

10/1/23

Lecture 9: Robust Stats, Iterative Filtering

Setup: q has mean $\mu \in \mathbb{R}^d$ and covariance $\Sigma \leq \text{Id}$,

Nature samples $x_1^*, \dots, x_n^* \sim q$, adversary

Corrupts arbitrary γ fraction, we are given
the corrupted samples $\{x_1, \dots, x_n\}$

Decompose $\{x_1, \dots, x_n\}$ into

$$S_{\text{good}} \cup S_{\text{bad}} \setminus S_r$$

where $|S_{\text{bad}}| = |S_r| = \gamma n$,

$S_{\text{good}} = \{x_1^*, \dots, x_n^*\}$ are the original i.i.d draws

from q , S_r are the points which were corrupted,
and S_{bad} are the points they've been replaced with.

Assume for now that for

$$\mu_g \stackrel{\triangle}{=} \frac{1}{|S_{\text{good}}|} \sum_{x \in S_{\text{good}}} x$$

$$\Sigma_g \stackrel{\triangle}{=} \frac{1}{|S_{\text{good}}|} \sum_{x \in S_{\text{good}}} (x - \mu_g)(x - \mu_g)^T,$$

$$\textcircled{1} \quad \|\mu_g - \mu\|_2 \leq \sqrt{\epsilon} \quad \text{and} \quad \textcircled{2} \quad \|\Sigma_g\|_{\text{op}} \leq 1$$

we'll use
Shorthand:

$$\sum_{\text{clean}} x_i : \sum_{x \in S_{\text{good}} \setminus S_r} x$$

$$\sum_{\text{bad}} x_i : \sum_{x \in S_{\text{bad}}} x$$

(we will see in pset 3 why this is a valid assumption. note: a little subtle b/c we don't assume the higher moments of q are bounded)

We will maintain weights $w = \{w_i\}_{i \in [n]}$ for the dataset that indicate how confident we are that x_i is clean.

$$0 \leq w_i \leq \frac{1}{n} \quad \forall i \in [n] \quad (\dagger)$$

Ideally we would want $w_i = \frac{1}{n} \cdot \mathbb{1}[i \in S_{\text{good}}]$.

NOTE: w_i is like " $\frac{1}{n}a_i$ " from SoS analysis.

But these are actual #'s now, not SoS variables

Define the weighted mean and weighted covariance by

$$\mu_w \stackrel{\Delta}{=} \frac{1}{\sum_i w_i} \sum_i w_i x_i$$

$$\Sigma_w \stackrel{\Delta}{=} \frac{1}{\sum_i w_i} \sum_i w_i (x_i - \mu_w)(x_i - \mu_w)^T$$

In general, $\{w_i\}$ can be viewed as "soft" indicators for a subset. We want to pick out a large subset of clean points, so we care about w_i 's that satisfy

$$\sum_i w_i \geq 1 - \gamma \quad (\text{**})$$

Note: set of w satisfying (b) and (bb) is convex hull K_γ of $\left\{ \frac{1}{n} \cdot \mathbf{1}_S : S \subseteq [n] \text{ s.t. } |S| \geq 1 - \gamma \right\}$

We will maintain that our $w \in K_\gamma$ throughout the algorithm.

Main lemma: (if μ_w wrong, $\|\sum_w\|$ large):

For any $w \in K_\gamma$,

$$\|\mu_g - \mu_w\|_2 \leq \sqrt{\gamma} \left(1 + \sqrt{\|\sum_w\|_{\text{op}}} \right)$$

First, let's see what to do with this....

Given current weights w , define scores

$$\tau_i \triangleq \langle v, x_i - \mu_w \rangle^2,$$

where v is top eigenvector of \sum_w . Let $\tau_{\max} = \max_{i: w_i > 0} \tau_i$

higher score $\tau_i \iff$ More likely x_i is corrupted

Lemma: Consider update rule

$$w'_i \leftarrow \left(1 - \frac{\tau_i}{\tau_{\max}}\right) w_i$$

Then if $\sum_{\text{clean } i} w_i \tau_i < \sum_{\text{bad } i} w_i \tau_i$, then

$$\sum_{\text{clean } i} (w_i - w'_i) < \sum_{\text{bad } i} (w_i - w'_i)$$

"Safety condition"

(More bad mass removed than good mass).

