

Week 10.

$$Y \sim X_1$$

$$Y \sim X_1 + X_2 + \dots + X_p$$

goal: ① Find linear relationship btw  $(Y, X)$   
 ② prediction  $\hat{Y} = \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p$

$Y \sim X_1$  (simple)  
 $Y \sim X_1 + X_2 + \dots + X_p$  (multiple)

• Model significance

$H_0: \beta_1 = 0$  vs.  $H_a: \beta_1 \neq 0$   
 $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$  vs.  $H_a: \text{at least one } \beta \neq 0$   
 model is not useful

• Individual term significance

$H_0: \beta_i = 0$  vs.  $H_a: \beta_i \neq 0$   
 $H_0: \beta_1 = 0$  vs.  $H_a: \beta_1 \neq 0$   
 $H_0: \beta_2 = 0$  vs.  $H_a: \beta_2 \neq 0$   
 $\vdots$   
 $H_0: \beta_p = 0$  vs.  $H_a: \beta_p \neq 0$

• Estimated regression line

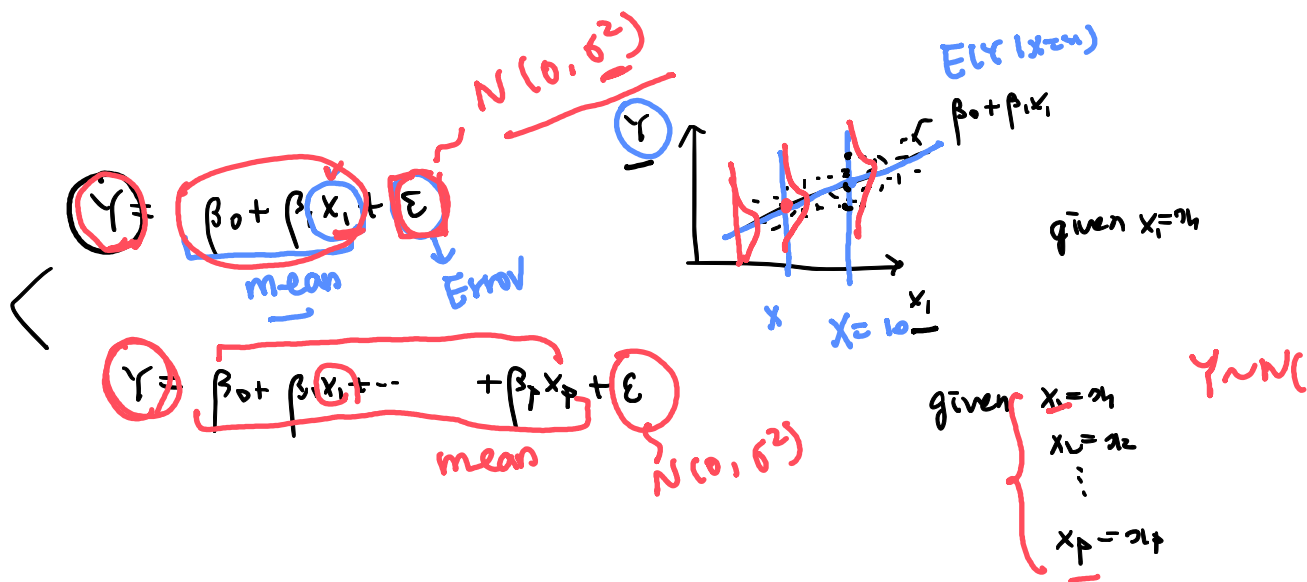
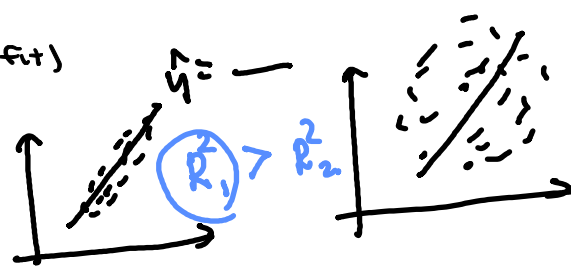
$$\hat{Y} = -6041 + 8.4 \cdot M + 35.2 \cdot S_0 + \dots$$

•  $R^2$

prediction power  
 $= \left( \frac{SS_{\text{model}}}{SS_{\text{total}}} \right)$  goodness-of-fit

• Model diagnostics

model assumption  
 influential cook's D.



## Multicollinearity

⇒ Remove highly cor variables

$y$

$(x_1 \dots x_{10})$

$n = 100$

$x_1$   
 $\vdots$   
 $x_{10}$

## Model selection

"optimal" or "best"

## Multicollinearity

(eg) US crime data

About predictors  
not model assumptions  
or  $Y$

Simulation studies

cirrhosis  
alcohol  
alcohol2 = 3 \* alcohol  
alcohol3 = 3 \* alcohol + 2

① cirrhosis ~ alcohol

② cirrhosis ~ alcohol + alcohol2

$cor(alcohol, alcohol2) = 1$

③ cirrhosis ~ alcohol + alcohol3

$cor(alcohol, alcohol3) = 0.99$

{ F-test : model is useful  
t-test

How to detect? what is the remedy?

