

Today

- 1) 3 exercises
 - 2) Gauss-Markov theorem
orthogonalization
 - 3) Model selection
 - 4) Shrinkage methods (Ch 3.4)
-

Ex 2.7

N pairs (x_i, y_i) ^{drawn} iid from

$$x_i \sim h(x)$$

$$y_i = f(x_i) + \varepsilon_i$$

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

Estimator for f linear in x :

$$\hat{f}(x_0) = \sum_{i=1}^N \ell_i(x_0; X) y_i$$

where weights $\ell_i(x_0; X) \perp y_i$, but depend on X , the entire training sequence of x_i

$$a) \hat{f}_{LR}(x_0) = X_0^T \hat{\beta} = x_0^T (X^T X)^{-1} X^T y$$

$$\Rightarrow \ell_i(x_0, X) = [x_0^T (X^T X)^{-1} X^T]_{i, \cdot}$$

$$\hat{f}_{kNN}(x_0) = \text{Ave} \{ y_i \mid x_i \in N_k(x_0) \}$$

where $N_k(x_0)$ is the set of the k nearest neighbors

$$\Rightarrow \ell_i(x_0, X) = \begin{cases} 1/k & \text{if } x_i \in N_k(x_0) \\ 0 & \text{otherwise} \end{cases}$$

b) Decompose the conditional mean-squared error

$$E_{Y|x} [(f(x_0) - \hat{f}(x_0))^2]$$

statistic fixed

$$\begin{aligned} & E_{Y|x} \left[\left(\underbrace{f(x_0) - E_{Y|x}[\hat{f}(x_0)]}_{\text{bias}} + \underbrace{E_{Y|x}[\hat{f}(x_0)] - \hat{f}(x_0)}_{\text{variance}} \right)^2 \right] \\ &= E_{Y|x} \left[\left(f(x_0) - E_{Y|x}[\hat{f}(x_0)] \right)^2 + \left(E_{Y|x}[\hat{f}(x_0)] - \hat{f}(x_0) \right)^2 \right. \\ &\quad \left. + 2 \left(f(x_0) - E_{Y|x}[\hat{f}(x_0)] \right) \left(E_{Y|x}[\hat{f}(x_0)] - \hat{f}(x_0) \right) \right] \\ &= \underbrace{\left(f(x_0) - E_{Y|x}[\hat{f}(x_0)] \right)^2}_{\text{bias}^2} + \underbrace{\text{Var}_{Y|x}(\hat{f}(x_0))}_{\text{variance}} + \\ &\quad + 2 \left(f(x_0) - E_{Y|x}[\hat{f}(x_0)] \right) E_{Y|x} \left[E_{Y|x}[\hat{f}(x_0)] - \hat{f}(x_0) \right] \\ &= \text{bias}^2 + \text{variance} \end{aligned}$$

c) unconditional

$$E_{Y,X} [(f(x_0) - \hat{f}(x_0))^2]$$

$$E_X \left[E_{Y|x} [(f(x_0) - \hat{f}(x_0))^2] \right]$$

d) conditional:

$$\left(f(x_0) - E_{Y|x}[\hat{f}(x_0)] \right)^2 + \text{Var}_{Y|x}(\hat{f}(x_0))$$

$$\left(f(x_0) - E_{Y|x} \left[\sum_{i=1}^n \ell_i(x_0; \mathcal{X}) y_i \right] \right)^2 + \text{Var}_{Y|x} \left[\sum_{i=1}^n \ell_i(x_0; \mathcal{X}) y_i \right]$$

$$\left(f(x_0) - \sum_{i=1}^n \ell_i(x_0; \mathcal{X}) f(x_0) \right)^2 + \sum_{i=1}^n \ell_i(x_0; \mathcal{X})^2 \sigma^2$$

unconditional

$$E_X \left[\left(f(x_0) - E_{Y|x} \left[\sum_{i=1}^n \ell_i(x_0; \mathcal{X}) y_i \right] \right)^2 \right] + E_X \left[\text{Var}_{Y|x} \left[\sum_{i=1}^n \ell_i(x_0; \mathcal{X}) y_i \right] \right]$$

Ex 3.1

$$F = \frac{(RSS_0 - RSS_1) / (p_1 - p_0)}{RSS_1 / (N - p_1 - 1)}$$

when we are testing only one parameter $F \approx z^2$

$$\rightarrow p_1 = p, \quad p_0 = p - 1$$

$$F = \frac{(RSS_0 - RSS_1)}{RSS_1 / (N - p - 1)} \sim F_{1, N-p-1}$$

$$z = \frac{\hat{\beta}_j}{\hat{\sigma} [X'X]_{jj}^{-1/2}} \sim t_{N-p-1}$$

$$(t_{N-p-1})^2 \stackrel{d}{=} F_{1, N-p-1} \Rightarrow z^2 \approx F$$

Gauss - Markov theorem

$\hat{\beta}_{LS}$ is BLUE

Best \leftarrow smallest variance

Linear $\leftarrow \hat{\theta} = a^T \hat{\beta} \leftrightarrow \theta = a^T \beta$

