

PROBLEM SET 2 SOLUTIONS

Problem 1: Ridge Leverage Scores Bound Low Rank Sensitivities

1. We will denote $A' = [A; \sqrt{\lambda}I]$. Let us consider the i^{th} leverage score of A' for $i \in [n]$. By the definition in the hint, we have that this is

$$a_i^{\top}(A'^{\top}A')^{-1}a'_i = a_i^{\top}(A'^{\top}A)^{-1}a_i.$$

So, we just need to prove that $A'^{\top}A' = A^{\top}A + \lambda I$.

We have that

$$A'^{\top}A' = [A; \sqrt{\lambda}I]^{\top}[A; \sqrt{\lambda}I] = A^{\top}A + (\sqrt{\lambda}I)^2 = A^{\top}A + \lambda I.$$

2. Using the previous part, we know that $\tau_i = \ell_i(A')$ where ℓ_i is the i^{th} leverage score. Using the hint, we can see that

$$\ell_i(A') = \sup_x \frac{(A'x)_i^2}{\|A'x\|_2^2} = \sup_x \frac{(Ax)_i^2}{\|A'x\|_2^2} = \sup_x \frac{(Ax)_i^2}{\|Ax\|_2^2 + \lambda\|x\|_2^2}.$$

3. Recall that the operator norm of a matrix is also the largest singular value of that matrix. So we have

$$\|A - A_{2k}\|_2^2 = \sigma_{2k+1}^2(A) = \sigma_{k+1}^2(A - A_k) \leq \frac{1}{k} \sum_{j=1}^k \sigma_j^2(A - A_k) \leq \frac{\|A - A_k\|_F^2}{k} \leq \lambda.$$

4. (a) We follow the steps in the hint. So, we know from part 2 that we have

$$\tau_i = \sup_x \frac{(Ax)_i^2}{\|Ax\|_2^2 + \lambda\|x\|_2^2}.$$

So, let us expand this. We have

$$\sup_x \frac{(Ax)_i^2}{\|Ax\|_2^2 + \lambda\|x\|_2^2} = \sup_x \frac{(Ax)_i^2}{\|A_{2k}x\|_2^2 + \|(A - A_{2k})x\|_2^2 + \lambda\|x\|_2^2}.$$

Using part 3, we now get that

$$\sup_x \frac{(Ax)_i^2}{\|A_{2k}x\|_2^2 + \|(A - A_{2k})x\|_2^2 + \lambda\|x\|_2^2} \geq \sup_x \frac{(Ax)_i^2}{\|A_{2k}x\|_2^2 + 2\lambda\|x\|_2^2}.$$

- (b) Now, we let $F \in \mathcal{F}_k$ be any rank k subspace. We set H to be the span of the rows of A_{2k} , F , and a_i . Clearly H is at most a $(3k + 1)$ dimensional subspace.

Let $x = P_H(I - P_F)g$ where g is a standard normal Gaussian vector and P_H denotes the projection onto H .

So, we have that $(Ax)_i = a_i^{\top}P_H(I - P_F)g = a_i^{\top}(I - P_F)g$. So, $(Ax)_i$ is distributed as a gaussian with variance $\|a_i^{\top}(I - P_F)\|_2^2$. This is a chi-squared random variable with one degree of freedom with expectation $\|a_i^{\top}(I - P_F)\|_2^2$. Using standard properties of the pdf of a chi-squared random variable we have

$$\Pr[(Ax)_i^2 \geq \|a_i^{\top}(I - P_F)\|_2^2/3] > 1/2.$$

- (c) We have that

$$E[\|A_{2k}x\|_2^2] = E[\|A_{2k}P_H(I - P_F)g\|_2^2] = E[\|A_{2k}(I - P_F)g\|_2^2] \leq \|A(I - P_F)\|_F^2.$$

We also have

$$\begin{aligned} E[\lambda\|x\|_2^2] &= E[\lambda\|P_H(I - P_F)g\|_2^2] \\ &= E[\lambda\|HH^\top(I - P_F)g\|_2^2] = E[\lambda\|H^\top(I - P_F)g\|_2^2]. \end{aligned}$$

H^\top has orthonormal rows, and projections do not increase norms. Therefore, each row of $H^\top(I - P_F)$ has length at most 1, and each row of $H^\top(I - P_F)g$ is a gaussian with variance at most 1. So we have

$$E[\lambda\|H^\top(I - P_F)g\|_2^2] \leq \lambda(3k + 1) \leq 4\|A - A_k\|_F^2 \leq 4\|A(I - P_F)\|_F^2.$$

So using Markov's bound we have

$$\Pr[\|A_{2k}x\|_2^2 + 2\lambda\|x\|_2^2 < 20\|A(I - P_F)\|_F^2] > 1/2.$$

(d) We have with positive probability that there exists an x such that

$$\tau_i \geq \sup_x \frac{(Ax)_i^2}{\|A_{2k}x\|_2^2 + 2\lambda\|x\|_2^2} \geq \frac{1}{60} \frac{\|a_i^\top(I - P_F)\|_2^2}{\|A(I - P_F)\|_F^2}.$$

We proved this for arbitrary F , and so we are done.

