2 = n$ ,  $\mathbf{z}^{(a)} \perp \mathbf{z}^{(b)}$   $a \neq b$

(Canonical Gaussian Model)

Equivalent to:

## Gaussian White Noise Model (Donoho, Johnstone, Kerkyacharian, Picard (1995))

$$X_j = \beta_j + \varepsilon_j, \quad j = 1, \dots, p, \quad \varepsilon_j \sim N\left(0, \frac{\sigma^2}{n}\right) \text{ iid}$$

$$X_j = \frac{[\mathbf{Z}^{(j)}]^T \mathbf{Y}}{n}$$

Assume

- i)  $\beta$  sparse: If  $S = \{j; \beta_j \neq 0\}$ ,  $|S| = s \ll p$ .
- ii) Signal strong:  $j \in S \implies |\beta_j| \geq \delta_n > 0$

Let,

$$\hat{X}_j \equiv X_j - h_K(X_j)$$

$$h_K \equiv \text{Huber function, } K = \sigma \sqrt{\frac{2 \log p}{n}}$$

**GWN Result:** If  $\delta_n \sqrt{\frac{n}{\log p}} \rightarrow \infty$ ,

$$\sum_{j=1}^p E(\hat{X}_j - \beta_j)^2 = \frac{s\sigma^2}{n} (1 + o(1))$$

(Best possible if  $S$  is known)

If  $\delta_n = \Omega \sqrt{\frac{\log p}{n}}$ ,  $s \rightarrow s \log p$ .

# The Lasso: Donoho, Saunders, Chen (1998), Tibshirani (1996)

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$$\hat{\beta}_L \equiv \arg \min \{ |\mathbf{Y} - \mathbf{Z}\beta|^2 + \lambda |\beta|_1 \}$$

For canonical model

$$\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(p)} \text{ orthonormal } |\mathbf{Z}^{(j)}|^2 = n, j = 1, \dots, p .$$

Then, for suitable  $\lambda(K)$

$$\hat{\beta}_{jL} = \hat{X}_j .$$

# Conclusion

- a) If  $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(p)}$  are nearly orthogonal
- b)  $\beta$  is sparse
- c) The signal is strong

$\hat{\beta}_L$  behaves as LS when we know  $S$ .

**Many Results:** Buhlmann, van de Geer, Tsybakov, Meinshausen, Yu, Fan and collaborators, have found minimal versions of a)-c) extended GWN result.

# Robust Case:

Bradic, Fan, Wang, JRSS(B) (2011)

Variable selection including robust objective functions give results of this type with

$$\sigma^2 \frac{s}{n} \rightarrow \frac{E\psi^2(e)}{[E\psi'(e)]^2} \frac{s}{n}$$

Open problems in paralleling the work done for LS + Lasso but see GLM results of van de Geer and others (Buhlmann, van de Geer(to appear)) Statistics for High Dimensional Data

# Behavior of $\|\hat{\beta} - \beta_0\|^2$

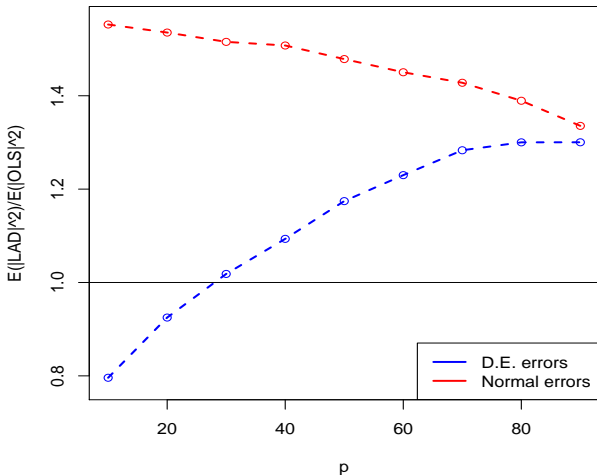
- A. What if a) or b) or c) conditions don't hold:  
*Another Lecture*
- B. What if  $\frac{p}{n} \rightarrow 0 < c < 1$  and robust and least squares are compared *without penalization*?

