

Randomized Algorithms for Least Squares

Krishna Balasubramanian

March 31, 2020

1 Introduction

In this notes, we will study randomized algorithm for computing ordinary least squares and ridge regression. OLS and Ridge regression are popular techniques for studying *linear relationships* between the covariates ($X^{(i)} \in \mathbb{R}^d$) and response ($Y^{(i)} \in \mathbb{R}$) given n samples. We denote the data matrix, with $X^{(i)}$ as its column, by $\mathbf{X} \in \mathbb{R}^{n \times d}$. Similarly, we denote the response vector, with $Y^{(i)}$ as its entries, by the vector $\mathbf{y} \in \mathbb{R}^n$. Then the OLS estimate of the linear relationship between \mathbf{X} and \mathbf{y} is given by the following problem:

$$\boldsymbol{\beta} = \arg \min_{\mathbf{c} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{c} - \mathbf{y}\|_2^2. \quad (1)$$

In statistics, typically the \mathbf{X} is assumed to be random matrix and a linear model assumption is explicitly made between \mathbf{X} and \mathbf{y} to study the statistical properties of the estimator $\boldsymbol{\beta}$. Our discussion below proceeds by conditioning on all the randomness in the data matrix and the response. That is, we assume \mathbf{X} and \mathbf{y} to be deterministic for the purpose of discussion.

The solution of the problem 1 satisfies the following normal equations:

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}. \quad (2)$$

When \mathbf{X} has full column-rank, then we can write the solution of 1 in closed form as follows:

$$\boldsymbol{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

If we have the following SVD for $\mathbf{X} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top$, then we have

$$\boldsymbol{\beta} = \mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{U}^\top \mathbf{y}.$$

Furthermore, the residual vector $\mathbf{r}_\beta = \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \in \mathbb{R}^n$ is orthogonal to the column space of \mathbf{X} . That is, we have

$$\mathbf{r}^\top \mathbf{X} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{X} = 0$$

Direct methods, that use the closed form solution have a complexity of $\mathcal{O}(nd^2)$. Iterative methods for directly solving the optimization problem above also have similar complexity. When considering the regime of d being small but $d^2 \leq n \leq e^d$, these methods are computationally prohibitive.

2 Sketching

The idea of sketching is simple: One can potentially construct a matrix $\Phi \in \mathbb{R}^{r \times n}$ and solve the following **sketched-OLS** problem instead of OLS:

$$\beta_s = \arg \min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi(\mathbf{X}\mathbf{c} - \mathbf{y})\|_2^2 = \arg \min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi\mathbf{X}\mathbf{c} - \Phi\mathbf{y}\|_2^2$$

Similar to the previous case, we have the closed form expression,

$$\beta_s = ((\Phi\mathbf{X})^\top \Phi\mathbf{X})^{-1} (\Phi\mathbf{X})^\top \Phi\mathbf{y}. \quad (3)$$

Remark 2.0.1. The matrix Φ is called as the **sketching matrix**. It reduces the dimensionality of the problem with respect to the sample size n . The other alternative to sketching is the trivial method: **just throw away part of the samples**. When we instead use the sketching idea and if the matrix Φ is constructed **carefully**, we would have a huge computational saving with very little loss in accuracy of the obtained solution.

A popular way to do Sketched-OLS is described in Algorithm 1. Basically, the sketching matrix in this algorithms is defined by $\Phi = \mathbf{S}^T \mathbf{H} \mathbf{D}$ where $\mathbf{S} \in \mathbb{R}^{n \times r}$ and $\mathbf{H}, \mathbf{D} \in \mathbb{R}^{n \times n}$. Notice that the algorithm fixes $r \approx d \log(n)/\epsilon$. Hence, even when $n \approx e^d$, we have $r \approx d^2/\epsilon$ which provides the computational benefits. Specifically, the sketched-OLS in Equation 3 could be computed in this case in $\mathcal{O}(rd^2) \approx d^4/\epsilon$ running time.

