

Resilience Beyond Mean Estimation

Logistics Pset 2 due today
Pset 3 should be out by tomorrow afternoon

Recaps

- Defined general $G \uparrow + G \downarrow \rightarrow$ bounds on modulus of continuity
- Second moment estimation: higher moments \Rightarrow resilience (nuclear norm reduction)
- Linear regression:
 - Hypercontractivity + bounded noise

Today

- Finish linear regression
 - State result for classification
 - Pay back debts: prove missing dual norm results
-

Linear regression: $Y = \langle \theta^*, X \rangle + \epsilon$, $\mathbb{E}[X\epsilon] = 0$

2 conditions

$\Rightarrow \theta^*$ is least squares estimator for (X, Y)

$$(1) \mathbb{E}_p[\langle X, v \rangle^4] \leq K \mathbb{E}_p[\langle X, v \rangle^2]^2 \quad \forall v \in \mathbb{R}^d$$

(hypercontractivity)

$$(2) \mathbb{E}_p[X \epsilon X^T] \leq \sigma^2 \cdot \mathbb{E}_p[XX^T]$$

(bounded noise)

Proposition If (1) and (2) hold, then p is

$(p, 5p, \epsilon)$ -resistant for $p = 2\sigma^2 \epsilon$ as long as

$$\epsilon \leq \frac{1}{8} \text{ and } \epsilon(k-1) \leq \frac{1}{6}$$

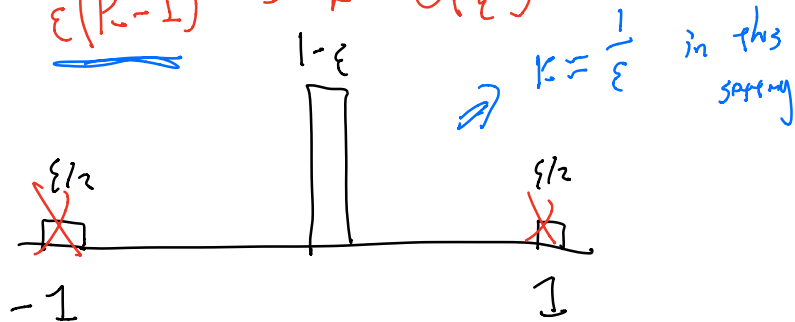
Interpretation, $L(p, \theta) = \mathbb{E}_p[\underbrace{(Y - \langle \theta, X \rangle)^2}_{\text{sq-and}}] - \mathbb{E}_p[(Y - \langle \theta^*, X \rangle)^2]$

Resistance \Rightarrow can robustly estimate θ w/ loss $\leq 5p$

$= 10\sigma^2 \epsilon$ \leftarrow roughly same sort of result as for mean estimation

\Rightarrow Definitely need $\epsilon \leq \frac{1}{2}$.

$k \geq 1$ $\epsilon(k-1) \rightarrow \frac{\text{Need}}{k} = \mathcal{O}\left(\frac{1}{\epsilon}\right)$



Proof Two conditions, G_d and G_T

$G_d: L(r, \theta^*) \leq p$, whenever $r \geq \frac{p}{1-\epsilon}$

$G_T: \text{If } L(r, \theta) \leq p, \text{ then } L(p, \theta) \leq 5p$

$L(r, \theta) = (\theta - \theta^*(r))^T S_r (\theta - \theta^*(r))$

$S_r = \mathbb{E}_r[XX^T]$

$$\mathcal{S}_\downarrow: (\theta^*(p) - \theta^*(r))^T \boxed{S_r} (\theta^*(p) - \theta^*(r)) \leq p \quad \hookrightarrow \|S_r^{1/2}(\theta^*(r) - \theta^*(p))\|_2$$

$$\mathcal{S}_\uparrow: \text{If } (\theta - \theta^*(r))^T \boxed{S_r} (\theta - \theta^*(r)) \leq p, \\ \text{then } (\theta - \theta^*(p))^T \boxed{S_p} (\theta - \theta^*(r)) \leq 5p.$$

Proof strategy

$$\textcircled{1} S_r \approx S_p \rightarrow \boxed{\frac{1}{2} S_p \leq S_r \leq \frac{3}{2} S_p.}$$

hypercontractivity

$$\textcircled{2} \|S_p^{1/2}(\theta^*(r) - \theta^*(p))\|_2 \text{ small by resilience}$$

bounded noise

$$\textcircled{1} \begin{cases} S_p = \mathbb{E}[xx^T] \\ S_r = \mathbb{E}_r[xx^T] \end{cases}$$

$$\text{For all } v, \quad \frac{1}{2} v^T S_p v \leq v^T S_r v$$

$$\rightarrow \frac{1}{2} \mathbb{E}_p[\langle x, v \rangle^2] \leq \mathbb{E}_r[\langle x, v \rangle^2]$$