# Surprising simulations

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

Ratio of Expected Squared Norms,  $n=100$ , 1000 simulations





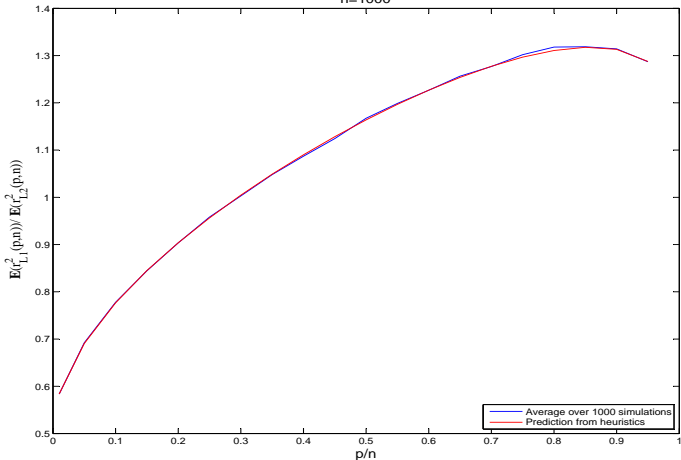
# Surprising simulations

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$E(r_{L_1}^2(p,n))/E(r_{L_2}^2(p,n))$  and  $r_{L_1}^2(\kappa)/r_{L_2}^2(\kappa)$  computed from system, double exponential errors, 1000 simulations

n=1000



# (Semi-heuristic) Results of el Karoui, Bean, Bickel, Lim and Yu

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$$Y_i = X_i^T \beta_0 + \epsilon_i \quad i = 1, \dots, n$$

- $\epsilon_i$  i.i.d.  $g \perp\!\!\!\perp \{X_i : i = 1, \dots, n\}$
- $X_i$  i.i.d.  $\mathcal{N}(0, \Sigma)$

Define  $\hat{\beta}(\rho; \beta_0, \Sigma) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n \rho(Y_i - X_i^T \beta)$

- $\rho$  convex
- $n \rightarrow \infty, \rho/n \rightarrow \kappa < 1$ .

# Key Lemma

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

$$\hat{\beta}(\rho; \beta_0, \Sigma) \stackrel{\mathcal{L}}{=} \beta_0 + \|\hat{\beta}(\rho; 0, I_p)\| \Sigma^{1/2} u,$$

where  $u$  is uniform on the  $p$  sphere of radius 1.

$\therefore$  Can assume  $\beta_0 = 0$ ,  $\Sigma = I_p$ .

# Special case of basic result

If  $r_\rho(\rho, n) \equiv \|\hat{\beta}(\rho; 0, I_\rho)\|$

Under suitable regularity conditions,

$r_\rho(\rho, n) \xrightarrow{P} r_\rho(\kappa)$  solving:

$$\mathbb{E} [\text{prox}_c(\rho)]'(\hat{z}_\epsilon) = 1 - \kappa$$

$$\mathbb{E} (\hat{z}_\epsilon - [\text{prox}_c(\rho)](\hat{z}_\epsilon))^2 = \kappa r_\rho^2(\kappa),$$

$$\hat{z}_\epsilon \stackrel{\mathcal{L}}{=} \epsilon + r_\rho(\kappa)Z, \quad \epsilon \perp\!\!\!\perp Z, \quad Z \sim \mathcal{N}(0, 1).$$

# Remarks

$$\text{prox}_c(\rho)(x) = \underset{y}{\operatorname{argmin}} \left( \rho(y) + \frac{(x - y)^2}{2c} \right)$$

