SDS 387
Linear Models

Fall 2024

Lecture 17 - Thu, Oct 29, 2024

Instructor: Prof. Ale Rinaldo

- Last time : convergence of gradient descent when $\text{rank}(\bar{\Phi}) = d < n$
  We saw that convergence to the minimum is slower when
  $$\hat{\Sigma} = \frac{\Phi^T \hat{\Phi}}{n} \quad \text{is rank-deficient} \quad (\text{has rank} < d).$$

- Classically, it is always assumed that $\text{rank}(\hat{\Sigma}) = d$. But what
  if $d_{\min}(\hat{\Sigma}) = 0$ ?

- Suppose that $\underset{n \times d}{\bar{\Phi}}$ has more columns than rows $(d > n)$.
  What happens?

  $\hat{\beta}$ is still obtained as solution to the normal equations

  $$\Phi^T \Phi \, \hat{\beta} = \Phi^T y \qquad \maltese$$

  But now there are infinitely many solutions! That is
  if say $\hat{\beta}$ solves $\maltese$, then $\hat{\beta} + u$ is also a
  solution for every $u \in \text{kernel}(\bar{\Phi})$

  $\boxed{1}$

- Furthermore, for any solution $\tilde{\beta}$ to $\bigstar$, we have that

$$\Phi \tilde{\beta} = Y$$

$$\hookrightarrow \quad \text{we } \underline{\text{interpolate}} \text{ the data (overfitting)}$$

<span style="color:blue">HW</span>

- Among the infinitely many solutions, one is somewhat
= canonical: it is the one with smallest Euclidean norm!
  It can be calculated using Moore-Penrose pseudo-inverse:

<span style="color:blue">→ pseudo-inverse</span>

$$\hat{\beta}_{MN} = (\Phi^T \Phi)^+ \Phi^T Y$$

$\downarrow$
min-norm

where for a matrix $A$ (m×n) its Moore-Penrose pseudo-inverse
is $A^+$ (n×m) a unique matrix satisfying the conditions:

i) $AA^+A = A$     ( $AA^+$ maps columns of $A$ to themselves, it is an identity on $C(A)$ )

ii) $A^+AA^+ = A^+$

iii) $AA^+$ (m×m) and $A^+A$ (n×n) are symmetric

Notice that $AA^+$ and $A^+A$ are idempotent (see properties i) and ii)) and symmetric. So, $AA^+$ is the orthogonal projector onto $C(A)$ and $A^+A$ is the orthogonal projection onto $C(A^T)$, the row space of $A$.

Useful identities $\qquad A^+ = (A^TA)^+A = A^T(AA^T)^+$

$$\hookrightarrow \hat{\beta}_{MN} = \Phi^+ Y$$

If $\quad A^TA \quad$ is $\quad$ invertible $\qquad A^+ = (A^TA)^{-1}A$

$\qquad AA^T \quad$ is $\quad$ invertible $\qquad A^+ = A^T(AA^T)^{-1}$

- If $\quad \underset{m \times n}{A} = \underset{m \times r}{U} \underset{r \times r}{\Sigma} \underset{r \times n}{V^T} \qquad r = \text{rank}(A) \qquad$ then

$$\underset{n \times m}{A^+} = V \Sigma^{-1} U^T$$

- So for regression :

$$\hat{\beta}_{MN} = \Phi^+ Y = \text{argmin} \left\{ \|\beta\| \text{ s.t. } \Phi\beta = Y \right\}$$

Fact : $\quad$ gradient descent initialized at $0$ (or at any point in the row-span of $\Phi$) converges to $\hat{\beta}_{MN}$

PROPERTIES OF OLS IN FIXED DESIGN SETTINGS AND ASSUMING WELL-SPECIFIE MODEL

From now on let's assume that the data are of the form

$$Y_i = \Phi_i^T \beta^* + \varepsilon_i \qquad \text{where}$$

$$i = 1, \ldots, n$$

where $\varepsilon_1, \ldots, \varepsilon_n \overset{iid}{\sim} (0, \sigma^2)$

$\Phi_1, \ldots, \Phi_n$ are deterministic vectors in $\mathbb{R}^b$

$\downarrow$

Assume

$\text{rank}(\Phi) = d$ $\longleftarrow$ $\underset{n \times d}{\Phi}$ has $\Phi_i^T$ as its $i$th row

