# Problem Set 4

Samuel B. Hopkins, Amit Rajaraman

Last updated November 11, 2024

Due: 11/27, 11:59pm.
Please typeset your solutions in LaTeX.

**Problem 1** (Sparse robust mean estimation). 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 $n = \widetilde{\Omega}(k^2(\log d)/\varepsilon^2)$ samples $v_1, \ldots, v_n \in \mathbb{R}^d$ are drawn from $\mathcal{N}(\mu, \mathrm{Id})$. Then an adversary alters $\varepsilon n$ of the samples and reorders them arbitrarily. We observe the resulting dataset $v'_1, \ldots, v'_n$. Our goal will be to give an algorithm for estimating $\mu$ 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 $w_1, \ldots, w_n$ and $d$-dimensional variables $z, z_1, \ldots, z_n$
$$w_i^2 = w_i$$
$$\sum_{i=1}^{n} w_i \geq (1 - \varepsilon)n$$
$$w_i(z_i - v'_i) = 0$$
$$\bar{z} = \frac{1}{n} \sum_{i=1}^{n} z_i \ , \ \Sigma = \frac{1}{n} \sum_{i=1}^{n} (z_i - \bar{z})(z_i - \bar{z})^T$$
$$-\frac{1}{k} \leq \Sigma_{ij} \leq \frac{1}{k} \quad \text{for all } i \neq j$$
$$-\frac{1}{k} \leq \Sigma_{ii} - 1 \leq \frac{1}{k} \quad \text{for all } i$$

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

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

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

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

(f) Use part (e) to deduce that
$$\mathcal{S} \vdash_4 \langle u, \overline{z} - \mu \rangle^2 \leq O(\varepsilon).$$

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

**Problem 2.** Unreleased.