Chapter 5

Convexity

Very rough draft

# 5.1 Approximation of convex functions

Convexity:: approximation

Expand Convexity simplifies arguments from Chapter 3. Reasons: local minima are global minima and pointwise convergence of a sequnce of convex functions implies uniform convergence on compacta. The first Lemma in this Section contains the standard result (Rockafellar 1970, Theorem 10.8) that pointwise convergence of a sequence of (deterministic) convex functions  $\{H_n\}$  to a limit function B on an open set G implies uniform convergence on each compact subset of G.

For a subset H of  $\mathbb{R}^k$  and  $\delta > 0$ , write  $H^{\delta}$  for the set  $\{y : d(y, H) \leq \delta\}$ . For a bounded subset T of  $\mathbb{R}^k$ , call a finite subset  $T_0 \subset T$  a  $\delta$ -net for T if  $d(t, T_0) \leq \delta$  for ever t in T. Equivalently, the union of closed balls  $\cup_{t \in T_0} B(t, \delta)$  covers T. Also equivalently,  $B(t, \delta) \cap T_0 \neq \emptyset$  for each t in T. Define  $T^{\delta} := \{t \in \mathbb{R}^k : d(t, T) \leq \delta\}.$ 

bound <1> **Lemma.** (based on Pollard 1991) Let K be a compact, convex subset of  $\mathbb{R}^d$ . Let  $H(\cdot)$  and  $A(\cdot)$  be real-valued functions defined on  $K^{4\delta}$ , with H convex and A satisfying the uniform continuity requirement

 $|A(s) - A(t)| \le \eta$  whenever  $|s - t| \le 2\delta$  and  $s, t \in K^{4\delta}$ .

Let  $T_0$  be a  $\delta$ -net for  $K^{3\delta}$ . Then

$$\sup_{s \in K} |H(s) - A(s)| \le 4\eta + 3 \max_{t \in T_0} |H(t) - A(t)|$$

PROOF Write  $\Delta$  for the maximum of |H(t) - A(t)| over the finite set  $T_0$ . The asserted inequality will follow from a pair of bounds for A:

> version: 29sept2010 printed: 6 October 2010

Asymptopia ©David Pollard

#### Proof of upper bound (2).

Let s be a point of  $K^{2\delta}$ . First I prove that s lies in the convex hull of the finite set

$$T(s) := \{t \in T : |s - t| \le 2\delta\} = T_0 \cap B(s, 2\delta)$$

Argue by contradiction. If s is not in the convex hull then, by the separating hyperplane theorem, there exists some unit vector u for which

$$u's > u't$$
 for all  $t$  in  $T(s)$ 

The closed ball  $B(s^*, \delta)$  around the point  $s^* = s + \delta u$  lies inside  $B(s, 2\delta) \cap K^{4\delta}$ and it contains no points from  $T_0$ , which contradicts the  $\delta$ -net property of  $T_0$ .

It follows that s must be a convex combination of points in  $T_s$ . That is,  $s = \sum_{t \in T(s)} \alpha_t t$  with  $\alpha_t \ge 0$  and  $\sum_t \alpha_t = 1$ . Then

$$\begin{split} H(s) &\leq \sum_{t} \alpha_{t} H(t) & \text{by convexity of } H \\ &\leq \sum_{t} \alpha_{t} \left( A(t) + \Delta \right) & \text{definition of } \Delta \\ &\leq \sum_{t} \alpha_{t} \left( A(s) + \eta + \Delta \right) & \text{because } |s - t| \leq 2\delta \text{ for } t \in T(s) \\ &\leq A(s) + \eta + \Delta. \end{split}$$

**Proof of lower bound** <3>**.** 

Let s be a point of K. Find a point of  $T_0$  with  $|s - t| \leq \delta$ . Define  $s^* = t - (s - t)$ . Then  $s^* \in T(s) \subseteq K^{2\delta}$  and  $t = (s + s^*)/2$ . By convexity of H,

$$H(t) \le (H(s) + H(s^*))/2$$

From <2>,

$$H(s^*) \le A(s^*) + \eta + \Delta,$$

and, by definition of  $\Delta$ ,

$$H(t) \ge A(t) - \Delta$$

Combine the last three inequalities to get

$$2A(t) - 2\Delta \le H(s) + A(s^*) + \eta + \Delta$$

Both  $A(s^*)$  and A(t) lie within  $\eta$  of A(s). Inequality <3> follows.

