

Lecture 10 — Februray 6th

Lecturer: Mert Pilanci

Scribe: Sakshi Namdeo & Ayan Mandal

10.1 Faster Least Squares Optimization

In this lecture, we study the Least Squares problem and see how we can use random projections to reduce the time complexity of the least squares problem. More specifically, we will study about the basic inequality method which can be used to bound the norm of the error between the actual least squares solution and the randomized least squares solution.

Consider the least squares problem where $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$.

$$\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 \quad (10.1.1)$$

The exact solution to the least squares problem is given as $x_{LS} = A^+b$. Using classical methods the solution to this problem is $O(nd^2)$.

Using sketching we can reduce the the dimension of A and therefore the complexity of the least squares problem. Using a Random left sketching matrix $S \in \mathbb{R}^{m \times n}$, where $m \ll n$, form SA and Sb . For an approximate solution we can solve the smaller problem-

$$\min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2 \quad (10.1.2)$$

The solution to the above least squares problem is given as $\tilde{x} = (SA)^+b$. Using classical methods the solution to this problem is $O(md^2)$.

10.1.1 Basic Inequality Method

Given the two least squares problems we want to estimate how close \tilde{x} is to x_{LS} . For this we define the quantity $\Delta = \tilde{x} - x_{LS}$. Since the least squares problem minimizes the cost function we can establish two optimality equations as below -

$$\|Ax_{LS} - b\|_2^2 \leq \|Ax' - b\|_2^2 \text{ for any } x', \text{ i.e. } A^T(Ax_{LS} - b) = 0 \quad (10.1.3)$$

$$\|S(A\tilde{x} - b)\|_2^2 \leq \|S(Ax_{LS} - b)\|_2^2 \quad (10.1.4)$$

Substituting $\tilde{x} = x_{LS} + \Delta$, we can rewrite the second equation as -

$$\|S(A(x_{LS} + \Delta) - b)\|_2^2 \leq \|S(Ax_{LS} - b)\|_2^2 \quad (10.1.5)$$

$$\implies \|S(Ax_{LS} - b)\|_2^2 + \|SA\Delta\|_2^2 - 2b^{\perp T} (S^T S - I)A\Delta \leq \|S(Ax_{LS} - b)\|_2^2; \text{ where } b^{\perp} = Ax_{LS} - b$$

$$\implies \|SA\Delta\|_2^2 \leq 2b^{\perp T} (S^T S - I)A\Delta \quad (10.1.6)$$

The above equation is called the basic inequality. We want to use it to say that $\|A\Delta\|_2^2$ is small. Note that if $S^T S = I$ then $\Delta = 0$. We know that for embeddings that satisfy approximate matrix multiplication $S^T S \approx I$.

To argue the closeness of \tilde{x} and x_{LS} first consider the following metric which measures the worst case distortion caused by S -

$$\max_{\Delta} \left| \frac{\|SA\Delta\|_2^2}{\|A\Delta\|_2^2} - 1 \right|; A = U\Sigma V^T \implies A\Delta = U\Sigma V^T \Delta = Uz \quad (10.1.7)$$

$$\begin{aligned} \max_{\Delta} \left| \frac{\|SA\Delta\|_2^2}{\|A\Delta\|_2^2} - 1 \right| &= \max_z \left| \frac{\|SUz\|_2^2}{\|z\|_2^2} - 1 \right| \\ &= \max_z \left| \frac{z^T U^T S^T S U z}{z^T U^T U z} - 1 \right| \\ &= \max_z \left| \frac{z^T}{\|z\|_2} (U^T S^T S U - I) \frac{z}{\|z\|_2} \right| \\ &= \sigma_{max}^2(U^T S^T S U - I) \end{aligned} \quad (10.1.8)$$

We know that the frobenius norm $\|\cdot\|_F = \sqrt{\sum \sigma_i^2}$. Therefore for any Q, $\sigma_{max}(Q) \leq \|Q\|_F$. Using the above and approximate matrix multiplication we can write

$$\sigma_{max}(U^T S^T S U - I) \leq \|U^T S^T S U - I\|_F \leq \epsilon \|U^T U\|_F^2, \quad (10.1.9)$$

where the second inequality follows from AMM. This is called a Subspace Embedding and σ_{max} can be controlled by choosing an appropriate m (Sampling size in S). In other words we can rescale ϵ to get $\sigma_{max}(U^T S^T S U - I) \leq \epsilon'$ for some m.

