

# Predictive Modeling

**Review** The best function of  $X$  to predict  $Y$  in terms of the **MSE** is  $f(X) = \mathbb{E}(Y|X)$ .

## Strategies

- (a) Model the **joint distribution** of  $(X, Y) \Rightarrow \mathbb{E}(Y|X)$ .
- (b) Model the **conditional distribution** of  $Y|X$  directly, often

$$Y|X \sim N(\mathbb{E}(Y|X), \sigma^2),$$

so the “only component left” to model is  $\mathbb{E}(Y|X)$  itself.

**Notation** Will begin to use **lower-case** letters for RVs as well; need **upper-case** letters for **matrices**, etc.

# Natural Starting Point

For  $y \in \mathbb{R}$  and  $\mathbf{x} \in \mathbb{R}^p$ , model

$$\mathbb{E}(y|\mathbf{x}) = \alpha + \boldsymbol{\beta}^\top \mathbf{x}$$

as a linear function of  $\mathbf{x}$ .

**Fact [Example 3.5 (p. 37)]** If  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{p+q}$  is modelled jointly as a (multivariate) normal distribution with

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix},$$

then it can be derived that

$$\mathbb{E}(\mathbf{y}|\mathbf{x}) = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x)$$

is linear in  $\mathbf{x}$ , and  $\text{Var}(\mathbf{y}|\mathbf{x}) = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$  “regardless of”  $\mathbf{x}$ .

# Random Vectors

For  $\mathbf{z} \in \mathbb{R}^m$ ,

$$\mathbb{E}(\mathbf{z}) = \begin{bmatrix} \mathbb{E}(z_1) \\ \mathbb{E}(z_2) \\ \vdots \\ \mathbb{E}(z_m) \end{bmatrix}_{m \times 1}$$

and

$$\text{Var}(\mathbf{z}) = \begin{bmatrix} \text{Var}(z_1) & \text{Cov}(z_1, z_2) & \cdots & \text{Cov}(z_1, z_m) \\ \text{Cov}(z_2, z_1) & \text{Var}(z_2) & \cdots & \text{Cov}(z_2, z_m) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(z_m, z_1) & \text{Cov}(z_m, z_2) & \cdots & \text{Var}(z_m) \end{bmatrix}_{m \times m} .$$

# Linear Regression

$$y_i = \overset{\beta_0}{\downarrow} \alpha + \underset{\mathbb{R}^{p-1}}{\uparrow} \boldsymbol{\beta}^\top \mathbf{x}_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, 2, \dots, n$$

$$\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\substack{n \times 1 \\ \downarrow \\ \mathbf{y}}} = \underbrace{\begin{bmatrix} 1 & \mathbf{x}_1^\top \\ 1 & \mathbf{x}_2^\top \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^\top \end{bmatrix}}_{\substack{n \times p \\ \downarrow \\ \mathbf{X}}} \underbrace{\begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix}}_{\substack{p \times 1 \\ \downarrow \\ \boldsymbol{\beta}}} + \underbrace{\begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}}_{\substack{n \times 1 \\ \downarrow \\ \boldsymbol{\varepsilon}}}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\Rightarrow \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

**Remarks** (i) Notice that  $\boldsymbol{\beta}^\top \mathbf{x}_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ . (ii) Often write compactly where each  $\mathbf{x}_i$  implicitly contains a “1”.

## Two Initial Questions

(a) How to estimate the (vector) parameter  $\beta$ ?

– Answer: least squares; maximum likelihood.

(b) How good is the estimator  $\hat{\beta}$ ?

– Answer: mean-squared error; bias-variance analysis.

# Least Squares (LSQ)

Calculus

$$\begin{aligned} \min_{\boldsymbol{\beta}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\ \Rightarrow \quad & \frac{d}{d\boldsymbol{\beta}}(\dots) = 0 \\ \Rightarrow \quad & \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0 \end{aligned}$$

Geometry

$$\begin{aligned} \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} & \perp \mathbf{X} \\ \Rightarrow \quad & \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0 \end{aligned}$$

$$\begin{aligned} \mathbf{X}^\top \mathbf{y} &= \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

# Bias-Variance Analysis

## Tools

$$\mathbb{E}(aX + b) = a\mathbb{E}(X) + b \quad \Rightarrow \quad \mathbb{E}(\mathbf{A}\mathbf{y} + \mathbf{b}) = \mathbf{A}\mathbb{E}(\mathbf{y}) + \mathbf{b}$$

$$\text{Var}(aX + b) = a^2\text{Var}(X) \quad \Rightarrow \quad \text{Var}(\mathbf{A}\mathbf{y} + \mathbf{b}) = \mathbf{A}\text{Var}(\mathbf{y})\mathbf{A}^\top$$

## Analysis

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \text{N}(\mathbf{0}, \sigma^2\mathbf{I})$$

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

$$\mathbb{E}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \underbrace{\mathbb{E}(\mathbf{y})}_{\mathbf{X}\boldsymbol{\beta}} = \boldsymbol{\beta}$$

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{A}} \underbrace{\text{Var}(\mathbf{y})}_{\sigma^2 \mathbf{I}} \underbrace{\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}}_{\mathbf{A}^\top} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

# Gauss-Markov Theorem

**Theorem** Suppose  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where

$$\mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0} \quad \text{and} \quad \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}.$$

Let

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

be the least-squares estimator, and let  $\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{y}$  be another linear and unbiased estimator of  $\boldsymbol{\beta}$  s.t.  $\mathbb{E}(\mathbf{A}\mathbf{y}) = \boldsymbol{\beta}$ . Then

$$\text{Var}(\boldsymbol{\ell}^\top \hat{\boldsymbol{\beta}}) \leq \text{Var}(\boldsymbol{\ell}^\top \tilde{\boldsymbol{\beta}})$$

for all  $\boldsymbol{\ell} \in \mathbb{R}^p$ .

