

Problem Set 1 – Solutions

Samuel B. Hopkins

November 30, 2024

Problem 1. In this problem, we will solve a sparse version of robust mean estimation. Let $\mu \in \mathbb{R}^d$ be an unknown k -sparse vector, in that only k of its entries are non-zero. First $n = \tilde{\Omega}(k^2(\log d)/\varepsilon^2)$ samples $v_1, \dots, v_n \in \mathbb{R}^d$ are drawn from $\mathcal{N}(\mu, \text{Id})$. Then an adversary alters εn of the samples and reorders them arbitrarily. We observe the resulting dataset v'_1, \dots, v'_n . Our goal will be to give an algorithm for estimating μ from these samples.

- (a) Let $\bar{v} = \frac{1}{n} \sum_{i=1}^n v_i$. Prove that with 0.99 probability, for all k -sparse vectors $u \in \mathbb{R}^d$ with $\|u\| = 1$,

$$\langle u, \bar{v} - \mu \rangle^2 \leq \varepsilon^2.$$

- (b) Define $\Sigma = \frac{1}{n} \sum_{i=1}^n (v_i - \bar{v})(v_i - \bar{v})^T$. Prove that with 0.99 probability, $|\Sigma_{ij}| \leq 1/k$ for $i \neq j$ and $|\Sigma_{ii} - 1| \leq 1/k$ for all $i, j \in [d]$.

- (c) Consider the following system, which we call \mathcal{S} , with scalar variables w_1, \dots, w_n and d -dimensional variables z_1, \dots, z_n

$$\begin{aligned} w_i^2 &= w_i \\ \sum_{i=1}^n w_i &\geq (1 - \varepsilon)n \\ w_i(z_i - v'_i) &= 0 \\ \bar{z} &= \frac{1}{n} \sum_{i=1}^n z_i, \quad \Sigma = \frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(z_i - \bar{z})^T \\ -\frac{1}{k} &\leq \Sigma_{ij} \leq \frac{1}{k} \quad \text{for all } i \neq j \\ -\frac{1}{k} &\leq \Sigma_{ii} - 1 \leq \frac{1}{k} \quad \text{for all } i \end{aligned}$$

Prove that with 0.99 probability, there is a feasible solution to this system where the w_i are indicators of the clean samples and the z_i are the actual clean samples.

From now on, assume that the events in (a), (b), (c) hold.

- (d) Now we consider the SoS relaxation of the system \mathcal{S} . Let $u \in \mathbb{R}^d$ be an arbitrary k -sparse vector with $\|u\| = 1$. Prove that

$$\mathcal{S} \vdash_2 \sum_{i=1}^n \langle u, z_i - v_i \rangle^2 \leq 10n(1 + \langle u, \bar{z} - \mu \rangle^2)$$

where recall v_i are the clean samples drawn from $N(\mu, I)$.

(e) Let $u \in \mathbb{R}^d$ be an arbitrary k -sparse vector with $\|u\| = 1$. Use part (c) to prove that

$$\mathcal{S} \vdash_4 \langle u, \bar{z} - \bar{v} \rangle^2 \leq 100\varepsilon(1 + \langle u, \bar{z} - \mu \rangle^2)$$

(f) Use part (e) to deduce that

$$\mathcal{S} \vdash_4 \langle u, \bar{z} - \mu \rangle^2 \leq O(\varepsilon).$$

Put everything together to show that there is a polynomial time algorithm that takes the samples v'_1, \dots, v'_n and with probability 0.9, outputs a k -sparse $\hat{\mu}$ such that $\|\mu - \hat{\mu}\| \leq O(\sqrt{\varepsilon})$.

Solution

(a) Note that $\bar{v} - \mu$ is distributed as $\mathcal{N}(0, (1/n)\text{Id})$. For ease of notation, let Z be distributed as $\mathcal{N}(0, (1/n)\text{Id})$, so the goal is to show that with probability 0.99, for all unit k -sparse vectors u ,

$$\langle u, Z \rangle^2 \leq \varepsilon.$$

By Cauchy-Schwarz,

$$\langle u, Z \rangle^2 \leq \sum_{i: u_i \neq 0} Z_i^2.$$

By the definition of sparsity, $\{i : u_i \neq 0\}$ is some set of size at most k . Fix some subset $S \subseteq [d]$ of size k . Then,

$$\Pr \left[\sum_{i \in S} Z_i^2 \geq \varepsilon^2 \right] = \Pr_{Y \sim \mathcal{N}(0, \text{Id}_k)} [\|Y\|^2 \geq n\varepsilon^2].$$