S_r suffices to show:

$$\left| \underbrace{\mathbb{E}_r[\langle x, v \rangle^2]}_{\text{mean of } \langle x, v \rangle^2 \text{ under } r} - \underbrace{\mathbb{E}_p[\langle x, v \rangle^2]}_{\text{mean under } p} \right| \leq \frac{1}{2} \mathbb{E}_p[\langle x, v \rangle^2]$$

Recall Chebyshev: If $\text{Var}[Z] \leq \sigma^2$, then Z is $(\epsilon, \frac{\sigma\sqrt{\epsilon}}{1-\epsilon})$ -resilient.

$$Z = \langle X, v \rangle^2 \downarrow$$

$$\text{Var}[Z] = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2 = \mathbb{E}[\langle X, v \rangle^4] - \mathbb{E}[\langle X, v \rangle^2]^2$$

$$\Rightarrow \text{Can take } \sigma = \sqrt{k} \cdot \mathbb{E}[\langle X, v \rangle^2] \leq (k-1) \mathbb{E}[\langle X, v \rangle^2]^2$$

$$|\mathbb{E}_r[\langle X, v \rangle^2] - \mathbb{E}_p[\langle X, v \rangle^2]| \leq 2 \cdot \frac{\sigma \sqrt{\epsilon}}{1-\epsilon} = \frac{2\sqrt{\epsilon(k-1)}}{1-\epsilon} \mathbb{E}[\langle X, v \rangle^2]$$

Need this to be $\leq \frac{1}{2}$.

Holds if $\epsilon \leq \frac{1}{8}$

and $\epsilon(k-1) \leq \frac{1}{6}$.

$$\Rightarrow |\mathbb{E}_r[\langle X, v \rangle^2] - \mathbb{E}_p[\langle X, v \rangle^2]| \leq \frac{1}{2} \mathbb{E}[\langle X, v \rangle^2]$$

$$\Rightarrow S_r \leq \frac{1}{2} S_p$$

$$\Rightarrow S_r \leq \frac{3}{2} S_p \quad \blacksquare$$

Claim.

Want to somehow show $\Theta^*(p) - \Theta^*(r)$ is small

$$\text{Claim. } \Theta^*(r) - \Theta^*(p) = S_r^{-1} \mathbb{E}_r[XZ]$$

$$S_r^{-1} \mathbb{E}_r[XZ]$$

$$= \mathbb{E}_r[XX^T]^{-1} \mathbb{E}_r[XZ]$$

$$\hookrightarrow Z = Y - \langle \Theta^*(p), X \rangle$$

$$= \mathbb{E}_r [X X^T]^{-1} \mathbb{E}_r [X Y - X X^T \theta^*(p)]$$

$$= \mathbb{E}_r [X X^T]^{-1} \mathbb{E}_r [X Y] - \theta^*(p)$$

cancel

$$= \theta^*(w) - \theta^*(p)$$

least squares formula

$$\hat{\theta}_{OLS} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right)$$

$\mathbb{E}_p [X Z] = 0$ by definition

$$\theta^*(p) - \theta^*(w) = S_r^{-1} \left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)$$

resilience (bounded noise)

$$(\theta^*(p) - \theta^*(w)) S_r (\theta^*(p) - \theta^*(w)) \leq p$$

$\hookrightarrow \|S_r^{1/2} (\theta^*(w) - \theta^*(p))\|_2$

$$\left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)^T S_r^{-1} \cdot S_r \cdot S_r^{-1} \left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)$$

$$= \left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)^T S_r^{-1} \left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)$$

$$\leq 2 \cdot \left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)^T S_p^{-1} \left(\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z] \right)$$

$$= 2 \|S_p^{-1/2} (\mathbb{E}_p [X Z] - \mathbb{E}_r [X Z])\|_2 = \mathbb{E}_p [w] - \mathbb{E}_r [w]$$

need to bound this

Regress to residual for variable $W = S_p^{-1} X Z$

$$\text{Var}_p[w] = \mathbb{E}_p[w w^T]$$

$$= \mathbb{E}_p[S_p^{-1/2} X Z^2 X^T S_p^{-1/2}]$$

$= O(\sigma^2 \cdot \epsilon)$

$\Rightarrow G \downarrow$

$$= S_p^{-1/2} \underbrace{\mathbb{E}_p[X Z^2 X^T]}_{\leq \sigma^2 \cdot \mathbb{E}_p[XX^T]} S_p^{-1/2}$$

$$= \sigma^2 \cdot S_p$$

$$\leq \sigma^2 \cdot S_p^{-1/2} S_p S_p^{1/2} = \sigma^2 \cdot I$$

\Rightarrow bounded variance \rightarrow residual

~~$G \downarrow$~~ : $(\theta^*(p) - \theta^*(w))^T \boxed{S_r} (\theta^*(p) - \theta^*(w)) \leq p$ $\hookrightarrow \|S_r^{1/2} (\theta^*(w) - \theta^*(p))\|_2$

