

Lecture 16 — October 13, 2021

Sebastien Roch, UW-Madison

Scribe: Yuchen Zeng, Sebastien Roch

1 Overview

In the last lecture we (a) studied least square estimator of linear regression model, and (b) proved the upper bound of mean squared error of estimated predicted variables.

In this lecture we (a) derive the concentration bound of the least square estimator of linear regression model, and (b) investigate a more general case: non-linear models.

2 Review of eigenvalues

We begin by reviewing a few basic factors about matrices.

Definition 1 (Induced Norm). *The 2-norm of a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ is*

$$\|\mathbf{A}\|_2 = \max_{\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^m} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\mathbf{x} \in \mathbb{S}^{m-1}} \|\mathbf{A}\mathbf{x}\|.$$

Theorem 2 (Spectral Theorem). *Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a symmetric matrix, that is, $\mathbf{A}^\top = \mathbf{A}$. Then \mathbf{A} has d orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$ with corresponding (not necessarily distinct) real eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. In matrix form, this is written as the matrix factorization*

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top = \sum_{i=1}^d \lambda_i \mathbf{q}_i \mathbf{q}_i^\top,$$

where \mathbf{Q} has columns $\mathbf{q}_1, \dots, \mathbf{q}_d$ and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$. We refer to this factorization as a spectral decomposition of \mathbf{A} .

Definition 3 (Rayleigh Quotient). *Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a symmetric matrix. The Rayleigh quotient is defined as*

$$\mathcal{R}_{\mathbf{A}}(\mathbf{u}) = \frac{\langle \mathbf{u}, \mathbf{A}\mathbf{u} \rangle}{\langle \mathbf{u}, \mathbf{u} \rangle},$$

which is defined for any $\mathbf{u} \neq \mathbf{0}$ in \mathbb{R}^d .

Theorem 4 (Courant-Fischer). *Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a symmetric matrix with spectral decomposition $\mathbf{A} = \sum_{i=1}^d \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$ where $\lambda_1 \geq \dots \geq \lambda_d$. For each $k = 1, \dots, d$, define the subspace*

$$\mathcal{V}_k(\mathbf{C}) = \text{span}(\mathbf{v}_1(\mathbf{C}), \dots, \mathbf{v}_k(\mathbf{C})) \text{ and } \mathcal{W}_{d-k+1}(\mathbf{C}) = \text{span}(\mathbf{v}_k(\mathbf{C}), \dots, \mathbf{v}_d(\mathbf{C})).$$

Then, for all $k = 1, \dots, d$,

$$\lambda_k = \min_{\mathbf{u} \in \mathcal{V}_k} \mathcal{R}_{\mathbf{A}}(\mathbf{u}) = \max_{\mathbf{u} \in \mathcal{W}_{d-k+1}} \mathcal{R}_{\mathbf{A}}(\mathbf{u}).$$

Furthermore we have the following min-max formulas, which do not depend on the choice of spectral decomposition, for all $k = 1, \dots, d$,

$$\lambda_k = \max_{\dim(\mathcal{V})=k} \min_{\mathbf{u} \in \mathcal{V}} \mathcal{R}_{\mathbf{A}}(\mathbf{u}) = \min_{\dim(\mathcal{W})=d-k+1} \max_{\mathbf{u} \in \mathcal{W}} \mathcal{R}_{\mathbf{A}}(\mathbf{u}).$$

Lemma 5 (Weyl's Inequality). *Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times d}$ be symmetric matrices. Then, for all $j = 1, \dots, d$,*

$$\max_{j \in [d]} |\lambda_j(\mathbf{B}) - \lambda_j(\mathbf{A})| \leq \|\mathbf{B} - \mathbf{A}\|_2,$$

where $\|\mathbf{C}\|_2$ is the induced 2-norm of \mathbf{C} .