The matrix \mathbf{S} and \mathbf{D} are described in Algorithm 1. **The matrix \mathbf{H} is called as a *normalized Hadamard transform matrix*.** First, the Hadamard transformation matrix is defined recursively as follows.

$$\tilde{\mathbf{H}}_2 = \begin{pmatrix} +1 & +1 \\ +1 & -1 \end{pmatrix} \quad \text{and} \quad \tilde{\mathbf{H}}_n = \begin{pmatrix} \tilde{\mathbf{H}}_{n/2} & \tilde{\mathbf{H}}_{n/2} \\ \tilde{\mathbf{H}}_{n/2} & -\tilde{\mathbf{H}}_{n/2} \end{pmatrix}$$

Then, the normalized Hadamard transformation matrix is defined as $\mathbf{H} = (1/\sqrt{n})\tilde{\mathbf{H}}_n$. The matrix $\mathbf{H}\mathbf{D}$ is called as **randomized Hadamard transform** as \mathbf{D} is random and has nice properties. Specifically, the matrix-vector multiplication $\mathbf{H}\mathbf{D}\mathbf{e}$ could be done in $n \log(n)$ time for any vector \mathbf{e} . **Note: A naive implementation might take $\mathcal{O}(n^2)$ time but the fast Hadamard transform approach outlined here has the above mentioned complexity.** **The matrix \mathbf{S} is a *sub-sampling* matrix which also contributes to the randomness in the algorithm.** Recall that, we assume that \mathbf{X} and \mathbf{y} are deterministic in our discussion. That is, you can think as follows: all of the discussion above is done condition on \mathbf{X} and \mathbf{y} if they are random. So in the algorithm as such, the the only source of randomness in through \mathbf{D} and \mathbf{S} . So, when you get to Fact 2.1.1 and encounter the phrase, **with high probability** we refer to this source of randomness only.

Algorithm 1 The **Sketched-OLS** algorithm

Input: $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$, and an error parameter $\epsilon \in (0, 1)$.

Output: $\beta_s \in \mathbb{R}^d$.

1. Let $r \approx \frac{d \log(n)}{\epsilon}$.
2. Let \mathbf{S} be an empty matrix.
3. **For** $t = 1, \dots, r$ (i.i.d. trials with replacement) **select uniformly at random** an integer from $\{1, 2, \dots, n\}$.
 - **If** i is selected, **then** append the column vector $(\sqrt{n/r}) \mathbf{e}_i$ to \mathbf{S} , where $\mathbf{e}_i \in \mathbb{R}^n$ is the i -th canonical vector.
4. Let $\mathbf{H} \in \mathbb{R}^{n \times n}$ be the normalized Hadamard transform matrix.
5. Let $\mathbf{D} \in \mathbb{R}^{n \times n}$ be a diagonal matrix with

$$\mathbf{D}_{ii} = \begin{cases} +1 & , \text{ with probability } 1/2 \\ -1 & , \text{ with probability } 1/2 \end{cases}$$

6. Compute and return $\beta_s = (\mathbf{S}^T \mathbf{H} \mathbf{D} \mathbf{X})^\dagger \mathbf{S}^T \mathbf{H} \mathbf{D} \mathbf{y}$, where for a matrix \mathbf{A} , $\mathbf{A}^\dagger = (\mathbf{A}^\top \mathbf{A})^{-1}$.
-

2.1 What is a *good* sketching matrix ?

We provide some general statements that need to be satisfied for a matrix Φ to be a “good” sketching matrix. Basically, any matrix Φ that satisfies the two assumptions stated in Lemma 1 is called as a good sketching matrix. Instead of stating the assumption upfront, we outline the proofs from which we can see where the stated assumptions arise from. To fix notation, we denote the residual vectors with respect to OLS solution as $\mathbf{r}_\beta = \mathbf{y} - \mathbf{X}\beta$ and the residual with the estimator β_s as $\mathbf{r}_{\beta_s} = \mathbf{y} - \mathbf{X}\beta_s$.

