SDS 387
Linear Models

Fall 2024

Lecture 22 - Tue, Nov 14, 2024

Instructor: Prof. Ale Rinaldo

- HW 4 is now due Thu, Nov 21, by midnight and project report
  due on Monday Nov 18, by midnight.

- Last time: we finally finished discussing OLS properties under
  fixed-design, well-specified model.

- We are now dropping both assumptions. The lack of linearity
  and randomness of the design matrix create extra-complications
  in particular an increase in variability.

- Let's assume for now that $\underset{n \times d}{\mathbb{B}}$ is random but the model
  is well-specified. This means that our observations are
  $$(Y_1, \Phi_1), \ldots, \left(Y_n, \Phi_n\right) \underset{\sim}{iid} \quad P_{Y, \Phi} \quad \text{in} \quad \mathbb{R} \times \mathbb{R}^d$$
  $$Y_n = \Phi_n^T \beta^* + \varepsilon_n \qquad \text{where} \quad \varepsilon_1, \ldots, \varepsilon_n \mid \Phi_1, \ldots, \Phi_n$$
  $$\underset{\sim}{iid} \quad (0, \sigma^2)$$

  Now the definition of the risk has to be modified:

$$\beta \in \mathbb{R}^d \longmapsto R(\beta) = \mathop{\mathbb{E}}_{Y, \Phi}\left[\left(Y - \Phi^\top \beta\right)^2\right]$$

<span style="color:blue">↓ expectation wrt joint distribution of $Y$ and $\Phi$</span>

**Prop 3.9 in Bach's book**     <span style="color:blue">$\longrightarrow (\beta - \beta^*)^\top \Sigma (\beta - \beta^*)$</span>

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

where $\quad \sigma^2 = R(\beta^*) = \inf_\beta R(\beta) \quad$ and $\quad \Sigma = \mathbb{E}\left[\Phi \Phi^\top\right].$

Pf/   For any $\beta \in \mathbb{R}^d$

$$R(\beta) = \mathbb{E}\left[\left(Y - \Phi^\top \beta\right)^2\right] = \mathbb{E}\left[\left(Y - \Phi^\top \beta^* + \Phi^\top(\beta^* - \beta)\right)^2\right]$$

$$= \mathbb{E}\left[\left(Y - \Phi^\top \beta^*\right)^2\right] + \mathbb{E}\left[\left(\Phi^\top(\beta^* - \beta)\right)^2\right]$$

$$+ 2\,\mathbb{E}\left[\left(Y - \Phi^\top \beta^*\right)\left(\Phi^\top(\beta^* - \beta)\right)\right]$$

$$\underbrace{\phantom{+ 2\,\mathbb{E}\left[\left(Y - \Phi^\top \beta^*\right)\left(\Phi^\top(\beta^* - \beta)\right)\right]}}_{}$$

<span style="color:blue">$= 0$</span>

<span style="color:blue">Exercise $\left(\mathbb{E}\left[Y - \Phi^\top \beta^* \mid \Phi\right] = 0\right)$</span>

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

∎

Because $\sigma^2$ is "intrinsic noise quantity", we will focus on the excess risk:

$$R(\beta) - \sigma^2 = \|\beta^* - \beta\|_\Sigma^2$$

• So, now assume we observe $(Y_1, \Phi_1), \ldots, (Y_n, \Phi_n)$

②

and compute the OLS $\hat{\beta} = \hat{\Sigma}^{-1} \sum_{i=1}^{n} \frac{y_i \cdot \Phi_n}{n}$

where $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} \Phi_n \Phi_n^T$ which <u>we assume to</u>

<u>be invertible with probability 1.</u>   A sufficient

condition for this is $\underset{d \times d}{\hat{\Sigma}}$ is of full rank

and $n \geq d$

so → the distribution of the $\Phi_n$'s does not concentrate on any affine linear subspace

<u>Prop. 3.10</u>   The expected excess risk

of $\hat{\beta}$ (OLS estimator) is:

$$\frac{\sigma^2}{n} \; \mathbb{E}\left[ tr\left( \Sigma \, \hat{\Sigma}^{-1} \right) \right]$$