$G \uparrow$: If $(\theta - \theta^*(w))^T \boxed{S_r} (\theta - \theta^*(w)) \leq p$,
then $(\theta - \theta^*(p))^T \boxed{S_p} (\theta - \theta^*(p)) \leq 5p$.

$$\sqrt{(\theta - \theta^*(p))^T S_p (\theta - \theta^*(p))}$$

$$\leq \sqrt{(\theta - \theta^*(w))^T \underbrace{S_p}_{S_r + G \uparrow} (\theta - \theta^*(w))} + \sqrt{(\theta^*(w) - \theta^*(p))^T \underbrace{S_p}_{S_r + G \downarrow} (\theta^*(w) - \theta^*(p))}$$

$\leq O(\sigma^2 \cdot \epsilon)$.



$$\begin{aligned} \mathbb{E}_p[X\varepsilon] &= \mathbb{E}_p[X(Y - \langle \theta^*, X \rangle)] \\ &= \mathbb{E}_p[XY] - \mathbb{E}_p[XX^T]\theta^* = 0 \end{aligned}$$

~~$\mathbb{E}_p[\varepsilon|X] = 0 \quad \forall X$~~ \leftarrow assuming linear model, non-specified

(X, Y)

Added features?

$(X + H, Y)$

$$H \sim N(0, \lambda \cdot I)$$



ridge regression w/ param λ

X hypercontractive $\Rightarrow X + H$

Linear regression

Recap: hypercent. $\Rightarrow \frac{1}{2} S_p \leq S_r \leq \frac{3}{2} S_p$

resistance of $S_p^{1/2} XZ \Rightarrow$ bounded noise
 \Rightarrow combined w/ triangle ineq.

Linear classification

$$(X, Y) \in \mathbb{R}^d \times \{-1, +1\}$$

$$L(p, \theta) = P_p [Y \neq \text{sign}(\langle X, \theta \rangle)]$$

linear classifier

0/1-loss

$$B(p, \theta) = \mathbb{E}_{X, Y} [\max(1 - y \langle \theta, X \rangle, 0)]$$

hinge loss

hinge $\geq 0/1$

$\hookrightarrow S_{\downarrow} / S_{\uparrow}$
"bordering"

Prop. Suppose p satisfies:

$$(1) \mathbb{E}_{(X, Y) \sim p} [\max(1 - y \langle \theta^*, X \rangle, 0)] \leq (1 - \epsilon) p_1$$

$$(2) \text{ If } P_{(X, Y) \sim p} [y \langle \theta, X \rangle \leq \frac{1}{2}] \leq \epsilon + 2(1 - \epsilon) p_1$$

then $P_{(X, Y) \sim p} [y \langle \theta, X \rangle \leq 0] \leq p_2$.

Then p is (p_1, p_2, ϵ) -resistant for L and B .

Dual norms.

Schatten p -norms

X , singular values $\sigma_1, \dots, \sigma_n$

$$\|X\|_{S(p)} = \sqrt[p]{\sigma_1^p + \dots + \sigma_n^p}$$

$S(1)$: nuclear norm

$S(\infty)$: operator norm

Thm. $S(p)$ and $S(q)$ are dual norms if $\frac{1}{p} + \frac{1}{q} = 1$.

Key lemma.

For any A, B w/ singular values σ_i (in sorted order),

$$\langle A, B \rangle \leq \sum_{i=1}^n \sigma_i \tau_i$$

① Reduce to PSD

② Prove f -PSD

① Lemma,

$$\langle A, B \rangle^2 \leq \langle (A^T A)^{1/2}, (B^T B)^{1/2} \rangle \langle (A A^T)^{1/2}, (B B^T)^{1/2} \rangle$$

$$\begin{bmatrix} \lambda (A^T A)^{1/2} & A \\ A^T & \lambda (A A^T)^{1/2} \end{bmatrix} \succeq 0$$

\Rightarrow take matrix
Jost product

$$\begin{bmatrix} \lambda (B^T B)^{1/2} & -B \\ -B^T & \lambda (B B^T)^{1/2} \end{bmatrix} \succeq 0$$

(analogous to Cauchy-Schwarz)

② Lemma

$$\sum_{i=1}^k B_{ii} \leq \sum_{i=1}^k \tau_i$$

\Leftarrow Cauchy interlacing theorem

$$\langle A, B \rangle = \langle \text{diag}(\sigma_1, \dots, \sigma_n), B \rangle$$

$$\sum_{i=1}^n \sigma_i B_{ii} \leq \sum_{i=1}^n \sigma_i \tau_i$$

"stochastic dominance"
(Abel summation)