Note: if we further assume that $\varepsilon_1, \ldots, \varepsilon_n \overset{iid}{\sim} N(0, \sigma^2)$ then the likelihood of $Y_1, \ldots, Y_n$ is

$$\left( \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \right)^d \exp\left\{ - \frac{\| Y - \Phi \beta^* \|^2}{2\sigma^2} \right\}$$

and

$\hat{\beta}$ is the MLE of $\beta^*$

$\swarrow$

OLS

Next, for any $\beta \in \mathbb{R}^d$, the risk of $\beta$ is

$$R(\beta) = \mathbb{E}_Y \left[ \frac{\| Y - \Phi \beta \|^2}{n} \right] = \mathbb{E}_\varepsilon \left[ \frac{1}{n} \| \Phi(\beta^* - \beta) + \varepsilon \|^2 \right]$$

$\underset{\text{expectation wrt } Y}{\vee}$

$$= (\beta^* - \beta)^T \underbrace{\frac{\Phi^T \Phi}{n}}_{} (\beta^* - \beta) + \mathbb{E}_\varepsilon \left[ \underbrace{\frac{\| \varepsilon \|^2}{n}}_{= \sigma^2} \right]$$

④

$$= \| \beta^* - \beta \|_{\hat{\Sigma}}^2 + \sigma^2 \qquad = \| \beta^* - \beta \|_{\hat{\Sigma}}^2 + R(\beta^*)$$

$$\underbrace{\phantom{xxxx}}_{\text{estimation error}} \qquad = R(\beta^*)$$

irreducible error
(smallest possible risk)

The quantity $R(\beta) - R(\beta^*) \geq 0$ is the <u>excess risk</u>

So, let's look at the expected excess risk of $\tilde{\beta} \to$ any estimator of $\beta^*$

$$\mathbb{E}\left[ R(\tilde{\beta}) \right] - R(\beta^*) = \mathbb{E}\left[ \left( \beta^* - \tilde{\beta} \right)^T \hat{\Sigma} \left( \beta^* - \tilde{\beta} \right) \right]$$

$$= \quad \text{add and subtract} \quad \mathbb{E}[\tilde{\beta}]$$

$$= \quad \left( \beta^* - \mathbb{E}[\tilde{\beta}] \right)^T \hat{\Sigma} \left( \beta^* - \mathbb{E}[\tilde{\beta}] \right) \quad +$$

$$\mathbb{E}\left[ \left( \tilde{\beta} - \mathbb{E}[\tilde{\beta}] \right)^T \hat{\Sigma} \left( \tilde{\beta} - \mathbb{E}[\tilde{\beta}] \right) \right]$$

$$+ \quad 2\, \mathbb{E}\left[ \underline{\left( \tilde{\beta} - \mathbb{E}[\tilde{\beta}] \right)^T \hat{\Sigma} \left( \beta^* - \mathbb{E}[\tilde{\beta}] \right)} \right]$$

$$= 0 \quad \text{by} \quad \text{linearity of } \mathbb{E}[\cdot]$$

$$= \quad \mathbb{E}\left[ \| \tilde{\beta} - \mathbb{E}[\tilde{\beta}] \|_{\hat{\Sigma}}^2 \right] \quad + \quad \| \beta^* - \mathbb{E}[\tilde{\beta}] \|_{\hat{\Sigma}}^2$$

$$\underbrace{\phantom{xxxxxxxxxxxx}}$$

Variance term for $\tilde{\beta}$

Bias term
($=0$ if $\mathbb{E}[\tilde{\beta}] = \beta^*$)