PF/   <u>Notation</u>   Let $\hat{\Phi}$ be $n \times d$ matrix with rows

given by $\Phi_1^T, \dots, \Phi_n^T$. This is the random

design matrix. So in particular

$$\hat{\Sigma} = \frac{1}{n} \hat{\Phi}^T \hat{\Phi}$$

Also let $Y$ and $\varepsilon$ be $n$-dimensional vectors

of responses and errors, so

because → the model is linear

$$\hat{\beta} = \hat{\Sigma}^{-1} \hat{\Phi}^T Y = \hat{\Sigma}^{-1} \frac{\hat{\Phi}^T}{n} \left( \hat{\Phi} \beta^* + \varepsilon \right)$$

$$= \beta^* + \frac{\hat{\Sigma}^{-1} \hat{\Phi}^T \varepsilon}{n}$$

③

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

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

$$= \mathbb{E}_{\varepsilon,\hat{\Phi}}\left[\text{tr}\left(\Sigma \;\; \hat{\Sigma}^{-1}\frac{\hat{\Phi}^T \varepsilon}{n}\frac{\varepsilon^T\hat{\Phi}}{n}\hat{\Sigma}^{-1}\right)\right]$$

$$= \mathbb{E}_{\hat{\Phi}}\left[\mathbb{E}_{\varepsilon|\hat{\Phi}}\left[\qquad \text{``}\qquad\right]\right]$$

$$= \mathbb{E}_{\hat{\Phi}}\left[\frac{\text{tr}}{n}\left(\Sigma \;\; \hat{\Sigma}^{-1}\frac{\hat{\Phi}^T}{n}\underbrace{\mathbb{E}_{\varepsilon|\hat{\Phi}}\left[\varepsilon\varepsilon^T\right]}_{\sigma^2 I_n}\hat{\Phi}\,\hat{\Sigma}^{-1}\right)\right]$$

$$= \frac{\sigma^2}{n}\mathbb{E}_{\hat{\Phi}}\left[\text{tr}\left(\Sigma\,\hat{\Sigma}^{-1}\underbrace{\frac{\hat{\Phi}^T\hat{\Phi}}{n}}_{\hat{\Sigma}}\hat{\Sigma}^{-1}\right)\right]$$

$$= \frac{\sigma^2}{n}\mathbb{E}\left[\text{tr}\left(\Sigma\,\hat{\Sigma}^{-1}\right)\right]$$

• What if the model is not well specified? That is what if $\mathbb{E}[Y \mid \Phi] \neq \Phi^T\beta^*$ (the regression function is not linear).

$\downarrow$ $\in \mathbb{R}$ $\downarrow$ $\in \mathbb{R}^d$

Then we can define as our parameter

$$\beta^* = \underset{\beta \in \mathbb{R}^d}{\text{argmin}}\; \mathbb{E}_{Y,\Phi}\left[(Y - \Phi^T\beta)^2\right]$$

④

$$= \underset{\beta \in \mathbb{R}}{\text{argmin}} \quad \mathbb{E}_{\underline{\Phi}} \left[ \left( \mathbb{E}[Y_i | \underline{\Phi}] - \underline{\Phi}^T \beta \right)^2 \right]$$

best linear approximation to the regression function in $L_2$

$$= \underline{\Sigma}^{-1} \; \mathbb{E}\left[ \underline{\Phi} \cdot Y \right]$$

this <u>unique</u> as long as $\underline{\Sigma} = \mathbb{E}\left[ \underline{\Phi} \, \underline{\Phi}^T \right]$ is <u>invertible</u> and $\mathbb{E}[Y^2] < \infty$.
$\qquad \qquad \qquad \qquad \text{dist}$

- $\beta^*$ is sometimes called the <u>projection</u> parameter

- More than one joint distribution of $(Y, \underline{\Phi})$ can have the same projection parameter!

- $\beta^*$ is the vector of coefficients of the $L_2$ projection of $Y$ onto the linear span of $\underline{\Phi}$ (ie the set of r.v.'s of the form $\{\underline{\Phi}^T \beta, \beta \in \mathbb{R}^d\}$)

  ↓ measure of linear association between $Y$ and the vector $\underline{\Phi}$.

- Then, one can show that, in this situation the risk

$$\beta \in \mathbb{R}^d \; \mapsto \; \mathbb{E}_{Y, \underline{\Phi}} \left[ (Y - \underline{\Phi}^T \beta)^2 \right] =$$

$$\underbrace{\mathbb{E}\left[ (Y - \mathbb{E}[Y | \underline{\Phi}])^2 \right]}_{\sigma^2 \text{ intrinsic/unavoidable variance}} + \underbrace{\mathbb{E}\left[ \left( \mathbb{E}[Y | \underline{\Phi}] - \underline{\Phi}^T \beta^* \right)^2 \right]}_{\substack{\text{non- linearity} \\ \text{which is} \\ = 0 \text{ when} \\ \text{model is linear}}} + \underbrace{\| \beta^* - \beta \|^2_{\underline{\Sigma}}}_{\substack{\text{function} \\ \text{form}}}$$

Recall $\mathbb{E}[Y | \underline{\Phi}] = \underset{f}{\text{argmin}} \quad \mathbb{E}\left[ (Y - f(\underline{\Phi}))^2 \right]$

⑤