Theorem 6 (Davis-Kahan). *Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times d}$ be symmetric matrices. For an $i \in \{1, \dots, d\}$, assume that*

$$\delta := \min_{j \neq i} |\lambda_i(\mathbf{A}) - \lambda_j(\mathbf{A})| > 0,$$

then

$$\min_{s \in \{+1, -1\}} \|\mathbf{v}_i(\mathbf{A}) - s\mathbf{v}_i(\mathbf{B})\|^2 \leq \frac{8 \|\mathbf{A} - \mathbf{B}\|_2^2}{\delta^2}.$$

3 Controlling the coefficients

Assume we have n observation pairs $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^p$, $y_i = \mathbf{x}_i^\top \boldsymbol{\theta}^* + \varepsilon_i$, $\boldsymbol{\theta}^* \in \mathbb{R}^p$ is unknown, and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$ is of mean zero. Denote the $\mathbf{y} = (y_1, \dots, y_n)$, \mathbf{X} as the matrix with \mathbf{x}_i^\top on row i . Therefore, the linear regression model can be written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon}.$$

The least square estimator of $\boldsymbol{\theta}^*$ is defined as $\widehat{\boldsymbol{\theta}}^{\text{LS}} \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$. Then $\widehat{\boldsymbol{\theta}}^{\text{LS}} = \mathbf{X}\mathbf{y}$ is the minimum norm solution in the set of least square estimators, which is generally unique.

In last lecture, we proved the following theorem, which provides an upper bound for the mean square error of the estimations.

Theorem 7. *Suppose $\boldsymbol{\varepsilon}$ is subgaussian, i.e., $\|\boldsymbol{\varepsilon}\|_{\psi_2} \leq K$. Then w.p.1 - δ , we have*

$$\text{MSE}(\mathbf{X}\widehat{\boldsymbol{\theta}}^{\text{LS}}) \lesssim \frac{K^2}{n} (\text{rank}(\mathbf{X}) + \log(\frac{1}{\delta})),$$

where $\text{MSE}(\mathbf{X}\widehat{\boldsymbol{\theta}}^{\text{LS}}) = \frac{1}{n} \|\mathbf{X}\widehat{\boldsymbol{\theta}}^{\text{LS}} - \mathbf{X}\boldsymbol{\theta}^*\|_2^2$.

We derive a corollary which provides an concentration bound for $\widehat{\boldsymbol{\theta}}^{\text{LS}}$, based on the theorem above.

Corollary 8. *Suppose $p \leq n$ and $\mathbf{B} = \frac{\mathbf{X}^\top \mathbf{X}}{n}$ has rank p . Then w.p.1 - δ ,*

$$\|\widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{\theta}^*\|_2^2 \lesssim \frac{1}{\lambda_{\min}(\mathbf{B})} \frac{K^2}{n} (p + \log(\frac{1}{\delta})).$$

Proof. By Courant-Fischer characterization, we have

$$\lambda_{\min}(\mathbf{B}) = \min_{\mathbf{u} \in \mathbb{R}^p} \frac{\langle \mathbf{u}, \mathbf{B}\mathbf{u} \rangle}{\langle \mathbf{u}, \mathbf{u} \rangle}. \quad (1)$$

Note that

$$\frac{1}{n} \left\| \mathbf{X}\widehat{\boldsymbol{\theta}}^{\text{LS}} - \mathbf{X}\boldsymbol{\theta}^{\star} \right\|_2^2 = (\widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{\theta}^{\star})^{\top} \frac{\mathbf{X}^{\top} \mathbf{X}}{n} (\widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{\theta}^{\star}) = (\widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{\theta}^{\star})^{\top} \mathbf{B} (\widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{\theta}^{\star}). \quad (2)$$

Combining (1) and (2) yields

$$\left\| \widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{\theta}^{\star} \right\|_2^2 \leq \frac{1}{\lambda_{\min}(\mathbf{B})} \text{MSE}(\mathbf{X}\widehat{\boldsymbol{\theta}}^{\text{LS}}).$$

Plugging the upper bound of $\text{MSE}(\mathbf{X}\widehat{\boldsymbol{\theta}}^{\text{LS}})$ in Theorem 7 completes the proof. \square

