

Iterative algorithms

- alternative to convex relaxation, with no shrinkage bias
- computationally efficient.

① Iterative hard thresholding

- Coordinate sparse linear regression:

$$Y = X\beta^* + \varepsilon \quad \text{with} \quad \|\beta^*\|_0 \text{ small and } \varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$$

- Reminder on Lasso and proximal optimisation:
- The lasso estimator

$$\hat{\beta}^{\text{Lasso}} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \} \quad (6.2)$$

has been derived as a convex relaxation of

$$\hat{\beta}^{\text{NS}} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda^2 \|\beta\|_0 \} \quad (6.1)$$

- proximal recipe:

→ approximate $F(\beta) = \|Y - X\beta\|^2$ by $F(\beta^t) + \langle \nabla F(\beta^t), \beta - \beta^t \rangle + \frac{1}{2\eta} \|\beta - \beta^t\|^2$

→ iterate

$$\beta^{t+1} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ F(\beta^t) + \langle \nabla F(\beta^t), \beta - \beta^t \rangle + \frac{1}{2\eta} \|\beta - \beta^t\|^2 + \lambda \|\beta\|_1 \right\}$$

$$= \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\beta - (\beta^t - 2\eta \nabla F(\beta^t))\|^2 + 2\lambda \|\beta\|_1 \right\}$$

$$= S_{\lambda\eta}(\beta^t - 2\eta \nabla F(\beta^t))$$

where

$$S_{\lambda}(\alpha) = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\alpha - \beta\|^2 + \lambda \|\beta\|_1 \right\}$$

$$= \begin{bmatrix} \alpha_1 (1 - \lambda/|\alpha_1|)_+ \\ \vdots \\ \alpha_p (1 - \lambda/|\alpha_p|)_+ \end{bmatrix} \in \mathbb{R}^p$$

} Soft thresholding operator

Iterative Hard Thresholding:

Why not applying the proximal method on the original problem (6.1) ?

$$\begin{aligned}
 \beta^{t+1} &\in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ F(\beta^t) + \langle \nabla F(\beta^t), \beta - \beta^t \rangle + \frac{1}{2\eta} \|\beta - \beta^t\|^2 + \lambda^2 \|\beta\|_0 \right\} \\
 &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \|\beta - (\beta^t - \eta \nabla F(\beta^t))\|^2 + 2\eta \lambda^2 \|\beta\|_0 \right\} \\
 &= H_{\lambda \sqrt{2\eta}} (\beta^t - \eta \nabla F(\beta^t))
 \end{aligned}$$

where

$$\begin{aligned}
 H_{\lambda}(\alpha) &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \|\alpha - \beta\|^2 + \lambda^2 \|\beta\|_0 \right\} \\
 \text{Exercise 2.8.1} &\rightarrow \left[\begin{array}{c} \alpha_1 \mathbb{1}_{|\alpha_1| > \lambda} \\ \vdots \\ \alpha_p \mathbb{1}_{|\alpha_p| > \lambda} \end{array} \right] \in \mathbb{R}^p \left. \vphantom{\operatorname{argmin}} \right\} \text{Hard Thresholding.}
 \end{aligned}$$

⌊ no shrinkage

⌊ (6.1) is non-convex so no convergence guarantee



- fix $\eta = 1/2, \hat{\beta}^0 = 0$
- choose λ_t decreasing from a high-value to the optimal level $\lambda_{\infty} = c \sigma \sqrt{\log p}$ for (6.1)

(IHT) $\left\{ \begin{array}{l} \cdot \eta = 1/2, \lambda_t = \alpha^{-t} A + B, \hat{\beta}^0 = 0 \\ \cdot \hat{\beta}^{t+1} = H_{\lambda_{t+1}} \left(\hat{\beta}^t - \frac{1}{2} \nabla F(\hat{\beta}^t) \right) = H_{\lambda_{t+1}} \left(\underbrace{(I - X^T X)}_{=: \Lambda} \hat{\beta}^t + X^T Y \right) \right.$

$= X^T (X \hat{\beta}^t - Y)$

② Risk analysis for IHT

a/ Discussion

- β_{Lasso}^t = approximation $\hat{\beta}^{\text{Lasso}}$ after t steps of Soft-Thresholding iterations.