Solves, if  $\rho$  is differentiable, strictly convex:

$$y + c\rho'(y) = x.$$

# Key ideas:

I. Leave out 1 predictor  $\{X_{pi} : i = 1, \dots, n\}$ .

$$r_{i,[p]} \equiv \epsilon_i - V_i^T \hat{\gamma}$$

where  $V_i \equiv (X_{1i}, \dots, X_{p-1,i})^T$ ,  $\hat{\gamma} \equiv$  estimate of  $(\beta_{01}, \dots, \beta_{0,p-1})^T$  without  $X^{(p)} \equiv (X_{p1}, \dots, X_{pn})^T$ . Then:

$$(*) \quad \hat{\beta}_p = \frac{\sum_i X_{pi} \psi(r_{i,[p]})}{\sum_i X_{pi}^2 \psi'(r_{i,[p]}) - v_p^T S_p^{-1} v_p} + o_p(n^{-1/2})$$

$$v_p = \sum_i \psi'(r_{i,[p]}) V_i X_{pi}$$

$$v_p^T S_p^{-1} v_p = [X^{(p)}]^T D^{1/2} \Pi_V D^{1/2} [X^{(p)}], \quad D_{ii} = \psi'(r_{i,[p]}).$$

$\Pi_V$  is a projection matrix of rank  $p - 1$ .

$$X^{(p)} \perp\!\!\!\perp r_{i,[p]}, V_i$$

# Key ideas:

For LS  $\psi(x) = x$  and (\*) holds exactly.

Asymptotically,  $r_{i,[p]} \sim g * \text{Gaussian}$

$\sqrt{n}\hat{\beta}_p \Rightarrow \mathcal{N}(0, \sigma^2(\rho, g, \kappa))$ .

$\sigma^2(\rho, g, \kappa) = \frac{\sigma^2}{1-\kappa}$  for LS.

II. Analysis requires leave out one  $(X_i, Y_i)$  as well.

# Projection Pursuit

(J) Kruskal (1969),(1972), Switzer (1970), Switzer and Wright (1971), Friedman Tukey (1974), *Huber* (1985), Diaconis and Freedman (1985)

Given:

$$X_1, \dots, X_n \quad p \times 1 \quad \text{iid}$$

Find “interesting” projections i.e.

$$a, |a| = 1 \ni$$

$$P_{n,a} \equiv \frac{1}{n} \sum_{i=1}^n \delta_{a^T X_i} \text{ is as non-normal as possible}$$



In expectation  $P_{n,a} \approx P_a \leftrightarrow f_a \equiv$  density of  $a^T X$

### Measures of Nonnormality

$$SK(P_{n,a}) \leftrightarrow \frac{E_a(X - E_a X)^3}{[E_a(X - E_a X)^2]^{\frac{3}{2}}} \equiv SK(P_a)$$

$$KURT(P_{n,a}) \leftrightarrow \frac{E_a(X - E_a X)^4}{[E_a(X - E_a X)^2]^2} - 3 \equiv K(P_a)$$

These are highly nonrobust to outliers.

# Robust and “efficient” measures

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

**Alternatives:** Robust Measures of Skewness , Kurtosis by Trimming

**Efficient Measure** (Estimate)

$$\int \log f_a f_a(x) dx + \log(2\pi e)^{\frac{1}{2}} [E_a(X - E_a X)^2]^{\frac{1}{2}}$$

**Procedure:** Maximize over  $a$

# A Rationale

ROBUSTNESS  
AND HIGH  
DIMEN-  
SIONAL  
DATA

Peter J. Bickel

Diaconis, Freedman (1985)

If  $p \rightarrow \infty$  with  $n \rightarrow \infty$  under weak conditions e.g.

$$\text{If } X_j \text{ iid } F, EX^2 < \infty, X = (X_1, \dots, X_p)^T,$$

Then, almost all  $P_{n,a}$  are asymptotically Gaussian, where “almost” all is with respect to Lebesgue measure on surface of unit sphere in  $R^p$ .