2

### 5.2 Minimizers of convex functions

Convexity::argmin

Suppose M is a strictly convex function on  $\mathbb{R}^k$  satisfing the growth condition

growth < 4 >

$$M(t)/|t| \to \infty$$
 as  $|t| \to \infty$ .

For each z in  $\mathbb{R}^k$ , the function M(t, z) := M(t) - z't achieves its minimum value  $M^*(z) = \inf_{t \in \mathbb{R}^k} M(t) - z't$  at a unique point,  $\psi_M(z)$ .

I claim that the map  $z \mapsto \psi_M(z)$  is continuous. I would also like to show that if  $M(t^*, z) \leq M^*(z) + \epsilon$  then  $t^*$  must be close to  $\psi_M(z)$  in some suitably uniform sense. What else would be needed to establish a more general form of Lemma  $\langle 5 \rangle$ ?

### 5.3 A limit theorem for M-estimators

Convexity::limit

stoch.approximation <5>

**Lemma.** Suppose  $\{H_n(t) : t \in \mathbb{R}^k\}$  is a stochastic process with convex sample paths. Assume

(i)  $H_n(\hat{t}_n) \leq \inf_{t \in \mathbb{R}^k} H_n(t) + o_p(1)$  for some random  $\hat{t}_n$ 

(ii) there exists a sequence of random vectors  $\{Z_n\}$  of order  $O_p(1)$  for which

$$H_n(t) + Z'_n t \to \frac{1}{2} |t|^2$$
 in probability, for each fixed  $t \in \mathbb{R}^k$ 

Then  $\hat{t}_n = Z_n + o_p(1)$ .

Proof Define

$$A_n(t) = \frac{1}{2}|t|^2 - Z'_n t = \frac{1}{2}|Z_n - t|^2 - \frac{1}{2}|Z_n|^2,$$

so that  $|H_n(t) - A_n(t)| \to 0$  in probability, for each fixed t. Note that  $A_n$  is minimized at  $Z_n$ , which will imply that  $H_n$  is minimized near  $Z_n$ , with high probability.

Suppose  $\epsilon > 0$  and  $\eta > 0$  are given, small quantities. By (ii), there exists an R > 0 for which  $\mathbb{P}\{|Z_n| > R\} < \epsilon$  eventually. Choose a  $\delta < 1/4$  such that

$$|A_n(s) - A_n(t)| \le \eta^2 \quad \text{when } |s - t| \le 2\delta \text{ and } |Z_n| \le R \text{ and } s, t \in B(0, R + 1)$$

Let  $T_0$  be a  $\delta$ -net for B(0, R+1). Then there exists a set  $\Omega_n$  with probability eventually greater than  $1 - 2\epsilon$ , on which

$$|Z_n| \le R$$
  

$$\Delta_n := \max_{t \in T_0} |H_n(t) - A_n(t)| \le \eta^2$$
  

$$H_n(\hat{t}_n) \le H_n(Z_n) + \eta^2$$

On  $\Omega_n$ , Lemma <1> implies that

$$|H_n(t) - A_n(t)| \le 7\eta^2 \quad \text{for all } |t| \le R.$$

In particular, for all unit vectors v,

$$H_n(Z_n) - 7\eta^2 \le A_n(Z_n) = A_n(Z_n + 6\eta) - 18\eta^2 \le H_n(Z_n + 6\eta v) - 11\eta^2$$

Convexity of  $r \mapsto H_n(Z_n + rv)$  then implies that

 $H_n(Z_n) + 4\eta^2 \le H_n(Z_n + rv)$  for all  $r \ge 6\eta$  and all unt vectors v,

which forces  $\hat{t}_n$  to lie within the ball of radius  $6\eta$  around  $Z_n$ .