Problem 2: Sketching for Second Order Methods

1. We expand the initial expression. So we have that

$$\begin{aligned} &\operatorname{argmin}_{x \in \mathcal{C}} \left(\frac{1}{2} \|SA(x - x^t)\|_2^2 - \langle A^\top(b - Ax^t), x \rangle \right) \\ &= \operatorname{argmin}_{x \in \mathcal{C}} \left(\frac{1}{2} x^\top A^\top S^\top S A x - \frac{1}{2} x^\top A^\top S^\top S A x^t - \frac{1}{2} (x^t)^\top A^\top S^\top S A x - x^\top A^\top b + x^\top A^\top A x^t \right) \\ &= \operatorname{argmin}_{x \in \mathcal{C}} \left(\frac{1}{2} \|SAx\|_2^2 - x^\top A^\top S^\top S A x^t - x^\top A^\top b + x^\top A^\top A x^t \right) \\ &= \operatorname{argmin}_{x \in \mathcal{C}} \left(\frac{1}{2} \|SAx\|_2^2 - x^\top A^\top (S^\top S A x^t + b - A x^t) \right) \\ &= \operatorname{argmin}_{x \in \mathcal{C}} \left(\frac{1}{2} \|SAx\|_2^2 - x^\top A^\top (b - (I - S^\top S) A x^t) \right) \\ &= \operatorname{argmin}_{x \in \mathcal{C}} \left(\frac{1}{2} \|SAx\|_2^2 - \langle A^\top (b - (I - S^\top S) A x^t), x \rangle \right) \end{aligned}$$

2. Let $\Delta = x^* - x^{t+1}$.

$$\langle (SA)^\top S A x^{t+1} - A^\top z, \Delta \rangle \geq 0$$

and

$$\langle A^\top A x^* - A^\top b, (-\Delta) \rangle \geq 0.$$

Expanding the former,

$$(x^{t+1})^\top (SA)^\top (SA) \Delta - b^\top A \Delta + (x^t)^\top A^\top A \Delta - (x^t)^\top A^\top S^\top S A \Delta \geq 0,$$

and expanding the latter,

$$-(x^*)^\top A^\top A \Delta + b^\top A \Delta \geq 0.$$

Adding them together, we have

$$(x^{t+1} - x^t)^\top (A^\top S^\top S A) \Delta \geq (x^* - x^t)^\top A^\top A \Delta.$$

Now we add $(x^t - x^*)^T A^T S^T S A \Delta$ to both sides, and we get

$$(x^{t+1} - x^*)^T A^T S^T S A \Delta \geq (x^* - x^t)^T A^T A \Delta + (x^t - x^*)^T A^T S^T S A \Delta$$

which is the same as

$$-\Delta^T A^T S^T S A \Delta \geq (x^* - x^t)^T A^T (I - S^T S) A \Delta,$$

and rearranging gives

$$|(x^* - x^t)^T A^T (I - S^T S) A \Delta| \geq \|S A \Delta\|_2^2.$$

3. Let $S = S^{t+1}$. Applying Cauchy-Schwarz to the upper bound in the previous part, we have

$$|(x^* - x^t)^T A^T (I - S^T S) A \Delta| \leq \|\Sigma V^T (x^* - x^t)\|_2 \cdot \|\Sigma V^T \Delta\|_2 \cdot \epsilon,$$

where $A = U \Sigma V^T$ in its SVD, and we have used $\|I - U^T S^T S U\|_2 \leq \epsilon$ since S is a subspace embedding for A .

Also, we have

$$\|S A \Delta\|_2^2 \geq (1 - \epsilon) \|A \Delta\|_2^2,$$

also because S is a subspace embedding.

Combining these two bounds and the previous part, we obtain

$$\|A \Delta\|_2 = O(\epsilon) \|A(x^* - x^t)\|_2.$$

The claim now follows inductively across the N iterations, applying the same analysis as above with t .