Theorem 2.1. *Let β be the OLS solution and β_s be the Sketched-OLS solution. Let the sketching matrix Φ satisfy the two conditions stated in Lemma 1. Then we have*

$$\begin{aligned} \|\mathbf{r}_{\beta_s}\|_2^2 &\leq (1 + \epsilon) \|\mathbf{r}_\beta\|_2^2 \\ \|\beta - \beta_s\|_2^2 &\leq \frac{\|\mathbf{r}_\beta\|_2^2 \epsilon}{\sigma_{\min}^2(\mathbf{X})} \end{aligned}$$

$$\begin{aligned}
\|\mathbf{r}_{\beta_s}\|_2^2 &= \|\mathbf{y} - \mathbf{X}\beta + \mathbf{X}\beta - \mathbf{X}\beta_s\|_2^2 \\
&= \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \|\mathbf{X}\beta - \mathbf{X}\beta_s\|_2^2 \\
&= \|\mathbf{r}_\beta\|_2^2 + \|\mathbf{U}\theta\|_2^2 \\
&\leq \|\mathbf{r}_\beta\|_2^2 + \epsilon\|\mathbf{r}_\beta\|_2^2 \quad [\text{by Equation 4}]. \\
&\leq (1 + \epsilon)\|\mathbf{r}_\beta\|_2^2.
\end{aligned}$$

Lemma 1. Consider the optimal value of the following two optimization problems:

$$\min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi \mathbf{X} \mathbf{c} - \Phi \mathbf{y}\|_2^2 \quad \min_{\mathbf{e} \in \mathbb{R}^d} \|\Phi \mathbf{U} \mathbf{e} - \Phi \mathbf{r}_\beta\|_2^2.$$

Then both values are the same, i.e., $\min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi \mathbf{X} \mathbf{c} - \Phi \mathbf{y}\|_2^2 = \min_{\mathbf{e} \in \mathbb{R}^d} \|\Phi \mathbf{U} \mathbf{e} - \Phi \mathbf{r}_\beta\|_2^2$. Furthermore, if a vector θ satisfies:

$$\mathbf{U}\theta = \mathbf{X}(\beta_s - \beta)$$

then, it achieves the minimum of the optimization problem in the right hand side.

Proof. To see the first statement, note that

$$\begin{aligned}
&\min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi \mathbf{X} \mathbf{c} - \Phi \mathbf{y}\|_2^2 \\
&= \min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi \mathbf{X} \mathbf{c} - \Phi(\mathbf{X}\beta + \mathbf{r}_\beta)\|_2^2 \\
&= \min_{\mathbf{d} \in \mathbb{R}^d} \|\Phi \mathbf{X}(\beta + \mathbf{d}) - \Phi(\mathbf{X}\beta + \mathbf{r}_\beta)\|_2^2 \\
&= \min_{\mathbf{d} \in \mathbb{R}^d} \|\Phi \mathbf{X} \mathbf{d} - \Phi \mathbf{r}_\beta\|_2^2 \\
&= \min_{\mathbf{e} \in \mathbb{R}^d} \|\Phi \mathbf{U} \mathbf{e} - \Phi \mathbf{r}_\beta\|_2^2
\end{aligned}$$

In the above derivation, in the last step, we are allowed to replace \mathbf{X} with \mathbf{U} as singular vectors in \mathbf{U} form a basis for the column space of \mathbf{X} . Hence, we just have different combination of the vectors in \mathbf{U} instead of vectors in \mathbf{X} . This proves the first statement.