**Translation** LSQ estimator is best (in terms of **MSE**) among all linear and unbiased estimators. [**Will prove but later.**]

# A Look at Assumptions

Assumptions	LSQ	B-V Analysis
$\mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}$	×	✓
$\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$	×	✓
$\boldsymbol{\varepsilon} \sim \text{N}(\cdot, \cdot)$	×	×

**Remark** Don't automatically trust the software! All regression packages/libraries implement these assumptions and they can be wrong in practice.

**Exercise** Consider the trivial model  $y_i = \beta + \varepsilon_i$ , where  $\mathbb{E}(\varepsilon_i) = 0$  and  $\text{Var}(\varepsilon_i) = \sigma^2$  for all  $i$  but  $\text{Corr}(\varepsilon_i, \varepsilon_j) = \rho \neq 0$ . Find  $\text{Var}(\hat{\beta})$  and compare with what the software would produce under the standard assumptions. [For example, try  $\sigma^2 = 1$ ,  $n = 1000$  and  $\rho = 0.1$ .]

# LSQ and MLE

Note

$$\|z\|^2 = \sum_i z_i^2 \quad \Rightarrow \quad \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2,$$

so LSQ  $\Leftrightarrow$  MLE under model  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$ ,  $\varepsilon_i \stackrel{iid}{\sim} \mathbf{N}(0, \sigma^2)$ :

$$\begin{aligned} \max_{\boldsymbol{\beta}} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{2\sigma^2}\right] &= \\ \left[\frac{1}{\sqrt{2\pi\sigma^2}}\right]^n \exp\left[-\frac{\sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{2\sigma^2}\right] & \\ \Leftrightarrow \min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2. & \end{aligned}$$

But the LSQ principle itself does NOT depend on the modeling assumption  $\varepsilon_i \stackrel{iid}{\sim} \mathbf{N}(0, \sigma^2)$ .

# Simple Linear Regression

**Exercise** Consider the special case of  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ .

(a) Show that

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}.$$

(b) Show that

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}.$$

(c) Show that  $\text{Cov}(\bar{y}, \hat{\beta}_1) = 0$ , and that

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum (x_i - \bar{x})^2}, \quad \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\sigma^2 \bar{x}}{\sum (x_i - \bar{x})^2}.$$

**Remark** These formulae contain many interesting **insights**.

(a)

- the fitted regression line  $\hat{y} = \hat{\beta}_0 + \hat{\beta}x$  goes through  $(\bar{x}, \bar{y})$
- if data are **centered** in the sense that

$$x_i \leftarrow x_i - \bar{x} \quad \text{and} \quad y_i \leftarrow y_i - \bar{y},$$

then  $\bar{x} = \bar{y} = 0$  and no intercept term  $\beta_0$  is needed

- if data are not only **centered** but also **scaled** in the sense that

$$\sum x_i^2 = 1 \quad \text{and} \quad \sum y_i^2 = 1,$$

then

$$|\hat{\beta}_1| = \left| \sum x_i y_i \right| \leq \sqrt{\sum x_i^2 \sum y_i^2} = 1 \quad (\text{why}),$$

aka “**regression to the mean**” [Fun Box 1, p. 74]

(b)

**Exercise** Two researchers are investigating the effect of **age** on **blood pressure** with simple linear regression. Amy collected her data by setting up a desk in a supermarket, and her subjects' ages are

$\{21, 27, 32, 35, 39, 50, 58, 65, 73, 79\}$ .

Bob collected his data set by visiting a university classroom, and his subjects' ages are

$\{19, 19, 20, 20, 20, 20, 21, 21, 21, 22\}$ .

Whose estimate of the “age effect on blood pressure” will be more accurate? Why?

(c)

**Prediction** Given a new observation  $x_0$ , predict  $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ .  
How accurate is the prediction?

$$\begin{aligned}\text{Var}(\hat{y}_0) &= \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0) \\ &= \text{Var}(\hat{\beta}_0) + x_0^2 \text{Var}(\hat{\beta}_1) + 2x_0 \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum (x_i - \bar{x})^2} + \frac{\sigma^2 x_0^2}{\sum (x_i - \bar{x})^2} - \frac{2x_0 \bar{x} \sigma^2}{\sum (x_i - \bar{x})^2} \\ &= \sigma^2 \left[ \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]\end{aligned}$$

**Lesson** Prediction most accurate when  $x_0 \approx \bar{x}$ , and progressively less accurate as  $x_0$  moves away from center of **training data**.

# Important Take-Home Messages

- You can always use your model — any model — to make a prediction, but you shouldn't always trust it.
- Not all predictions made by the same model are equal; some are more reliable than others.
- Extrapolation is dangerous in general.

↓  
applicable **in general**; NOT just for simple linear regression