# Some Comfort

**Theorem** If  $F$  is Gaussian iid and  $\frac{p}{n} \rightarrow 0$ , then

$$\sup_a \sup_x |P_a(-\infty, x] - P_{na}(-\infty, x]| \xrightarrow{P} 0$$

# Some Caution

But, even if  $F$  is Gaussian iid, if  $\frac{p}{n} \rightarrow c > 0$ ,

$$\max_a \int a^T (\mathbf{x} - \mu(P_{n,a}))^2 dP_{n,a} \not\rightarrow 1$$

(Wigner, Geman)

# How bad can things get?

Suppose  $\frac{p}{n} \rightarrow \infty$

**Theorem** (B, Nadler)

Let  $X_1, \dots, X_n$  i.i.d.  $N(\mathbf{0}, I_p)$ . Let  $G$  be *any* cdf such that  $G - \Phi$  doesn't change sign.. Let  $\hat{F}_a$  denote the empirical cdf of  $a^T X_1, \dots, a^T X_n$ ,  $|a| = 1$ . Then:

$$P \left[ \inf_a \|\hat{F}_a - G\|_\infty \rightarrow 0 \right] = 1$$

where  $\|f\|_\infty = \sup_x |f(x)|$ .

# Idea of proof:

a) Let  $\psi : R \rightarrow R$  monotone increasing bounded.

$$\text{Let } \Psi_a \equiv \frac{1}{n} \sum_{i=1}^n \psi(a^T X_i)$$

$N = \lambda^p$  for  $\lambda$  to be chosen,  $\lambda > 1$ .

$A = \{a_j : 1 \leq j \leq N\}$  points in  $\mathcal{S}_p$

Then, for any  $\varepsilon > 0$ ,

$$P_{\Phi}[K_0 - \varepsilon \leq \Psi_{a_j} \leq K_0 + \varepsilon \text{ for some } j, 1 \leq j \leq N] \rightarrow 1$$

b) Let  $\psi^{(1)}, \dots, \psi^{(m)}$  be as above,  $\varepsilon > 0$ . For,

$$K_j \text{ arbitrary, } j = 1, \dots, m,$$
$$\text{sgn} \left( K_j - \int \psi^{(j)}(\xi) \phi(\xi) d\xi \right) \text{ constant}$$

Then,

$$P_\Phi [K_j - \varepsilon \leq \Psi_a^{(j)} \leq K_j, 1 \leq j \leq m \text{ for some } a \in \mathcal{A}] \rightarrow 1$$



c) Let  $\psi^{(j)}(u) = 1(x_j, \infty)$

$$K_j = \overline{G}(x_j)$$

By b)

$$P[\exists j \in \mathcal{A} \ni |\hat{F}_{a_j}(x_k) - \overline{G}(x_k)| \leq \epsilon \text{ for all } 1 \leq k \leq m] \rightarrow 1$$

# Lemma

$\exists \lambda > 1, \varepsilon > 0, a_1, \dots, a_N \in \mathcal{S}_p$ , such that, for  $N = \lambda^p$ ,

$\exists a_1, \dots, a_N$  so that  $|(a_j, a_{j'})| \leq 1 - \varepsilon$  for all  $1 \leq j \neq j' \leq N$ ,

Then,

$$(a_1^T X_1, \dots, a_N^T X_1)^T \sim \mathcal{N}_N(\mathbf{0}, R)$$

$$R = \|\rho_{ij}\|_{N \times N}, \text{ where } |\rho_{ij}| \leq 1 - \varepsilon, i \neq j$$

If  $X \sim N(\mathbf{0}, R_0)$

$$R_0 \equiv (1 - \varepsilon)\mathbf{1}\mathbf{1}^T + \varepsilon I_d$$

$$\mathbf{1} \equiv (1, \dots, 1)^T,$$

$$X = (1 - \varepsilon)\mathbf{z}_0\mathbf{1} + (1 - (1 - \varepsilon)^2)^{\frac{1}{2}}\mathbf{Z}$$

$$\mathbf{z}_0 \sim \mathcal{N}_1(0, 1) \perp \mathbf{Z} \sim \mathcal{N}_N(\mathbf{0}, I_d)$$

# Slepian's inequality

Extended by Joag-dev et al (1983) Ann. Prob.

