

M2: Multiple linear regression [2.3, 3, B4] 27.08.2017

Notation:

Independent units of observation

$$i = 1, \dots, n \quad \mathbb{R} (y_i, x_i^T) \quad p\text{-vector}$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \underbrace{\mathbf{Y}}_{n \times 1} \quad \begin{array}{l} \leftarrow \text{random variable} \\ \text{vector of responses} \quad \leftarrow \text{rent} \\ \text{wzges} \end{array}$$

$$\begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = \underbrace{\mathbf{X}}_{n \times p} \quad \begin{array}{l} \leftarrow \text{fixed} \\ \text{design matrix} \\ \text{rows = observations} \\ \text{cols = variables} \end{array}$$

parameters of interest

β (px1) vector of regression coeff (parameters)
 ↗ unknown
 ↗ main focus to find β estimate

ε : (n x 1) random vector
 ↗ unknown & unobserved

Model:

$$Y = \mathbf{X}\beta + \varepsilon \quad \varepsilon \sim N_n(0, \sigma^2 I)$$

$$y_i = x_i^T \beta + \varepsilon_i \quad \begin{array}{l} \varepsilon_i \sim N(0, \sigma^2) \\ \varepsilon_i, \varepsilon_j \text{ independent} \end{array}$$

1

Other model assumptions:

Full rank: $\text{rank}(\mathbf{X}) = p$

$\mathbf{X}_{n \times p}$ where $n \gg p$

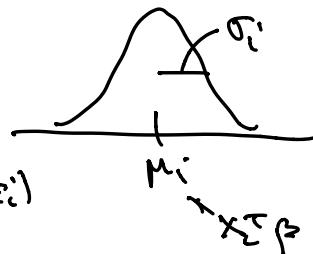
number of lin. indep. columns.

Identifiability problems when not full rank.

Distribution of Y_i

$$Y_i = \underbrace{\mathbf{x}_i^T \boldsymbol{\beta}}_{\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi}} + \varepsilon_i \stackrel{N(0, \sigma^2)}{\sim}$$

$$Y_i \sim N(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2)$$



$$E(Y_i) = E(\mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i) = \mathbf{x}_i^T \boldsymbol{\beta} + E(\varepsilon_i)$$

$$\text{Var}(Y_i) = \text{Var}(\mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i) = 0 + \sigma^2$$

Y_i, Y_j will be independent since ε_i and ε_j are indep.

$$\mathbf{Y}_{n \times 1} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

but, we could have written

$$Y_i \sim N(\mu_i, \sigma^2)$$

↑
par. of interest

↑ nuisance
distribution of response

$$\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta} \quad \leftarrow \text{linear predictor}$$

(in pred.)

$$\mu_i = \eta_i \quad \leftarrow \text{connection } \xrightarrow{\text{mean}} \mu_i \text{ and } \eta_i$$