I would like to strengthen the Lemma by replacing (ii) by: there exists a sequence of random vectors  $\{Z_n\}$  of order  $O_p(1)$  for which

$$H_n(t) + Z'_n t \to M(t)$$
 in probability, for each fixed  $t \in \mathbb{R}^k$ ,

where M is a strictly convex function satisfying the growth condition  $\langle 4 \rangle$ . I would like to conclude something like

$$\hat{t}_n = \Psi_M(Z_n) + o_p(1)$$

If we had  $Z_n \rightsquigarrow Z$  then we would get  $\hat{t}_n \rightsquigarrow \psi_M(Z)$ , a very clean limit theorem.

random.sum <6>

**Theorem.** Let  $G_n(\theta) = \sum_{i \leq n} g_{n,i}(\theta)$  be a random process with convex sample paths defined on  $\mathbb{R}^k$ . Suppose there exist points  $\theta_n$  in  $\mathbb{R}^k$  and positive definite matrices  $J_n$  for which:

(i)  $\hat{\theta}_n$  is an estimator for which  $G_n(\hat{\theta}_n) \leq \inf_{\theta \in \Theta} G_n(\theta) + o_p(1)$ 

- (ii) the function  $\theta \mapsto \mathbb{P}G_n(\theta)$  is minimized at  $\theta_n$
- (iii) each  $\{g_{n,i}\}$  has a linear approximation near  $\theta_n$ ,

$$g_{n,i}(\theta) = g_{n,i}(\theta_n) + (\theta - \theta_n)' D_{n,i} + r_{n,i}(\theta)$$

where the  $D_{n,i}$  are random vectors with zero expected value

(iv)  $\mathbb{P}G_n(\theta_n + J_n^{-1/2}t) = \frac{1}{2}|t|^2 + o(1)$  for each t(v)  $var\left(\sum_{i \le n} r_{n,i}(\theta_n + J_n^{-1/2}t)\right) = o(1)$  for each t(vi)  $Z_n := -J_n^{-1/2} \sum_{i \le n} D_{n,i} = O_p(1)$ Then  $J_n^{1/2}(\hat{\theta}_n - \theta_n) = Z_n + o_p(1).$ 

Try to replace (iv) by  $\mathbb{P}G_n(\theta_n + J_n^{-1/2}t) = M(t) + o(1)$  for each t, with M as in Section 2. Maybe the Theorem will then be strong enough to cover the lasso (Example <17>).

PROOF The standardized random vector  $\hat{t}_n = J_n^{1/2}(\hat{\theta}_n - \theta_n)$  comes within  $o_p(1)$  of minimizing the process

$$H_n(t) := G_n(\theta_n + J_n^{-1/2}t) - G_n(\theta_n)$$
  
=  $-t'Z_n + \sum_{i < n} r_{n,i}(\theta_n + J_n^{-1/2}t).$ 

For a fixed t, assumption (v) ensures that the sum contributed by the  $\{r_{n,i}\}$  lies within  $o_p(1)$  of its expectation,

$$\mathbb{P}\sum_{i \leq n} r_{n,i}(\theta_n + J_n^{-1/2}t)$$
  
=  $\mathbb{P}G_n(\theta_n + J_n^{-1/2}t) - \mathbb{P}G_n(\theta_n)$  because  $\mathbb{P}D_{n,i} = 0$   
=  $\frac{1}{2}|t|^2 + o(1)$  by (iv)

Thus  $H_n(t) + t'Z_n = \frac{1}{2}|t|^2 + o_p(1)$  for each fixed t. An appeal to Lemma <??> completes the proof.