To see the next statement, take the objective function of the optimization problems in the right hand side and note that

$$\begin{aligned}
\|\Phi \mathbf{U} \theta - \Phi \mathbf{r}_\beta\|_2^2 &= \|\Phi \mathbf{X}(\beta_s - \beta) - \Phi \mathbf{r}_\beta\|_2^2 \\
&= \|\Phi \mathbf{X} \beta_s - \Phi \mathbf{X} \beta - \Phi \mathbf{r}_\beta\|_2^2 \\
&= \|\Phi \mathbf{X} \beta_s - \Phi(\mathbf{X} \beta + \mathbf{r}_\beta)\|_2^2 \\
&= \|\Phi \mathbf{X} \beta_s - \Phi \mathbf{y}\|_2^2 \\
&= \min_{\mathbf{c} \in \mathbb{R}^d} \|\Phi \mathbf{X} \mathbf{c} - \Phi \mathbf{y}\|_2^2 \\
&= \min_{\mathbf{e} \in \mathbb{R}^d} \|\Phi \mathbf{U} \mathbf{e} - \Phi \mathbf{r}_\beta\|_2^2
\end{aligned}$$

Hence, we have shown that the vector $\boldsymbol{\theta}$ indeed is the vector that achieves the minimum value of the optimization problem on the right hand side. \square

Lemma 2. *Assume the following*

- *Isometry assumption 1:*

$$\sigma_{\min}^2(\Phi\mathbf{U}) \geq \frac{1}{\sqrt{2}}$$

- *Isometry assumption 2:*

$$\|(\Phi\mathbf{U})^\top \Phi \mathbf{r}_\beta\|_2^2 \leq \frac{\epsilon \|\mathbf{r}_\beta\|_2^2}{2}$$

Then, we have the following bound for θ , from Lemma 1.

$$\|\boldsymbol{\theta}\|_2^2 \leq \epsilon \|\mathbf{r}_\beta\|_2^2 \quad (4)$$

Proof. Because $\boldsymbol{\theta}$ is the solution of the optimization problem $\min_{\mathbf{e} \in \mathbb{R}^d} \|\Phi\mathbf{U}\mathbf{e} - \Phi\mathbf{r}_\beta\|_2^2$, it should satisfy the normal equation:

$$(\Phi\mathbf{U})^\top \Phi\mathbf{U}\boldsymbol{\theta} = (\Phi\mathbf{U})^\top \Phi\mathbf{r}_\beta$$

Hence, under our assumption, we have

$$\frac{\|\boldsymbol{\theta}\|_2^2}{2} \leq \|(\Phi\mathbf{U})^\top \Phi\mathbf{U}\boldsymbol{\theta}\|_2^2 = \|(\Phi\mathbf{U})^\top \Phi\mathbf{r}_\beta\|_2^2 \leq \frac{\epsilon \|\mathbf{r}_\beta\|_2^2}{2}$$

\square

Remark 2.1.1. Any sketching matrix Φ that satisfies the two assumptions stated in Lemma 2 is called as a “good” sketching matrix. The numbers in the assumptions are fixed just arbitrarily, i.e., instead of $\sigma_{\min}^2(\Phi\mathbf{U}) \geq \frac{1}{\sqrt{2}}$, we can have $\sigma_{\min}^2(\Phi\mathbf{U}) > 0$ and things will still work.

Remark 2.1.2. The assumptions are called as **isometry** assumptions because of a reason. A matrix or a transformation is called as isometric, if it preserves norms. That is, $\|\mathbf{A}\mathbf{e}\|_2^2 \approx \|\mathbf{e}\|_2^2$ for all \mathbf{e} or **most** vectors \mathbf{e} . Sketching is exactly doing something like this, while reducing the dimension as well. The two assumptions are basically a way to characterize when the sketching matrix is isometric.

Fact 2.1.1. The Sketching matrix $\Phi = \mathbf{S}^T \mathbf{H} \mathbf{D}$ constructed in the previous section, satisfies the properties of a good **sketching matrix with high probability** (as the matrix is a random matrix). If you are curious how to show this fact, you may refer to [Mah16].

References

[Mah16] Michael W Mahoney, *Lecture notes on randomized linear algebra*, arXiv preprint arXiv:1608.04481 (2016).