$\downarrow$

bias - variance decomposition of excess risk

⑤

If we choose to use $\hat{\beta}$ the OLS estimator, then

i) $\mathbb{E}[\hat{\beta}] = \beta^*$

ii) $\text{Var}[\hat{\beta}] = \frac{\sigma^2}{n}\hat{\Sigma}^{-1}$

Pf/ $\mathbb{E}[\hat{\beta}] = (\Phi^T\Phi)^{-1}\Phi^T\underbrace{\mathbb{E}[Y]}_{\Phi\beta^*} = \beta^*$    because $\Phi^T\Phi$ is invertible

$$\text{Var}[\hat{\beta}] = \text{Var}\left[(\Phi^T\Phi)^{-1}\Phi^T Y\right] \qquad \left(\begin{array}{l}\text{Var}[AY] = \\ A\,\text{Var}[Y]\,A^T\end{array}\right)$$

$$= (\Phi^T\Phi)^{-1}\Phi^T\underbrace{\text{Var}[Y]}_{\sigma^2 I}\Phi(\Phi^T\Phi)^{-1}$$

$$= \sigma^2(\Phi^T\Phi)^{-1} = \frac{\sigma^2}{n}\hat{\Sigma}^{-1} \quad \blacksquare$$

Using these facts, we can establish that

$$\mathbb{E}\left[R(\hat{\beta})\right] - R(\beta^*)] = \sigma^2\frac{d}{n} \quad \to 0 \quad \text{if} \quad d = o(n)$$

Pf/ Because $\mathbb{E}[\hat{\beta}] = \beta^*$ we only need to analyze the variance term:

$$\mathbb{E}\left[\|\hat{\beta} - \beta^*\|_{\hat{\Sigma}}^2\right] = \mathbb{E}\left[\|(\Phi^T\Phi)^{-1}\Phi^T(\Phi\beta^* + \varepsilon) - \beta^*\|_{\hat{\Sigma}}^2\right]$$

$$= \mathbb{E}\left[\|\beta^* + (\Phi^T\Phi)^{-1}\Phi^T\varepsilon - \beta^*\|_{\hat{\Sigma}}^2\right]$$

$$= \mathbb{E}\left[\|\frac{\hat{\Sigma}^{-1}\Phi^T}{n}\varepsilon\|_{\hat{\Sigma}}^2\right]$$

$$= \mathbb{E}\left[ \varepsilon^T \frac{\Phi}{n} \hat{\Sigma}^{-1} \hat{\Sigma} \hat{\Sigma}^{-1} \frac{\Phi^T}{n} \varepsilon \right]$$

$$= \frac{1}{n} \mathbb{E}\left[ \varepsilon^T \Phi (\Phi^T \Phi)^{-1} \Phi^T \varepsilon \right]$$

$$= \mathbb{E}\left[ \frac{1}{n} \varepsilon^T H \varepsilon \right]$$

$$= \frac{1}{n} \mathbb{E}\left[ tr\left( H \varepsilon \varepsilon^T \right) \right]$$

$$= \frac{1}{n} tr\left( H \underbrace{\mathbb{E}[\varepsilon \varepsilon^T]}_{\sigma^2 I_n} \right)$$

$$= \frac{\sigma^2}{n} tr(H)$$

$$\stackrel{\text{Exercise!}}{=} \frac{\sigma^2}{n} d \quad \text{☏}$$

<u>Remarks</u>

i) This rate is "optimal", $\left(\text{in } d, n \text{ and } \sigma^2\right)$

ii) an analogour bound holds with high probability but it requires more advanced tools

iii) This bound implies that the risk of $\hat{\beta}$ is

$$\mathbb{E}\left[R(\hat{\beta})\right] = \mathbb{E}_{Y_{new} \in \mathbb{R}^n} \left[ \frac{\| Y_{new} - \Phi \hat{\beta} \|^2}{n} \right]$$

fresh ↙
new set of samples

$$= \sigma^2 \left( 1 + \frac{d}{n} \right)$$

This is called the out-of-sample risk or expected test error

⑦

If we insted we compute the in-sample
expected risk :

$$\mathbb{E}\left[\hat{R}(\hat{\beta})\right] = \mathbb{E}_Y\left[\frac{\|Y - \Phi\hat{\beta}\|^2}{n}\right]$$

↙ 
expected training
error

expectation wrt to
same data used to compute $\hat{\beta}$

$$= \sigma^2\left(1 - \frac{d}{n}\right)$$

⑧