**Remarks.** The theorem is almost true if  $\theta$  ranges over only a subset  $\Theta$  of  $\mathbb{R}^d$ , instead of over the whole of  $\mathbb{R}^d$ .

# 5.4 Examples

# Convexity::examples

| median $<7>$                       | <b>Example.</b> Suppose $X_1, X_2, \ldots$ are independent observations from a probability distribution $P$ on the real line with sample median $\theta_0$ . Suppose also that there is a neighborhood of $\theta_0$ in which $P$ has a continuous, strictly positive density (with respect to Lebesgue measure). Let $\hat{\theta}_n$ be the sample median. Show that $2f(\theta_0)(\hat{\theta}_n - \theta_0) \rightsquigarrow N(0, 1)$ . |
|------------------------------------|---|
|                                    |   |
| spatial.median $\langle 8 \rangle$ | <b>Example.</b> Generalize Example $<7>$ to cover argmin $\sum_i  t - X_i $ in higher dimensions.   |
| LAD.regression $\langle 9 \rangle$ | <b>Example.</b> Theorem 1 of Pollard (1991) for L1 regression.  |
| $\exp.fam < 10>$                   | <b>Example.</b> Modify the following result from Dou, Pollard, and Zhou (2009), but with improved argument from talk at Wellner conference:   |
|                                    | Let $\{Q_{\lambda} : \lambda \in \mathbb{R}\}$ be an exponential family of probability measures with densities $dQ_{\lambda}/dQ_0 = f_{\lambda}(y) = \exp(\lambda y - \psi(\lambda))$ . Remember that $e^{\psi(\lambda)} = Q_0 e^{\lambda y}$ and that the distribution $Q_{\lambda}$ has mean $\psi^{(1)}(\lambda)$ and variance $\psi^{(2)}(\lambda)$ . We assume:  |
| psi3                               | $(\psi 3)$ There exists an increasing real function G on $\mathbb{R}^+$ such that   |
|                                    | $ \psi^{(3)}(\lambda+h)  \le \psi^{(2)}(\lambda)G( h )$ for all $\lambda$ and $h$   |
|                                    | Without loss of generality we assume $G(0) \ge 1$ .   |
| psi2                               | ( $\psi$ 2) For each $\epsilon > 0$ there exists a finite constant $C_{\epsilon}$ for which $\psi^{(2)}(\lambda) \leq C_{\epsilon} \exp(\epsilon \lambda^2)$ for all $\lambda \in \mathbb{R}$ . Equivalently, $\psi^{(2)}(\lambda) \leq \exp(o(\lambda^2))$ as $ \lambda  \to \infty$ .   |
|                                    | As shown in Section $\ref{eq:constraint},$ these assumptions on the $\psi$ function imply that  |
| hell < 11 >                        | $h^2(Q_{\lambda}, Q_{\lambda+\delta}) \le \delta^2 \psi^{(2)}(\lambda) \left(1 +  \delta \right) G( \delta )  \text{for all } \lambda, \delta \in \mathbb{R}.$  |
|                                    | <b>Remark.</b> We may assume that $\psi^{(2)}(\lambda) > 0$ for every real $\lambda$ . Otherwise we would have $0 = \psi^{(2)}(\lambda_0) = \operatorname{var}_{\lambda_0}(y) = \nu f_{\lambda_0}(y)(y - \psi^{(1)}(\lambda_0))^2$ for some $\lambda_0$ , which would make $y = \psi^{(1)}(\lambda_0)$ for $\nu$ almost all $y$ and $Q_{\lambda} \equiv Q_{\lambda_0}$ for every $\lambda$ .  |

The theory in this section combine ideas from Portnoy (1988) and from Hjort and Pollard (1993). We write our results in a notation that makes the applications in Section ?? and ?? more straightforward. The notational cost is that the parameters are indexed by  $\{0, 1, \ldots, N\}$ . To avoid an excess of parentheses we write  $N_+$  for N + 1. In the applications N changes with the sample size n and  $\mathbb{Q}$  is replaced by  $\mathbb{Q}_{n,a,\mathbb{B},N}$  or  $\mathbb{Q}_{n,a,\mathbb{B},N}$ .