$$\begin{aligned} \implies \max_{\Delta} \left| \frac{\|SA\Delta\|_2^2}{\|A\Delta\|_2^2} - 1 \right| &\leq \epsilon' \\ \implies -\epsilon' \|A\Delta\|_2^2 &\leq \|SA\Delta\|_2^2 - \|A\Delta\|_2^2 \leq \epsilon' \|A\Delta\|_2^2 \\ \implies (1 - \epsilon') \|A\Delta\|_2^2 &\leq \|SA\Delta\|_2^2 \end{aligned} \quad (10.1.10)$$

We have now derived the following inequalities -

$$\|SA\Delta\|_2^2 \leq 2b^{\perp T} (S^T S - I) A\Delta \quad (10.1.11)$$

$$(1 - \epsilon') \|A\Delta\|_2^2 \leq \|SA\Delta\|_2^2 \quad (10.1.12)$$

Consider the RHS of the first inequality and $A = U\Sigma V^T \implies UU^T A = UU^T U\Sigma V^T = U\Sigma V^T = A$.

$$\begin{aligned} 2b^{\perp T} (S^T S - I) A\Delta &= 2b^{\perp T} (S^T S - I) UU^T A\Delta \\ &\leq 2\|b^{\perp T} (S^T S - I) UU^T\|_2 \|A\Delta\| \\ &= 2\|b^{\perp T} S^T S UU^T - b^{\perp T} UU^T\|_2 \|A\Delta\| \\ &\leq \frac{\epsilon}{\sqrt{m}} \|b^{\perp}\|_F \|UU^T\|_F \|A\Delta\|; \text{ using Approximate Matrix Multiplication} \\ &\leq \frac{\epsilon}{\sqrt{m}} f(x_{LS}) \sqrt{d} \|A\Delta\|_2; \text{ as } \|b^{\perp}\|_F = \|Ax_{LS} - b\|_2 = f(x_{LS}) \text{ and } \|UU^T\|_F = \sqrt{d} \end{aligned} \quad (10.1.13)$$

From equations 10.1.6, 10.1.12 and 10.1.13 we can write -

$$(1 - \epsilon') \|A\Delta\|_2^2 \leq \|SA\Delta\|_2^2 \leq 2b^{\perp T} (S^T S - I) A\Delta \leq \frac{\epsilon}{\sqrt{m}} f(x_{LS}) \sqrt{d} \|A\Delta\|_2 \quad (10.1.14)$$

$$\implies \|A\Delta\|_2 \leq \frac{\epsilon'}{1 - \epsilon'} f(x_{LS}) \quad (10.1.15)$$

Special Case: Consider a linear system $Ax^* = b$ and an approximate linear system $SA\tilde{x} = Sb$ where x^* is the unique solution without sketching and \tilde{x} is the least squares solution with a sketching matrix. Then,

$$Ax^* - A\tilde{x} = A\Delta \quad (10.1.16)$$

Also consider,

$$\begin{aligned} SA\tilde{x} = Sb \ \& \ SAx^* = Sb \\ \implies SA\Delta &= 0 \end{aligned} \quad (10.1.17)$$

This also comes from the fact that if we sample enough number of times, then the error probability $\rightarrow 0$. From Equation 10.1.14, we get $\|A\Delta\|_2^2 \leq \frac{1}{1-\epsilon'} \|SA\Delta\|_2^2$ (10.1.18) Plugging the LHS in Equation 10.1.15, we get

$$\|A\Delta\|_2^2 = \frac{\epsilon}{1 - \epsilon'} f(x_{LS}) \frac{\sqrt{d}}{\sqrt{m}} \quad (10.1.19)$$

under the assumption $U^T U = U^T S^T S U$.

10.2 Leverage Scores and Basic Inequality

Leverage scores are an important statistical measure to determine correlation between singular vectors of a matrix and the standard basis. They give a feel of how much influence or leverage a row has on the best least squares fit. This also gives information to the extent to which a data point is an outlier or errors. This method finds relevance especially while working in higher dimensional space to get high quality numerical implementation of randomized matrix algorithms.

Consider the Approximate Matrix Multiplication problem for $U^T U$ where

$$\|U^T S^T S U - U^T U\|_F = \|U^T S^T S U - I\|_F \leq \epsilon \quad (10.2.1)$$

for some $\epsilon \in \{0, 1\}$ implies a Least Squares approximation. The S is chosen as a randomized sampling matrix so that the or L-2 norms are preserved with an ϵ -error with high probability.