• classical approach:

$$\|\beta_{\text{Lasso}}^t - \beta^*\| \leq \underbrace{\|\beta_{\text{Lasso}}^t - \hat{\beta}^{\text{Lasso}}\|}_{\text{optimisation error}} + \underbrace{\|\hat{\beta}^{\text{Lasso}} - \beta^*\|}_{\text{statistical error}}$$

→ the two errors are analysed apart

• IHT:

- $\hat{\beta}^t \rightarrow \dots \Rightarrow$ cannot split into optimisation / statistical errors

• recipe: $\Lambda = I - X^T X$

$$\hat{\beta}^{t+1} := H_{\lambda_{t+1}}(\Lambda \hat{\beta}^t + X^T Y) \stackrel{Y = X\beta^* + \varepsilon}{=} H_{\lambda_{t+1}}\left(\underbrace{\beta^*}_{\text{target}} + \underbrace{\Lambda(\hat{\beta}^t - \beta^*)}_{\text{contraction?}} + \underbrace{X^T \varepsilon}_{\text{statistical noise}}\right)$$



With a suitable tuning of λ_t (to keep $\hat{\beta}^t$ sparse)

$\hat{\beta}^t - \beta^*$ will ^{be} sparse and Λ will act as a contraction.

Hence for t large

$$\hat{\beta}^t \approx H_{\lambda_{\infty}}(\beta^* + Z)$$

optimal estimation of β^* from $\beta^* + Z$

b/ Risk bound

Theorem 6.1 (with $c=1$) Deterministic bound

Assumptions: $\Lambda = I - X^T X$

- for some $0 < \delta < 1/3$:

$$\max_{\substack{|S| \leq k \\ S \subset \{1, \dots, p\}}} |\Lambda_{SS}|_{op} \leq \delta, \text{ with } k = 3|\beta^*|_0 \quad (6.10)$$

$$1 < a \leq \frac{1}{3\delta}, \quad A \geq \frac{\|\beta^*\|}{3|\beta^*|_0^{1/2}}, \quad B > \frac{a}{a-1} |X^T \varepsilon|_{\infty} \quad (6.11)$$

Then: writing $m^* = \text{supp}(\beta^*)$

$$(i) \quad |\hat{\beta}_{\bar{m}^*}^t|_0 \leq |\beta^*|_0 \quad \text{where } \bar{m}^* = \{1, \dots, p\} \setminus m^*$$

$$(ii) \quad \|\hat{\beta}^t - \beta^*\| \leq 3\lambda_t \sqrt{|\beta^*|_0}$$

Discussion:

① since $|X^T \varepsilon|_{\infty} \leq \sigma \sqrt{2 \log p}$, choosing $B \leq c \sigma \sqrt{\log p}$, we have $\lambda_t \geq c \sigma \sqrt{\log p}$ and $\|\hat{\beta}^t - \beta^*\|^2 \leq c' \sigma^2 |\beta^*|_0 \log p$ as $t \rightarrow \infty$

\leadsto see Corollary 6.3 below

② the condition (6.10) ensures the contraction of $\Lambda(\hat{\beta}^t - \beta^*)$ for $\hat{\beta}^t$ sparse enough

③ (i) states that $\hat{\beta}^t$ remains sparse for $\lambda_t = a^{-t} A + B$.

④ the condition (6.10) ensures the restricted isometry property

$$(1-\delta) \|\beta\|^2 \leq \|X\beta\|^2 \leq (1+\delta) \|\beta\|^2 \text{ for all } \beta \text{ with } |\beta|_0 \leq k$$

(check it!)

Proof of Theorem 6.1:

notation: $Z = X^T \varepsilon$, $b^{t+1} = \beta^* + \lambda (\hat{\beta}^t - \beta^*) + Z$ so that $\hat{\beta}^{t+1} = H_{\lambda_{t+1}}(b^{t+1})$;
 $k^* = |\beta^*|_0$.

Lemma 6.2: contraction property.

If (i) and (ii) hold at step t , then $\forall |S| \leq k^*$

$$\|(b^{t+1} - \beta^*)_S\| \leq \delta \|\beta^* - \hat{\beta}^t\| + \sqrt{|S|} \|Z\|_\infty$$

Proof Lemma 6.2:

• Set $\mathcal{S}^t = m^* \cup \text{supp}(\hat{\beta}_{\bar{m}^*}^t)$ and $S' = \mathcal{S}^t \cup S$.

$|S'| \leq |\mathcal{S}^t| + |S| \leq 3k^* = \bar{k}$ by (i) at step t .