Suppose  $\xi_1, \ldots, \xi_n$  are (nonrandom) vectors in  $\mathbb{R}^{N_+}$ . Suppose  $\mathbb{Q} = \bigotimes_{i \leq n} Q_{\lambda_i}$  with  $\lambda_i = \xi'_i \gamma$  for a fixed  $\gamma = (\gamma_0, \gamma_1, \ldots, \gamma_N)$  in  $\mathbb{R}^{N_+}$ . Under  $\mathbb{Q}$ , the coordinate maps  $y_1, \ldots, y_n$  are independent random variables with  $y_i \sim Q_{\lambda_i}$ .

The log-likelihood for fitting the model is

$$L_n(g) = \sum_{i \le n} (\xi'_i g) y_i - \psi(\xi'_i g) \quad \text{for } g \in \mathbb{R}^{N_+},$$

which is maximized (over  $\mathbb{R}^{N_+}$ ) at the MLE  $\widehat{g}$  (=  $\widehat{g}_n$ ).

**Remark.** As a small amount of extra bookkeeping in the following argument would show, we do not need  $\hat{g}$  to exactly maximize  $L_n$ . It would suffice to have  $L_n(\hat{g})$  suitably close to  $\sup_g L_n(g)$ . In particular, we need not be concerned with questions regarding existence or uniqueness of the argmax.

Define

(i) 
$$J_n = \sum_{i \le n} \xi_i \xi'_i \psi^{(2)}(\lambda_i)$$
, an  $N_+ \times N_+$  matrix

(ii)  $w_i := J_n^{-1/2} \xi_i$ , an element of  $\mathbb{R}^{N_+}$ 

(iii)  $W_n = \sum_{i \leq n} w_i \left( y_i - \psi^{(1)}(\lambda_i) \right)$ , an element of  $\mathbb{R}^{N_+}$ 

Notice that  $\mathbb{Q}W_n = 0$  and  $\operatorname{var}_{\mathbb{Q}}(W_n) = \sum_{i \leq n} w_i w_i' \psi^{(2)}(\lambda_i) = I_{N_+}$  and

$$\mathbb{Q}|W_n|^2 = \operatorname{trace}\left(\operatorname{var}_{\mathbb{Q}}(W_n)\right) = N_+.$$

mle.approx <12> Lemma. Suppose  $0 < \epsilon_1 \leq 1/2$  and  $0 < \epsilon_2 < 1$  and

$$\max_{i \le n} |w_i| \le \frac{\epsilon_1 \epsilon_2}{2G(1)N_+} \quad \text{with } G \text{ as in Assumption } (\psi 3).$$

Then  $\widehat{g} = \gamma + J_n^{-1/2} (W_n + r_n)$  with  $|r_n| \le \epsilon_1$  on the set  $\{|W_n| \le \sqrt{N_+/\epsilon_2}\}$ , which has  $\mathbb{Q}$ -probability greater than  $1 - \epsilon_2$ .

PROOF The equality  $\mathbb{Q}|W_n|^2 = N_+$  and Tchebychev give  $\mathbb{Q}\{|W_n| > \sqrt{N_+/\epsilon_2}\} \le \epsilon_2$ .