"progress condition" and $\text{nnz}(w') < \text{nnz}(w)$.

↑
nonzero entries

$$\text{PF: } \sum_{\text{clear:}} (w_i - w'_i) = \sum_{\text{clear:}} \frac{\tau_i w_i}{\tau_{\max}}$$

$$\text{and } \sum_{\text{bad:}} (w_i - w'_i) = \sum_{\text{bad:}} \frac{\tau_i w_i}{\tau_{\max}}$$

"Progress Condition" immediate from def. \square

We will maintain invariant that

$$\sum_{\text{clear:}} \left(\frac{1}{n} - w_i \right) < \sum_{\text{bad:}} \left(\frac{1}{n} - w_i \right), \quad (\text{INV})$$

i.e. less good mass removed than bad mass.

As long as $\|\sum w\|$ still large, we update the weights via the rule in the Lemma.

Obs 1: If (INV) always maintained, then algorithm runs for at most $2^{\gamma n}$ iterations.

PF: After $2^{\gamma n}$ iterations, we must have removed at least 2^{γ} mass total, thus $\geq \gamma$ mass from bad points (by INV), resulting in

$\sum_{\text{bad:}} w_i = 0$. Because INV always maintained, and there is no bad mass left, alg. must terminate. \square

Obs 2: If $\|\sum_w\| \leq 1$, then safe to output μ_w by Main Lemma.

Obs 3: If (INV) always maintained, we $\propto_{2\gamma}$.

$$\text{Pf: } \sum_{\text{clean:}} (\frac{1}{n} - w_i) < \sum_{\text{bad:}} (\frac{1}{n} - w_i) \leq \gamma, \text{ so } \sum_{\text{all:}} w_i > 1 - 2\gamma. \square$$

By Obs's 1, 2, suffices to show:

Lemma: Suppose $\|\sum_w\| >> 1$, and (INV)

holds, then $\sum_{\text{clean:}} w_i \tau_i < \sum_{\text{bad:}} w_i \tau_i$.

Pf: Note

$$\begin{aligned} \sum_{\text{all:}} w_i \tau_i &= \sum_i w_i \langle v, x_i - \mu_w \rangle^2 \\ &= v^T \left[\sum_i w_i (x_i - \mu_w)(x_i - \mu_w)^T \right] v \\ &= v^T \sum_w v = \|\sum_w\|_{\text{op}} \end{aligned}$$

so suffices to show $\sum_{\text{clean:}} w_i \tau_i \leq \frac{1}{2} \|\sum_w\|_{\text{op}}$.

Note $\sum w_i \tau_i \leq \frac{1}{n} \sum_{x \in S_{\text{good}}} \langle v, x - \mu_w \rangle^2$
 clean:

$$= \frac{1}{n} \sum_{x \in S_{\text{good}}} \langle v, (x - \mu_g) + (\mu_g - \mu_w) \rangle^2$$

$$\leq \underbrace{\frac{2}{n} \sum_{x \in S_{\text{good}}} \langle v, x - \mu_g \rangle^2}_{(A)} + \underbrace{\frac{2}{n} \sum_{x \in S_{\text{good}}} \langle v, \mu_g - \mu_w \rangle^2}_{(B)}$$

$$(A) \leq 2v^\top \left(\frac{1}{|S_{\text{good}}|} \sum_{x \in S_{\text{good}}} (x - \mu_g)(x - \mu_g)^\top \right)v \leq 1 \text{ by}$$

initial assumption on P. 1.

$$(B) \leq 2 \langle v, \mu_g - \mu_w \rangle^2$$

$$\leq 2 \|\mu_g - \mu_w\|^2$$

(Main Lemma)

$$\leq \sqrt{\gamma} \left(1 + \sqrt{\|\sum_w\|_{\text{op}}} \right)$$

Provided γ sufficiently small constant,
if $\|\Sigma_w\| \geq C$ for sufficiently large constant,

$$\textcircled{A} + \textcircled{B} \leq 2 + O\left(\sqrt{\gamma}\left(1 + \sqrt{\|\Sigma_w\|_{op}}\right)\right)$$

$$\leq \frac{\|\Sigma_w\|_{op}}{2}.$$

□

Algorithm :