• $\|(b^{t+1} - \beta^*)_S\| \leq \|\lambda(\hat{\beta}^t - \beta^*)_S\| + \|Z_S\|$

$$\leq \|\lambda(\hat{\beta}^t - \beta^*)_{S'}\| + |S|^{1/2} \|Z\|_\infty$$

$$\leq \|\lambda_{S'}(\hat{\beta}^t - \beta^*)_{S'}\| + |S|^{1/2} \|Z\|_\infty$$

$$(6.10) \rightarrow \leq \delta \|\hat{\beta}^t - \beta^*_{S'}\| + |S|^{1/2} \|Z\|_\infty$$

$$= \delta \|\hat{\beta}^t - \beta^*\| + |S|^{1/2} \|Z\|_\infty.$$

□

• We analyse $\hat{\beta}^{t+1} - \beta^*$ on m^* and \bar{m}^* apart.

• On m^* :

$$\|\hat{\beta}_{m^*}^{t+1} - \beta_{m^*}^*\| \leq \underbrace{\|b_{m^*}^{t+1} - H_{\lambda_{t+1}}(b_{m^*}^{t+1})\|}_{\text{definition of } H_{\lambda_{t+1}}} + \underbrace{\|b_{m^*}^{t+1} - \beta_{m^*}^*\|}_{\Delta}$$

$$\leq \sqrt{|m^*|} \lambda_{t+1}$$

$$\leq \delta \|\hat{\beta}^t - \beta^*\| + \sqrt{|m^*|} \|Z\|_\infty$$

LEM. 6.2

$$(ii) \leq 3\lambda_t \sqrt{|m^*|}$$

$$\leq \sqrt{k^*} (\lambda_{t+1} + 3\delta \lambda_t + \|Z\|_\infty) \leq 2\lambda_{t+1} \sqrt{k^*} \quad (6.17)$$

$$(6.11) \rightarrow \leq \frac{1}{a} \lambda_t + \frac{a-1}{a} B = \lambda_{t+1}$$

on \bar{m}^* :

• We first prove $|\hat{\beta}_{\bar{m}^*}^{t+1}|_0 < k^*$.

For any $S \subset \text{supp}(\hat{\beta}_{\bar{m}^*}^{t+1})$ with $|S| \leq k^*$ we have

$$\lambda_{t+1} \sqrt{|S|} \stackrel{\text{def. of } H_{\lambda_{t+1}}}{\leq} \|\hat{\beta}_S^{t+1}\| = \|\hat{b}_S^{t+1}\| = \|\hat{b}_S^{t+1} - \beta_S^*\|$$

\uparrow $S \subset \text{supp}(\hat{\beta}_{\bar{m}^*}^{t+1})$ \uparrow $S \subset \bar{m}^*$

$$\stackrel{\text{Lem. 6.2}}{\leq} \sqrt{|S|} \|Z\|_\infty + \delta \|\beta^* - \hat{\beta}^t\| \quad (6.18)$$

$$(ii) \leq 3\lambda_t \sqrt{k^*}$$

$$(6.11) \leq \sqrt{k^*} \left(\frac{a-1}{a} B + \underbrace{3\delta\lambda_t}_{\leq \frac{1}{a}\lambda_t} \right)$$

and $|S| \leq k^*$

(6.18 bis)

$$\lambda_{t+1} = \frac{1}{a}\lambda_t + \frac{a-1}{a}B \stackrel{\leq}{\rightarrow} \lambda_{t+1}$$

Hence, $|S| < k^*$ which means that we cannot find

$S \subset \text{supp}(\hat{\beta}_{\bar{m}^*}^{t+1})$ with $|S| = k^*$.

We have proved $|\hat{\beta}_{\bar{m}^*}^{t+1}|_0 < k^*$, i.e. (i) for step $t+1$.

• upper-bound on $\|(\hat{\beta}_{\bar{m}^*}^{t+1} - \beta_{\bar{m}^*}^*)\|$: (6.18) with $S = \text{supp}(\hat{\beta}_{\bar{m}^*}^{t+1})$ gives

$$\|\hat{\beta}_{\bar{m}^*}^{t+1} - \beta_{\bar{m}^*}^*\| = \|\hat{\beta}_{\bar{m}^*}^{t+1}\| \stackrel{(6.18 \text{ bis})}{\leq} \lambda_{t+1} \sqrt{k^*}$$

conclusion: with (6.17)

$$\|\hat{\beta}^{t+1} - \beta^*\| \leq \|\hat{\beta}_{\bar{m}^*}^{t+1} - \beta_{\bar{m}^*}^*\| + \|\hat{\beta}_{\bar{m}^*}^{t+1} - \beta_{\bar{m}^*}^*\|$$

$$\leq 3\lambda_{t+1} \sqrt{k^*}$$

which is (ii) for step $t+1$.