Problem 3: Block Leverage Scores

1. We have that

$$\mathcal{L}_i(A) = \text{Tr}(A^i (A^\top A)^{-1} (A^i)^\top) = \sum_j A_j^i (A^\top A)^{-1} (A_j^i)^\top = \sum_j \ell_j(A).$$

2. We saw in class how to estimate ℓ_i to within a $(1 \pm \epsilon)$ approximation by taking $\tilde{\ell}_i = |e_i \text{ARG}|_2^2$. So the algorithm $\mathcal{L}_i(A)$ is to compute $\tilde{\ell}_i$ for each $i \in T$ where T is the set of rows of A in A^i and then add them up.

So using the previous part, we have that

$$\tilde{\mathcal{L}}_i(A) = \sum_{i \in T} \tilde{\ell}_i = \sum_{i \in T} \ell_i (1 \pm \epsilon) = (1 \pm \epsilon) \mathcal{L}_i(A).$$

In terms of runtime, note that we only have to compute product ARG once, which takes $(\text{nnz}(A) + d^2) \log n$ time. Computing $|e_i \text{ARG}|_2^2$ for each $i \in [n]$ takes $O(\text{nnz}(A) \log n)$ time.

3. Let $U \Sigma V^\top$ be the SVD decomposition of A . So, we can also rewrite A^i as $U^i \Sigma V^\top$ where U^i consists of the first i rows of U . Since our quantity is scale invariant, finding

$$\sup_X \frac{\|A^i X\|_F^2}{\|A X\|_2^2}$$

is the same as maximizing $\|U^i \Sigma V^\top X\|_F^2$ subject to the constraint that $\|U \Sigma V^\top X\|_2^2 = 1$. Note that U has orthonormal columns and therefore the constraint becomes $\|\Sigma V^\top X\|_2^2 = 1$. Now, let us say that $Y = \Sigma V^\top X$. So, we want to maximize $\|U^i Y\|_F^2$ subject to $\|Y\|_2^2 = 1$. This is maximized when we have $Y = I$. So we have that

$$\sup_X \frac{\|A^i X\|_F^2}{\|A X\|_2^2} = \sum_{r=1}^i \|U_r\|_2^2,$$

giving us the result.

4. We have that

$$\mathcal{L}_i(A) = \sum_j \ell_j(A) = \sum_j |U_j|_2^2 = \|U^i\|_F^2 \geq \|U^i\|_2^2.$$

5. We want to show that $|SAx|_2^2 = (1 \pm \varepsilon)|Ax|_2^2$ for all x . Doing the standard change of variable, this is equivalent to showing that

$$|SUY|_2^2 = (1 \pm \varepsilon)|y|_2^2$$

for all y and U with orthonormal columns. Like in class we will show with high probability that

$$\|U^\top S^\top S U - I\|_2 \leq \varepsilon.$$

We will use Matrix Chernoff. Let us set up our random variables and bound the required terms.

For $j \in [t]$, let i_j be the index that was picked in the j -th trial. Let

$$X_j = I_d - \frac{U_{i_j}^\top U_{i_j}}{q_{i_j}}$$

where U_{i_j} is the block matrix corresponding to the sampled rows in the j -th trial.

We can see that all the X_j are independent copies of a symmetric matrix. Now, let us verify the expectation. So we have

$$E[X_j] = I_d - \sum_i q_i \cdot \frac{U_i^\top U_i}{q_i} = I_d - \sum_i U_i^\top U_i = I_d - U^\top U = I_d - I_d = 0^d.$$

We also have

$$\|X_j\|_2 \leq \|I_d\|_2 + \frac{\|U_{i_j}^\top U_{i_j}\|_2}{q_{i_j}} \leq 1 + \max_i \frac{\|U_i\|_2^2}{q_i} = 1 + \max_i \frac{|U_i|_2^2 \cdot d}{\beta \mathcal{L}_i} \leq 1 + \frac{d}{\beta}$$

where the last inequality follows from the previous part.

Finally we have

$$\begin{aligned} E[X^\top X] &= I_d - 2E\left[\frac{U_{i_j}^\top U_{i_j}}{q_{i_j}}\right] + E\left[\frac{U_{i_j}^\top U_{i_j} U_{i_j}^\top U_{i_j}}{q_{i_j}^2}\right] \\ &= \sum_i \frac{U_i^\top U_i U_i^\top U_i}{q_i} - I_d \preceq \sum_i \frac{\|U_i^\top U_i\|_2 U_i^\top U_i}{q_i} - I_d \preceq \frac{d}{\beta} \sum_i U_i^\top U_i - I_d \leq \left(\frac{d}{\beta} - 1\right)I_d. \end{aligned}$$

So we therefore have that

$$\|E[X^\top X]\|_2 \leq \frac{d}{\beta} - 1.$$

The rest follows from the lecture.

References