4 Modeling non-linearities and an oracle bound

In this section, we extend the results to some nonlinear models. Consider a dictionary $\mathcal{H} = \{\phi_1, \dots, \phi_M\}$, where $\phi_i : \mathbb{R}^p \rightarrow \mathbb{R}$ is a map, for $i = 1, \dots, M$. Then we can consider functions of \mathbf{x} which can be written as $\sum_{i=1}^M \theta_i \phi_i(\mathbf{x})$. Given n observation pairs $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, we minimize the mean square error and obtain the least square estimator

$$\widehat{\boldsymbol{\theta}}^{\text{LS}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^M} \|\mathbf{y} - \Phi\boldsymbol{\theta}\|_2^2,$$

where the i th row of Φ is $\Phi_i^{\top} = (\phi_1(\mathbf{x}_i), \dots, \phi_M(\mathbf{x}_i))$.

Next, we provide a quick example to motivate the analysis on nonlinear model.

Example 1. Suppose we observe a quadratic pattern between \mathbf{x} and y . We select the features to be $\phi_1(x) = 1, \phi_2(x) = x, \phi_3(x) = x^2$. Then we have $\widehat{\boldsymbol{\theta}}^{\text{LS}} \in \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n (y_i - (\theta_1 + \theta_2 \mathbf{x}_i + \theta_3 \mathbf{x}_i^2))^2$.

Lastly, we present a theorem, which provides an upper bound for mean square error of the estimation under non-linear models' cases—without assuming that the data is generated by the class of models used to fit it. This type of bound is called an *oracle bound*, and we will have more to say about them later in the course.

Theorem 9 (See Theorem 3.3 in [1]). Assume the data $y_i = f(\mathbf{x}_i) + \varepsilon_i$ with $\|\varepsilon\|_{\psi_2} \leq K$, for all $i = 1, \dots, n$. Fix the dictionary \mathcal{H} . Then w.p. $1 - \delta$,

$$\text{MSE}(\Phi\widehat{\boldsymbol{\theta}}^{\text{LS}}) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \text{MSE}(\Phi\boldsymbol{\theta}) + \frac{cK^2}{n} (M + \log(\frac{1}{\delta})),$$

where $\text{MSE}(\Phi\boldsymbol{\theta}) = \frac{1}{n} \|\Phi\boldsymbol{\theta} - \mathbf{f}\|_2^2$, and $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$.

Proof. A solution to $\inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \text{MSE}(\Phi\boldsymbol{\theta})$ is given by $\bar{\boldsymbol{\theta}} = \Phi^+ \mathbf{f}$, which corresponds to projecting \mathbf{f} onto the column space of Φ . By the properties of the orthogonal projection, for any other vector $\Phi\boldsymbol{\theta}$ in the span of the columns of Φ , we have the orthogonality

$$\langle \Phi\boldsymbol{\theta} - \Phi\bar{\boldsymbol{\theta}}, \mathbf{f} - \Phi\bar{\boldsymbol{\theta}} \rangle = 0.$$

That applies in particular to the choice $\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}^{\text{LS}}$. Hence, by Pythagoras,

$$\begin{aligned} \frac{1}{n} \left\| \Phi \widehat{\boldsymbol{\theta}}^{\text{LS}} - \boldsymbol{f} \right\|_2^2 &= \frac{1}{n} \left\| \Phi \widehat{\boldsymbol{\theta}}^{\text{LS}} - \Phi \bar{\boldsymbol{\theta}} + \Phi \bar{\boldsymbol{\theta}} - \boldsymbol{f} \right\|_2^2 \\ &= \frac{1}{n} \left\| \Phi \bar{\boldsymbol{\theta}} - \boldsymbol{f} \right\|_2^2 + \frac{1}{n} \left\| \Phi \widehat{\boldsymbol{\theta}}^{\text{LS}} - \Phi \bar{\boldsymbol{\theta}} \right\|_2^2. \end{aligned}$$

The claim then follows from the definition of $\bar{\boldsymbol{\theta}}$ and our previous bound on the MSE of linear regression. \square

References

- [1] Philippe Rigollet and Jan-Christian Hutter, *High Dimensional Statistics*, Lecture Notes, 2019.
<http://www-math.mit.edu/~rigollet/PDFs/RigNotes17.pdf>