Let  $\mathbf{Z}^{(j)} \sim \mathcal{N}(0, R^{(j)})$   $j = 0, 1$ ,

$$R_{N \times N}^{(j)} \equiv \rho_{ab}^{(j)} = \delta_{ab} + (1 - \delta_{ab})\rho_{ab}^{(j)},$$

and  $\rho_{ab}^{(0)} \leq \rho_{ab}^{(1)}$  for all  $a, b$

Let  $\Psi : R^N \rightarrow R$ , bounded,  
 $\frac{\partial^2 \psi}{\partial x_a \partial x_b} \geq 0$  all  $a \neq b$ . Then,

$$E\Psi(\mathbf{Z}^{(0)}) \leq E\Psi(\mathbf{Z}^{(1)})$$

(Valid if  $\Delta_{a,b}^2 \psi \geq 0$ , where

$$\Delta_{a,b}^2 \psi = \psi(x_a + h_a, x_b + h_b, x_c, c \neq a, b) - \psi(x_a + h_a, x_b, x_c, c \neq a, b) - \psi(x_a, x_b + h_b, x_c, c \neq a, b) + \psi(x_c, c = 1, \dots, N).$$

Let  $\psi_j$ ,  $j = 1, \dots, m$  be bounded non-decreasing function

Consider  $a^T X_1, \dots, a^T X_n$ ,  $a \in \mathcal{A}$ ,  $a_1, \dots, a_N$  as in Lemma.

Consider  $\mathcal{Y}_{m \times n}$  where  $\mathcal{Y}_{ij} = a_i^T \vec{X}_j$  and  $\mathcal{X}_c(\vec{\mathcal{Y}})$  where

$$\begin{aligned} \mathcal{X}_c(u_{ik} : i = 1, \dots, n, k = 1, \dots, N) \\ \equiv \prod_{k=1}^N \left[ 1 - \prod_{\ell=1}^m 1(\mathcal{X}^{(\ell)}(u_{1k}, \dots, u_{nk}) \geq c_\ell) \right], \end{aligned}$$

and  $\mathcal{X}^{(\ell)}(v_1, \dots, v_n) = \frac{1}{n} \sum_{i=1}^n \psi_\ell(u_i)$ .

$\mathcal{X}$  satisfies our hypotheses.

f) Apply large deviation theory to

$$\frac{1}{n} \sum_{i=1}^n \psi(Z_{ij}^{(0)})$$

where  $\mathbf{Z}_1^0, \dots, \mathbf{Z}_n^0$

$$\mathbf{Z}_i^{(0)} \equiv (Z_{i1}^{(0)}, \dots, Z_{iN}^{(0)})^T$$

are iid  $N_N(\mathbf{0}, R_{N \times N}^{(0)})$

$$R_{N \times N}^{(0)} = \|\delta_{ab} + (1 - \delta_{ab})(1 - \varepsilon)\|_{N \times N}$$

to obtain

$$P\left[\frac{1}{n} \sum_{i=1}^n Z_{ij}^{(0)} \notin [K_0 - \delta, K_0 + \delta] \text{ for any } 1 \leq j \leq N\right] \rightarrow 0$$

- g) Apply Slepian's inequality to get a).  
Generalize to b), c) using Joag-dev's inequality.  
Choose the  $\{x_j\}$  to be dense to get (\*)

# Discussion

## II: Huber (1985)

- 1) “Perhaps the practical conclusion to be drawn is that we shall have to acquiesce to the fact that PP will in practice reveal not only true but also spurious structure and that we must weed out the latter by other methods.”
- 2) What structures survive if we consider a random set of  $m$  projections of the data?  
E.g. Suppose the true population is

$$(1 - \varepsilon)N(\mathbf{0}, I_p) + \varepsilon N(\mathbf{0}, \Sigma)$$

$\Sigma$  of rank  $\ll p$

If we take  $m = o(n)$  projections, what chance do we have of finding  $N(0, \Sigma)$  structure?

- 3) Conjecture: Result holds for all  $G$ .



## I. Proof now available

### Questions:

- 1) “Optimal”  $\rho$ ,  $g$ ,  $\kappa$  known: DONE.
- 2) Robust  $\rho$ , optimizing on  $g$  given  $g$  in a small neighborhood of  $\phi$ : OPEN.