Reparametrize by defining  $t = J_n^{1/2}(g - \gamma)$ . The concave function

$$\mathcal{L}_n(t) := L_n(\gamma + J_n^{-1/2}t) - L_n(\gamma) = \sum_{i \le n} y_i w_i' t + \psi(\lambda_i) - \psi(\lambda_i + w_i' t)$$

is maximized at  $\hat{t}_n = J_n^{1/2}(\hat{g} - \gamma)$ . It has derivative

$$\dot{\mathcal{L}}_n(t) = \sum_{i \le n} w_i \left( y_i - \psi^{(1)}(\lambda_i + w'_i t) \right).$$

For a fixed unit vector  $u \in \mathbb{R}^{N_+}$  and a fixed  $t \in \mathbb{R}^{N_+}$ , consider the real-valued function of the real variable s,

$$H(s) := u'\dot{\mathcal{L}}_n(st) = \sum_{i \le n} u'w_i \left( y_i - \psi^{(1)}(\lambda_i + sw'_i t) \right),$$

which has derivatives

$$\dot{H}(s) = -\sum_{i \le n} (u'w_i)(w'_it)\psi^{(2)}(\lambda_i + sw'_it)$$
$$\ddot{H}(s) = -\sum_{i \le n} (u'w_i)(w'_it)^2\psi^{(3)}(\lambda_i + sw'_it).$$

Notice that  $H(0) = u'W_n$  and  $\dot{H}(0) = -u'\sum_{i\leq n} w_i w'_i \psi^{(2)}(\lambda_i)t = -u't$ . Write  $M_n$  for  $\max_{i\leq n} |w_i|$ . By virtue of Assumption ( $\psi$ 3),

$$\begin{split} |\ddot{H}(s)| &\leq \sum_{i \leq n} |u'w_i| (w'_i t)^2 \psi^{(2)}(\lambda_i) G\left(|sw'_i t|\right) \\ &\leq M_n G\left(M_n |st|\right) t' \sum_{i \leq n} w_i w'_i \psi^{(2)}(\lambda_i) t \\ &= M_n G\left(M_n |st|\right) |t|^2. \end{split}$$

By Taylor expansion, for some  $0 < s^* < 1$ ,

$$|H(1) - H(0) - \dot{H}(0)| \le \frac{1}{2} |\ddot{H}(s^*)| \le \frac{1}{2} M_n G(M_n|t|) |t|^2.$$

That is,

u'll<13>

$$\left| u' \left( \dot{\mathcal{L}}_n(t) - W_n + t \right) \right| \le \frac{1}{2} M_n G\left( M_n |t| \right) |t|^2.$$

Approximation <13> will control the behavior of  $\widetilde{\mathcal{L}}(s) := \mathcal{L}_n(W_n + su)$ , a concave function of the real argument s, for each unit vector u. By concavity, the derivative

$$\widetilde{\mathcal{L}}(s) = u' \dot{\mathcal{L}}_n(W_n + su) = -s + R(s)$$

is a decreasing function of s with

$$|R(s)| \le \frac{1}{2}M_n G(M_n|W_n + su|)|W_n + su|^2$$

On the set  $\{|W_n| \le \sqrt{N_+/\epsilon_2}\}$  we have

$$|W_n \pm \epsilon_1 u| \le \sqrt{N_+/\epsilon_2} + \epsilon_1$$

Thus

$$M_n|W_n \pm \epsilon_1 u| \le \frac{\epsilon_1 \epsilon_2}{2G(1)N_+} \left(\sqrt{N_+/\epsilon_2} + \epsilon_1\right) < 1,$$

implying

$$R(\pm\epsilon_1)| \leq \frac{1}{2}M_nG(1)|W_n \pm \epsilon_1 u|^2$$
$$\leq \frac{\epsilon_1\epsilon_2}{G(1)N_+} \left(N_+/\epsilon_2 + \epsilon_1^2\right)$$
$$\leq \epsilon_1 \left(1 + \epsilon_1^2\epsilon_2/N_+\right) < \frac{5}{8}\epsilon_1$$

Deduce that

$$\hat{\mathcal{L}}(\epsilon_1) = -\epsilon_1 + R(\epsilon_1) \le -\frac{3}{8}\epsilon_1$$
$$\hat{\mathcal{L}}(-\epsilon_1) = \epsilon_1 + R(-\epsilon_1) \ge \frac{3}{8}\epsilon_1$$

The concave function  $s \mapsto \mathcal{L}_n(W_n + su)$  must achieve its maximum for some s in the interval  $[-\epsilon_1, \epsilon_1]$ , for each unit vector u. It follows that  $|\hat{t}_n - W_n| \le \epsilon_1$ .