Standard concentration bounds (e.g. Bernstein's inequality used with the subexponentiality of χ^2 random variables) imply that this quantity is bounded by $\exp\left(-\Omega\left(\frac{\varepsilon^2 n}{k}\right)\right) = \exp(-\Omega(k \log d))$. We can now take a union bound over all subsets S of size at most k – there are $\exp(O(k \log(d)))$ such subsets, completing the proof if we take the constant factor in n sufficiently large.

(b) We have

$$\begin{aligned} \Sigma_{ij} &= \frac{1}{n} \sum_{r=1}^n (v_r - \bar{v})_i (v_r - \bar{v})_j \\ &= \frac{1}{n} \sum_{r=1}^n (v_r - \mu)_i (v_r - \mu)_j + (\mu - \bar{v})_i (\mu - \bar{v})_j + (v_r - \bar{v})_i (\mu - \bar{v})_j + (\mu - \bar{v})_i (v_r - \bar{v})_j \\ &= (\bar{v} - \mu)_i (\bar{v} - \mu)_j + \frac{1}{n} \sum_{r=1}^n (v_r - \mu)_i (v_r - \mu)_j \end{aligned}$$

Since the distinct coordinates of $(\bar{v} - \mu)$ are distributed as independent copies of $\mathcal{N}(0, 1/n)$, with probability at least 0.95 (say), all their absolute values are less than $1/2k$. Indeed,

$$\Pr_{X \sim \mathcal{N}(0, 1/n)} \left[|X| \geq \frac{1}{2k} \right] \leq \exp(-\Omega(n/4k^2)) = O\left(\frac{1}{d}\right),$$

and we can take a union bound over the d coordinates (X above is one of the coordinates of $\bar{v} - \mu$). Given that this event happens, note that the first term of Σ_{ij} can be bounded using Cauchy-Schwarz by $\frac{1}{2} \left((\bar{v} - \mu)_i^2 + (\bar{v} - \mu)_j^2 \right) \leq 1/2k$, so it now suffices to bound the second term. That is, shifting the points v_r back by μ to get w_r , we wish to show that given $w_1, \dots, w_n \sim \mathcal{N}(0, \text{Id})$,

$$\left\| \frac{1}{n} \sum_{r=1}^n w_r w_r^\top - \text{Id} \right\|_\infty \leq \frac{1}{2k}$$

with high probability. For the diagonal entries, this probability is

$$\Pr \left[\frac{1}{n} \sum_{r=1}^n (w_r)_i^2 - 1 \geq \frac{1}{2k} \right] = \Pr_{X \sim \mathcal{N}(0, \text{Id}_n)} \left[\|X\|^2 - n \geq \frac{n}{2k} \right].$$

Applying standard concentration bounds again, this is $\exp(-\Omega(n/k^2)) = O(1/d^2)$. For off-diagonal entries ij , this probability is

$$\Pr \left[\frac{1}{n} \sum_{r=1}^n (w_r)_i (w_r)_j \geq \frac{1}{2k} \right] = \Pr_{X, X' \sim \mathcal{N}(0, \text{Id}_n)} \left[\langle X, X' \rangle \geq \frac{n}{2k} \right] = \Pr_{\substack{X \sim \mathcal{N}(0, \text{Id}_n) \\ X' \sim \mathcal{N}(0, 1)}} \left[\|X\| X' \geq \frac{n}{2k} \right],$$

where the final inequality is because conditioned on X , the projection of X' on the direction of X is distributed as a standard one-dimensional Gaussian. Standard Gaussian tail bounds imply that this is bounded as