- $w_i \leftarrow \frac{1}{n} \quad \forall i \in [n]$
- While $\|\Sigma_w\|_{op} \geq C$:
 - $v \leftarrow$ top eigenvector of Σ_w
 - $\tau_i \leftarrow \langle v, x_i - \mu_w \rangle^2 \quad \forall i \in [n]$
 - $\tau_{max} \leftarrow \max_{i: w_i > 0} \tau_i$
 - $w'_i \leftarrow w_i \left(1 - \frac{\tau_i}{\tau_{max}}\right)$
- Output μ_w

Remains to prove Main Lemma:

Pf: (will look like an SOS proof)

$$\left(\sum_{\text{all } i} w_i \right) \| \mu_w - \mu_g \|^2$$

(for $x \in S_{\text{good}} \setminus [n]$,
define $w_i = 0$)

$$= \sum_{\text{all } i} w_i \langle \mu_w - \mu_g, \mu_w - \mu_g \rangle$$

$$= \sum_{\text{all } i} w_i \langle x_i - \mu_g, \mu_w - \mu_g \rangle$$

$$= \sum_{\text{bad } i} " " + \sum_{i \in S_{\text{good}}} "$$

$$= \sum_{\text{bad } i} w_i \| \mu_w - \mu_g \|^2 +$$

$$\sum_{i \in S_{\text{good}}} \frac{1}{n} \langle x_i - \mu_g, \mu_w - \mu_g \rangle$$

(bbb)

$$+ \underbrace{\sum_{\text{bad } i} w_i \langle x_i - \mu_w, \mu_w - \mu_g \rangle}_{\text{I}} + \underbrace{\sum_{i \in S_{\text{good}}} (w_i - \frac{1}{n}) \langle x_i - \mu_g, \mu_w - \mu_g \rangle}_{\text{II}}$$

$$\text{I}^2 = \left(\sum_{\text{bad } i} w_i \langle x_i - \mu_w, \mu_w - \mu_g \rangle \right)^2$$

$$\leq \left(\sum_{\text{bad } i} w_i \right) \cdot \left(\sum_{\text{bad } i} w_i \langle x_i - \mu_w, \mu_w - \mu_g \rangle^2 \right)$$

$$\leq \gamma \cdot \underbrace{\sum_{\text{all } i} "}_{\text{''}}$$

$$= \gamma \cdot (\mu_w - \mu_g)^T \underbrace{\left[\sum_{\text{all } i} w_i (x_i - \mu_w)(x_i - \mu_w)^T \right]}_{\sum_w} (\mu_w - \mu_g)$$

$$\leq \gamma \|\sum_w\|_{\text{op}} \cdot \|\mu_w - \mu_g\|^2$$

$$\text{so } \textcircled{I} \leq \sqrt{\gamma} \cdot \sqrt{\|\sum_w\|_{\text{op}}} \cdot \|\mu_w - \mu_g\|$$

$$\textcircled{II} = \left(\sum_{i \in S_{\text{good}}} n \left(w_i - \frac{1}{n} \right)^2 \right) \cdot \left(\frac{1}{n} \sum_{i \in S_{\text{good}}} \langle x_i - \mu_g, \mu_w - \mu_g \rangle^2 \right)$$

$|n(w_i - \frac{1}{n})| \leq 1,$

and
 $w_i = 0 \text{ for } i \in S_{\text{good}} \setminus \{n\}$

clean: $\leq \sum |w_i - \frac{1}{n}| + \sum_{i \in S_{\text{good}} \setminus \{n\}} \frac{1}{n} (\mu_w - \mu_g)^T \sum_g (\mu_w - \mu_g) \leq \|\mu_w - \mu_g\|^2$

$\leq 2\gamma \text{ (b/c } w \in K_\gamma \text{ and } |S_{\text{good}} \setminus \{n\}| \leq \gamma n)$

$$\text{so } \textcircled{II} \leq \sqrt{2\gamma} \|\mu_w - \mu_g\|$$

Substituting into (500), rearranging,
and dividing by $\|\mu_w - \mu_g\|$, we get

$$\left(\sum_{\text{clean } i} w_i \right) \|\mu_w - \mu_g\| \leq \sqrt{\gamma} \left(1 + \sqrt{\|\sum w_i\|_{\text{op}}} \right). \quad \square$$