AnBn < 14>

**Corollary.** Suppose  $\xi_i = D\eta_i$  for some nonsingular matrix D, so that

$$J_n = nDA_nD \qquad where A_n := \frac{1}{n} \sum_{i \le n} \eta_i \eta'_i \psi^{(2)}(\lambda_i)$$

If  $B_n$  is another nonsingular matrix for which

 $assumption A{<}15{>}$ 

$$|A_n - B_n||_2 \le (2 ||B_n^{-1}||_2)^{-1}$$

and if

 $\max.eta < 16 >$ 

$$\max_{i \le n} |\eta_i| \le \frac{\epsilon \sqrt{n/N_+}}{G(1)\sqrt{32 \|B_n^{-1}\|_2}} \qquad for \ some \ 0 < \epsilon < 1$$

then for each set of vectors  $\kappa_0, \ldots, \kappa_N$  in  $\mathbb{R}^{N_+}$  there is a set  $\mathcal{Y}_{\kappa,\epsilon}$  with  $\mathbb{Q}\mathcal{Y}_{\kappa,\epsilon}^c < 2\epsilon$  on which

$$\sum\nolimits_{0 \leq j \leq N} |\kappa_j'(\widehat{g} - \gamma)|^2 \leq \frac{6 \left\|B_n^{-1}\right\|_2}{n\epsilon} \sum\nolimits_{0 \leq j \leq N} |D^{-1}\kappa_j|^2.$$

**Remark.** For our applications of the Corollary in Sections ?? and ??, we need  $D = \text{diag}(D_0, D_1, \dots, D_N)$  and  $\kappa_j = e_j$ , the unit vector with a 1 in its *j*th position, for  $j \leq m$  and  $\kappa_j = 0$  for j > m. In our companion paper we will need the more general  $\kappa_j$ 's.

**PROOF** First we establish a bound on the spectral distance between  $A_n^{-1}$ and  $B_n^{-1}$ . Define  $H = B_n^{-1}A_n - I$ . Then  $||H||_2 \le ||B_n^{-1}||_2 ||A_n - B_n||_2 \le 1/2$ , which justifies the expansion

$$\left\|A_{n}^{-1}-B_{n}^{-1}\right\|_{2} = \left\|\left((I+H)^{-1}-I\right)B_{n}^{-1}\right\|_{2} \le \sum_{j\ge 1} \|H\|_{2}^{k} \|B_{n}^{-1}\|_{2} \le \|B_{n}^{-1}\|_{2}.$$

As a consequence,  $\|A_n^{-1}\|_2 \leq 2 \|B_n^{-1}\|_2$ . Choose  $\epsilon_1 = 1/2$  and  $\epsilon_2 = \epsilon$  in Lemma 12. The bound on  $\max_{i \leq n} |\eta_i|$ gives the bound on  $\max_{i \leq n} |w_i|$  needed by the Lemma:

$$n|w_i|^2 = \eta'_i D(J_n/n)^{-1} D\eta_i = \eta'_i A_n^{-1} \eta_i \le \left\|A_n^{-1}\right\|_2 |\eta_i|^2.$$