□

Discussion: if we set $\hat{F} = \lceil \log_a(A/B) \rceil$, then

$$\left\{ \begin{array}{l} \lambda_{\hat{E}} = a^{-\hat{F}} A + B \leq 2B. \\ \text{so } \|\hat{\beta}^{\hat{E}} - \beta^*\| \leq 6B \sqrt{|\beta^*|_0}. \end{array} \right.$$

Corollary 6.3: error bound for LHT

- Assume that $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$, $\|X_j\| = 1$, for $j=1, \dots, p$ and (6.10) hold.
- Set for $K > 1$ and $1 < a \leq 1/(3\delta)$

$$A = \frac{\|X^T Y\| + \sigma |X|_{\log p} (\sqrt{2} + \sqrt{2K \log p})}{3(1-\delta)} \quad \text{and} \quad B = \frac{a\sigma}{a-1} \sqrt{2K \log p}$$

- Then, with probability $\geq 1 - \frac{2}{p^{K-1}}$

$$(i) \quad \|\hat{\beta}^{\hat{E}} - \beta^*\|^2 \leq 2K \left(\frac{6a}{a-1}\right)^2 \sigma^2 |\beta^*|_0 \log(p)$$

$$(ii) \quad \|X \hat{\beta}^{\hat{E}} - X \beta^*\|^2 \leq 2(1+\delta)K \left(\frac{6a}{a-1}\right)^2 \sigma^2 |\beta^*|_0 \log(p)$$

Proof of Corollary 6.3:

- All we need is to prove that (6.11) holds. Indeed, then (i) holds and $|\hat{\beta}_{\hat{m}^{\hat{E}}}^{\hat{E}}|_0 \leq |\beta^*|_0$ so the restricted isometry property (6.10) ensures (ii). (from (i))

- In the analysis of Lasso estimator, we have already seen that $\mathbb{P}[B > \frac{a}{a-1} |X^T \varepsilon|_{\infty}] \geq 1 - \frac{1}{p^{K-1}}$

- So, it remains to prove that

$$\mathbb{P}\left[A \geq \frac{\|\beta^*\|}{3\sqrt{|\beta^*|_0}}\right] \geq 1 - \frac{1}{p^{K-1}} \quad (*)$$

- Gaussian concentration: $\varepsilon \rightarrow \|X^T X \beta^* + X^T \varepsilon\|$ is $|X|_{op}$ -Lipschitz so there exist $\zeta, \zeta' \sim \text{Exp}(1)$ such that

$$\mathbb{E}[\|X^T Y\|] - \sigma |X|_{op} \sqrt{2\zeta} \leq \|X^T Y\| \leq \mathbb{E}[\|X^T Y\|] + \sigma |X|_{op} \sqrt{2\zeta'}$$

As in Lecture 1:

$$\begin{aligned} \mathbb{E}[\|X^T Y\|^2] &\leq \mathbb{E}\left[\left(\mathbb{E}[\|X^T Y\|] + \sigma |X|_{op} \sqrt{2\zeta'}\right)^2\right] \\ &\stackrel{\text{Jensen}}{\leq} \left(\mathbb{E}[\|X^T Y\|] + \sigma |X|_{op} \sqrt{2}\right)^2 \end{aligned}$$

So

$$\begin{aligned} \|X^T Y\| + \sigma |X|_{op} (\sqrt{2} + \sqrt{2\zeta'}) &\geq \mathbb{E}[\|X^T Y\|] + \sigma |X|_{op} \sqrt{2} \\ &\geq \mathbb{E}[\|X^T Y\|^2]^{1/2} \\ &= \sqrt{\|X^T X \beta^*\|^2 + \underbrace{\mathbb{E}[\|X^T \varepsilon\|^2]}_{\geq 0}} \geq \|(X^T X \beta^*)_{m^*}\| \\ &\geq \|\beta_{m^*}^*\| - \|\Lambda_{m^*} \beta_{m^*}^*\| = \|\beta_{m^*}^*\| - \|\Lambda_{m^* m^*} \beta_{m^*}^*\| \\ (6.10) &\geq (1-\delta) \|\beta_{m^*}^*\| = (1-\delta) \|\beta^*\| \end{aligned}$$

Hence (*) holds.

□

Take home message:

- we have obtained some results similar as for the Lasso
- the number \hat{E} of iterations is controlled.

③ Group sparsity: Iterative Group Thresholding

→ Section 6.2 of lecture notes.