$$\begin{aligned} \Pr_{\substack{X \sim \mathcal{N}(0, \text{Id}_n) \\ X' \sim \mathcal{N}(0, 1)}} \left[X' \geq \frac{n}{2k\|X\|} \right] &\leq \mathbb{E}_{X \sim \mathcal{N}(0, \text{Id}_n)} \left[\exp \left(-\Omega \left(\frac{n^2}{k^2 \|X\|^2} \right) \right) \right] \\ &\leq \Pr_{X \sim \mathcal{N}(0, \text{Id}_n)} [\|X\|^2 \geq 2n] + \exp \left(-\Omega \left(\frac{n^2}{k^2 n} \right) \right) \\ &\leq \Pr_{X \sim \mathcal{N}(0, \text{Id}_n)} [\|X\|^2 \geq 2n] + \exp \left(-\Omega \left(\frac{n}{k^2} \right) \right) = O(1/d^2), \end{aligned}$$

where the second inequality is because $\exp(-z)$ is at most 1 for $z \geq 0$, and the final inequality follows exactly like the earlier concentration bound and the fact that $n = \Omega(k^2 \log d)$.

Taking a union bound over all d^2 entries of the matrix completes the proof.

- (c) This part is immediate from (a) and (b). First condition on these events (since both happen with probability 0.99). For simplicity, suppose that the adversary does not reorder the sample, so if a sample is not corrupted, then $v_i = v'_i$. Then, one can choose $w_i = \mathbf{1}[v_i = v'_i]$, and $z_i = v_i$ for all i . The first constraint is trivially satisfied since each $w_i \in \{0, 1\}$, the second because at most εn corruptions are introduced, the third by the definition of w_i and z_i , and the last two by (b).
- (d) By Cauchy-Schwarz, we have

$$\mathcal{S} \vdash_2 \frac{1}{n} \sum_{i=1}^n \langle u, z_i - v_i \rangle^2 \leq 4 \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \langle u, z_i - \bar{z} \rangle^2}_{\text{(I)}} + \underbrace{\langle u, \bar{z} - \mu \rangle^2}_{\text{(II)}} + \underbrace{\frac{1}{n} \sum_{i=1}^n \langle u, \bar{v} - v_i \rangle^2}_{\text{(III)}} \right).$$

First, note that (a) implies that (II) is bounded by ε^2 . To bound (III), set $\Sigma' = (1/n) \sum_{i=1}^n (v_i - \bar{v})(v_i - \bar{v})^\top$. By (b), $|\Sigma' - \text{Id}|_\infty \leq 1/k$. Then, setting $S = \{i : u_i \neq 0\}$ (with $|S| \leq k$),

$$\begin{aligned}
(\text{III}) &= u^\top \Sigma' u \\
&= \sum_{i \in S} u_i^2 \Sigma'_{ii} + \sum_{\substack{i, j \in S \\ i \neq j}} u_i u_j \Sigma'_{ij} \\
&\leq \left(1 + \frac{1}{k}\right) + \sum_{\substack{i, j \in S \\ i \neq j}} \frac{u_i^2 + u_j^2}{2} \cdot \frac{1}{k} \\
&\leq \left(1 + \frac{1}{k}\right) + \sum_{i \in S} \frac{1}{k} = 2 + \frac{1}{k}.
\end{aligned}$$

Similarly, by the SoS constraints, setting $\Sigma'' = (1/n) \sum_{i=1}^n (z_i - \bar{z})(z_i - \bar{z})^\top$, we have $\mathcal{S} \vdash_2 \|\Sigma'' - \text{Id}\|_\infty \leq 1/k$, so

$$\mathcal{S} \vdash_2 (\text{I}) = u^\top \Sigma'' u \leq 2 + \frac{1}{k}.$$

Therefore, putting the pieces together, we get that

$$\mathcal{S} \vdash_2 \frac{1}{n} \sum_{i=1}^n \langle u, z_i - v_i \rangle^2 \leq 4 \left(4 + \frac{2}{k} + \varepsilon^2 + \langle u, z - \mu \rangle^2\right) = O(1) \cdot (1 + \langle u, z - \mu \rangle^2),$$

as desired.