10.2.1 Importance Sampling

Importance Sampling is a method of sampling from a distribution that over-weights the important regions and hence, helps in variance reduction. We can use importance sampling for for sensitivity analysis and as foundation for some methods of computing normalizing constants of probability densities. In our construction of Importance Sampling, the weights are proportional to the row norms of the matrix U for a data matrix $A = U\Sigma^T$.

The leverage scores are calculated as:

$$\ell_i = \|u_i\|_2^2 \quad (10.2.2)$$

where $\|u_i\|_2^2$ are the norms of row of U . Another way of looking at leverage scores is

$$\sum_i \ell_i = \sum_i U^T U \implies \text{tr} I_d = d \quad (10.2.3)$$

when A is full column rank.

Thus, the sampling probabilities associated with which row can be computed by

$$p_i = \frac{1}{d} \|u_i\|_2^2 \quad (10.2.4)$$

which is the sampling probability distribution defined with replacement and $\sum_i p_i = 1$. This can be defined for a matrix A which maybe non uniform as well.

Note: If A is defined as $A = [I, 0]$ and $m = cd \log d$ and the row probabilities are proportional to the row norms $p_i = \frac{1}{d} \|u_i\|_2^2$ then,

$$\|A\tilde{x} - Ax^*\|_2^2 \leq \epsilon \quad (10.2.5)$$

and is related to a cross validation set of hold out size = 1

Note: Hat Matrix is a projection matrix that maps the vector of response values to the vector of the predicted values. Consider a generalised linear model

$$\begin{aligned} A(A^T A)^{-1} A^T &= U U^T \\ \implies \ell_i^T U U^T \ell_i &= \|u_i\|_2^2 \end{aligned}$$

The computational complexity of this is in the order of $\mathcal{O}(nd^2)$ for computing the leverage scores ℓ_i via SVD.

10.3 Fast Johnson Lindenstrauss Transform

Given $0 < \epsilon < 1$, a set X of m points in \mathbb{R}^N , and a number $n > 8 \ln(m)/\epsilon^2$, there is a linear map $f : \mathbb{R}^N \rightarrow \mathbb{R}^n$ such that

$$(1 - \epsilon) \|u - v\|_2^2 \leq \|f(u) - f(v)\|_2^2 \leq (1 + \epsilon) \|u - v\|_2^2 \quad (10.3.1)$$

for all $u, v \in X$.

An alternative statement for any $0 < \epsilon, \delta < 1/2$ and positive integer d , there exists a

distribution $\mathcal{R}^{n \times d}$ from which the matrix A is drawn such that $n = \mathcal{O}(\log(1/\delta)/\epsilon^2)$ and for any unit length vector $x \in \mathcal{R}^d$ then

$$\mathbb{P}(\|Ax\|_2^2 - 1 > \epsilon) < \delta \quad (10.3.2)$$

This can be obtained from the distributional version by setting $x = (u - v)/\|u - v\|_2$ and $\delta < 1/n^2$ for some pair $u, v \in X$.

The Fast JL Transform promises of computing matrix vector products in less than $\mathcal{O}(nd)$ by deriving distributions across the matrix under consideration.

Hadamard Matrix Let H be a $n \times n$ Hadamard Matrix which is given as

$$H_n = \begin{bmatrix} H_{n-1} & H_{n-1} \\ H_{n-1} & -H_{n-1} \end{bmatrix}$$

and $H_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$. The rows and columns of a Hadamard matrix are orthogonal to each other. and finds application in image signal processing, optical multiplexing, error correcting coding and design & analysis of statistics.

10.3.1 Constructing FJLT Sketching Matrix

To construct a FJLT Sketching Matrix S

Step 1: Construct $H = n \times n$ Hadamard Matrix

Step 2: Generate matrix $D = \mathcal{R}^{n \times n}$ diagonal matrix of random ± 1 sampled from a uniform distribution; known as Rademacher matrix

Step 3: Consider $P = \mathcal{R}^{m \times n} \in \{0, 1\}$ a uniform sub-sampling matrix scaled by \sqrt{n}/\sqrt{m} .

In the problem construction, let $\mathbb{E}[S^T S] = I$. Then,

$$S = PHD \quad (10.3.3)$$

Thus, for $Sx = PHDx$, the P matrix samples the frequencies and the HDx computation can be performed in $\mathcal{O}(n \log n)$ time using Fast Hadamard Transform.

Note: If the Hadamard Matrix H is replaced with a Discrete Fourier Transform Matrix F , then our construction of the sketching matrix becomes $S = PFD$. It works well on matrices with probability density concentrated in a small region as follows from the Heisenberg's principle.