Define  $K_j := J_n^{-1/2} \kappa_j$ , so that  $|\kappa'_j(\widehat{g} - \gamma)|^2 \le 2(K'_j W_n)^2 + 2(K'_j r_n)^2$ . By Cauchy-Schwarz,

$$\sum_{j} (K'_{j} r_{n})^{2} \leq \sum_{j} |K_{j}|^{2} |r_{n}|^{2} = U_{\kappa} |r_{n}|^{2}$$

where

$$U_{\kappa} := \sum_{j} \kappa'_{j} J_{n}^{-1} \kappa_{j} = \sum_{j} n^{-1} (D^{-1} \kappa_{j})' A_{n}^{-1} D^{-1} \kappa_{j}$$
$$\leq 2n^{-1} \left\| B_{n}^{-1} \right\|_{2} \sum_{j} |D^{-1} \kappa_{j}|^{2}.$$

For the contribution  $V_{\kappa} := \sum_{j} |K'_{j}W_{n}|^{2}$  the Cauchy-Schwarz bound is too crude. Instead, notice that  $\mathbb{Q}V_{\kappa} = U_{\kappa}$ , which ensures that the complement of the set

$$\mathcal{Y}_{\kappa,\epsilon} := \{ |W_n| \le \sqrt{N_+/\epsilon} \} \cap \{ V_\kappa \le U_\kappa/\epsilon \}$$

has  $\mathbb{Q}$  probability less that  $2\epsilon$ . On the set  $\mathcal{Y}_{\kappa,\epsilon}$ ,

$$\sum_{0 \le j \le N} |\kappa'_j(\widehat{g} - \gamma)|^2 \le 2V_\kappa + 2U_\kappa |r_n|^2 \le 3U_\kappa /\epsilon.$$

The asserted bound follows. 

lasso < 17 >

Example. Try to apply Theorem <6> to cover Theorem 2 of Knight and Fu (2000). Comment that more recent results with  $p = p_n$  are much closer to the way the lasso is used. 

# Notes

Convexity: Huber (1964) Brown (1985) and Appendix A of Maritz (1981). Andersen and Gill (1982) Haberman (1989) Niemiro (1992) Bickel, Klaassen, Ritov, and Wellner (1993, pp 328, 473, 519) Ritov (1987) Jurečková (1977)

# References

- Andersen, P. K. and R. D. Gill (1982). Cox's regression model for counting processes: a large sample study. Annals of Statistics 10, 1100–1120.
- Bickel, P. J., C. A. J. Klaassen, Y. Ritov, and J. A. Wellner (1993). Efficient and Adaptive Estimation for Semiparametric Models. Baltimore: Johns Hopkins University Press.
- Brown, B. M. (1985). Multiparameter linearization theorems. Journal of the Royal Statistical Society, Series B 47, 323–331.
- Dou, W., D. Pollard, and H. H. Zhou (2009). Functional regression for general exponential families. Technical report, Yale University Statistics Department.
- Haberman, S. J. (1989). Concavity and estimation. Annals of Statistics 17, 1631–1661.
- Hjort, N. L. and D. Pollard (1993). Asymptotics for minimisers of convex processes. Technical report, Yale University.
- Huber, P. J. (1964). Robust estimation of a location parameter. Annals of Mathematical Statistics 35, 73–101.
- Jurečková, J. (1977). Asymptotic relations of m-estimates and r-estimates in linear regression model. *Annals of Statistics* 5, 464–472.
- Knight, K. and W. Fu (2000). Asymptotics for lasso-type estimators. Annals of Statistics 28(5), 1356–1378.
- Maritz, J. S. (1981). *Distribution-Free Statistical Methods*. Chapman and Hall.
- Niemiro, W. (1992). Asymptotics for m-estimators defined by convex minimization. Annals of Statistics 20, 1514–1533.

- Pollard, D. (1991). Asymptotics for least absolute deviation regression estimators. *Econometric Theory* 7, 186–199.
- Portnoy, S. (1988). Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. Annals of Statistics 16(1), 356–366.
- Ritov, Y. (1987). Tightness of monotone random fields. Journal of the Royal Statistical Society, Series B 49, 331–333.
- Rockafellar, R. T. (1970). *Convex Analysis*. Princeton, New Jersey: Princeton Univ. Press.