(e) Note that $\mathcal{S} \vdash_2 (z_i - v_i)w_i \mathbf{1}_{v_i=v'_i} = 0$. Consequently

$$\begin{aligned}
\mathcal{S} \vdash_4 \langle u, \bar{z} - \bar{v} \rangle^2 &= \left(\frac{1}{n} \sum_{i=1}^n \langle u, z_i - v_i \rangle \right)^2 \\
&= \left(\frac{1}{n} \sum_{i=1}^n \langle u, (z_i - v_i)(1 - \mathbf{1}_{v_i=v'_i} w_i) \rangle \right)^2 \\
&\leq \left(\frac{1}{n} \sum_{i=1}^n \langle u, z_i - v_i \rangle^2 \right) \left(\frac{1}{n} \sum_{i=1}^n (1 - \mathbf{1}_{v_i=v'_i} w_i)^2 \right),
\end{aligned}$$

where the final inequality follows by Cauchy-Schwarz. Now,

$$\mathcal{S} \vdash_2 \sum_{i=1}^n (1 - \mathbf{1}_{v_i=v'_i} w_i)^2 = \sum_{i=1}^n (1 - \mathbf{1}_{v_i=v'_i} w_i) = \sum_{i=1}^n (1 - w_i) + (1 - \mathbf{1}_{v_i=v'_i}) \underbrace{w_i}_{\mathcal{S} \vdash_2 w_i \leq 1} \leq 2\varepsilon n,$$

where we used the second constraint in the sum-of-squares system to bound the first sum by εn , and that at most εn of the v_i were corrupted to bound the second term by εn . Substituting this back in the earlier string of equations, and using (d), we get that

$$\langle u, \bar{z} - \bar{v} \rangle^2 \leq O(\varepsilon)(1 + \langle u, \bar{z} - \mu \rangle^2)$$

as desired.

(f) Indeed, we have

$$\begin{aligned}\mathcal{S} \vdash_4 \langle u, \bar{z} - \mu \rangle^2 &\leq 2\langle u, \bar{z} - \bar{v} \rangle^2 + 2\langle u, \bar{v} - \mu \rangle^2 \\ &\leq O(\varepsilon)(1 + \langle u, \bar{z} - \mu \rangle^2) + O(\varepsilon^2) \\ \langle u, \bar{z} - \mu \rangle^2 &\leq O(\varepsilon)\end{aligned}$$

as desired. Here, the first inequality is Cauchy-Schwarz, and the second by (a) and (e).

The final algorithm is as follows. We find a pseudoexpectation $\tilde{\mathbb{E}}$ that is feasible for the sum-of-squares relaxation \mathcal{S} , and output $\hat{\mu}$, such that $\hat{\mu}_i = (\tilde{\mathbb{E}} \bar{z})_i$ for the k coordinates with largest $(\tilde{\mathbb{E}} \bar{z})_i$. Just as in usual robust mean estimation, we have that for any k -sparse unit vector u ,

$$\varepsilon \geq \tilde{\mathbb{E}} \langle u, \bar{z} - \mu \rangle^2 \geq \langle u, \tilde{\mathbb{E}} \bar{z} - \mu \rangle^2,$$

so $\langle u, \tilde{\mathbb{E}} \bar{z} - \mu \rangle = O(\sqrt{\varepsilon})$. Let $S = \{i : \hat{\mu}_i \neq 0\}$ and $T = \{i : \mu_i \neq 0\}$. For simplicity, assume that $|S| = |T| = k$ (in general one can add in arbitrary coordinates where $\hat{\mu}_i$ or μ_i is 0). Set $z' = (\tilde{\mathbb{E}} \bar{z})|_{S \setminus T}$, $z'' = (\tilde{\mathbb{E}} \bar{z})|_{S \cap T}$, and $z''' = (\tilde{\mathbb{E}} \bar{z})|_{T \setminus S}$ – observe that $\hat{\mu} = z' + z''$. Similarly, set $\mu'' = \mu|_{S \cap T}$ and $\mu''' = \mu|_{T \setminus S}$. Choosing u as the unit vectors in the directions of z' , $z'' - \mu''$, and $z''' - \mu'''$, we get that each of $\|z'\|$, $\|z'' - \mu''\|$, $\|z''' - \mu'''\|$ is $O(\sqrt{\varepsilon})$. Note that since z consists of the k largest coordinates of $\tilde{\mathbb{E}} \bar{z}$, and $|S \setminus T| = |T \setminus S|$, $\|z'''\| \leq \|z'\| = O(\sqrt{\varepsilon})$. Therefore,

$$\|\hat{\mu} - \mu\| = \|z' + z'' - \mu'' - \mu'''\| \leq \|z'\| + \|z'' - \mu''\| + \|z''' - \mu'''\| + \|z'''\| \leq O(\sqrt{\varepsilon})$$

as desired. □