Unbiased $\leftarrow E[\hat{\theta}] = \theta$

Estimator

consider X fixed

$$\begin{aligned} E[a^T \hat{\beta}] &= E[a^T (X^T X)^{-1} X^T y] \\ &= a^T (X^T X)^{-1} X^T X \beta \\ &= a^T \beta \end{aligned}$$

$$Th: \text{Var}(a^T \hat{\beta}) \leq \text{Var}(c^T \hat{\beta}) \rightarrow \text{ex 3.3 (a)}$$

$$\begin{aligned} MSE(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta)^2] \\ &= \text{Var}(\hat{\theta}) + (E[\hat{\theta}] - \theta)^2 \end{aligned}$$

smallest with LS \downarrow \parallel
0

• $MSE(\hat{\theta})$ is strongly related with the prediction error

Let consider $Y_0 = f(x_0) + \varepsilon$ $\varepsilon \sim N(0, \sigma^2)$

$$E[(Y_0 - \hat{f}(x_0))^2] = E[Y_0^2 + \hat{f}(x_0)^2 - 2Y_0 \hat{f}(x_0)]$$

$$= E[Y_0^2] + E[\hat{f}(x_0)^2] - 2E[Y_0 \hat{f}(x_0)]$$

$$= Var(Y_0) + \cancel{E[Y_0]^2} + Var(\hat{f}(x_0)) + E[\hat{f}(x_0)]^2 - 2f(x_0)E[\hat{f}(x_0)]$$

$$= \sigma^2 + Var(\hat{f}(x_0)) + E[\hat{f}(x_0)]^2 - 2f(x_0)E[\hat{f}(x_0)] + f(x_0)^2$$

$$= \sigma^2 + Var(\hat{f}(x_0)) + \text{bias}^2$$

MSE

$$E[Y_0] = E[f(x_0) + \varepsilon] = E[f(x_0)] + E[\varepsilon]$$

$$\downarrow$$

$$= f(x_0) + 0$$

$$E[Y_0^2] = E[(f(x_0) + \varepsilon)^2]$$

$$\downarrow$$

$$= Var(Y_0) + E[Y_0]^2$$

$$\downarrow$$

$$= \sigma^2 + f(x_0)^2$$

Suppose $Y = X\beta + \varepsilon$

univariable

if $Y = X_1\beta_1 + \varepsilon$

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{\langle x, y \rangle}{\langle x, x \rangle}$$

multivariable

if $Y = X_1\beta_1 + X_2\beta_2 + \varepsilon$

only if X is orthogonal, $\hat{\beta}_j = \frac{\langle x_j, y \rangle}{\langle x_j, x_j \rangle}$

alternative formula for variance of $\hat{\beta}_j$

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\langle e_j, e_j \rangle} = \frac{\sigma^2}{\|e_j\|^2}$$

→ variables may share the same information useful for explaining/predicting y

→ the estimates may be unstable due to high variance

↓
model selection

$$z_{\text{score}} = \frac{\hat{\beta}}{\text{sd}(\hat{\beta})} \quad \begin{array}{l} \rightarrow \text{smaller} \\ \rightarrow \text{larger} \end{array}$$

Model Selection

- prediction accuracy
- interpretability
- portability

• best subset technique

$$x = (x_1, x_2, x_3)$$

models : 0 variables

$$y = \beta_0 + \varepsilon$$

1 variable

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

$$y = \beta_0 + \beta_2 x_2 + \varepsilon$$

$$y = \beta_0 + \beta_3 x_3 + \varepsilon$$

2 variables

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

$$y = \beta_0 + \beta_1 x_1 + \beta_3 x_3 + \varepsilon$$

$$y = \beta_0 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

$$\binom{3}{2} = \frac{3 \cdot 2}{2} = 3$$

$$\rightarrow \binom{6}{2} = \frac{6 \cdot 5}{2} = 15 \quad \text{picture}$$

• Stepwise techniques

- forward selection

start : $y = \beta_0 + \varepsilon$

$$x_1, x_2, x_3$$

1st step : $y = \beta_0 + \beta_2 x_2 + \varepsilon$

$$x_1, x_3$$

2nd step : $y = \beta_0 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$

β_3 together with β_0 and β_2

- backward elimination

start: full model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

$p \gg n$ case
backward elimination is
not possible

1st step $y = \beta_0 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$

⋮

- stepwise selection
step back "

$$RSS(\theta) + J(\theta) + \underline{\underline{2|p|}}$$

stagewise regression

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

stepwise: the estimate
of β take into account
all x

stagewise: introduce
a new β , its estimate
is only based on x_j