Visualization : pairwise scatter plot  
 Quant. measure: VIF

$vif(i) = \frac{R_i^2}{R_j^2}$

$x_1$	$x_2$	$x_3$	$x_4$	...	$x_9$	$x_{10}$
3.6	2.7	4.5	2.3	...	9.5	9.7

$R_1^2$  (large)  $R_2^2$  (small)

$lm(x_1) \sim (x_2 + \dots + x_{10})$   
 $lm(x_2) \sim x_1 + x_3 + \dots + x_{10}$

$VIF \approx 1$

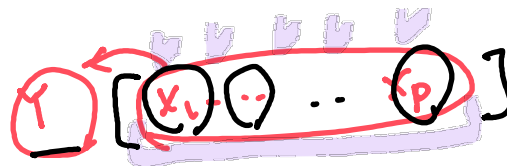
$x_j$ : high var w/ others  
 $x_j$ : small

⇒ Remove one by one

① Remove predictor w/ highest VIF (w/ threshold)

②  $Y \sim \dots - P02$

③



# Model Selection

$10 [X_1 \dots X_{10}] Y \sim \epsilon_1 \dots \epsilon_{10}$   
 $\epsilon_i \sim N(0, \sigma^2)$   
 $Y = X\beta + \epsilon$

Automatic Selection: backward / forward / stepwise

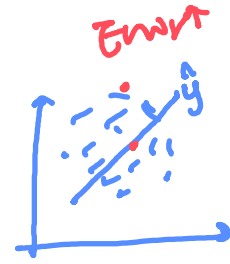
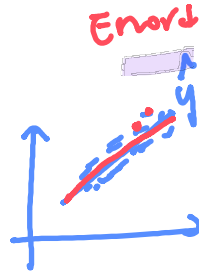
criteria  
 $\Rightarrow$  p-value  
 $\Rightarrow$  BIC

Best subset approach: goodness-of-fit measure criteria

{ adjusted  $R^2$  / AIC / BIC / Mallows  $C_p$  }

(90)  
 $p=100, R^2=88\%$   
 $p=10, R^2=85\%$   
 $p=2, R^2=50\%$

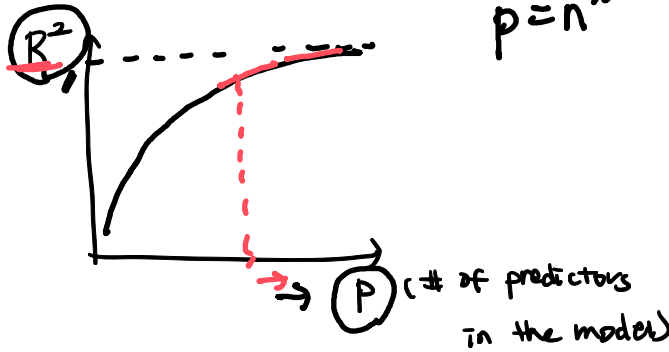
good prediction  $\rightarrow$  simplicity



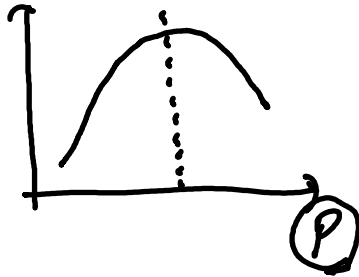
penalized goodness-of-fit  $\sum (y_i - \hat{y}_i)^2$   
 small

$\square$  adj- $R^2$

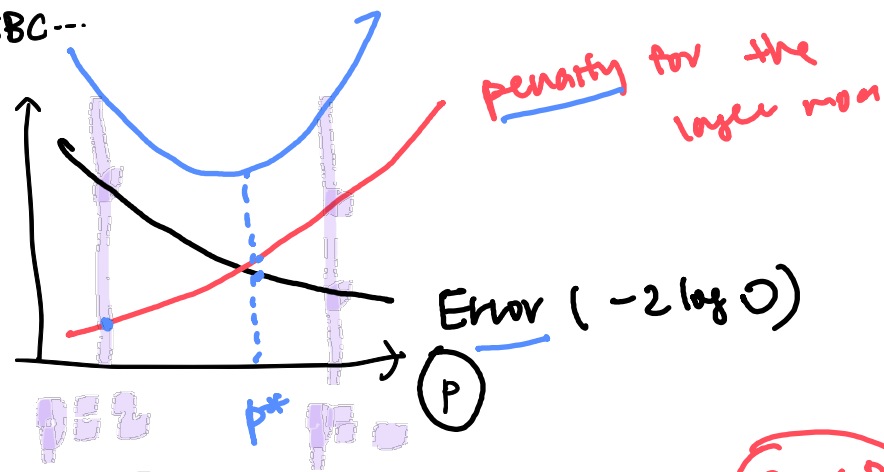
100 "p=n" 100  
 $R^2=1$



adj- $R^2$



④ AIC, BIC, SBC...



$$= \underbrace{-2 \log L(M)}_{\sum (y_i^{\wedge} - y_i)^2} + \text{penalty}$$

model 1	$C_p: 5$	$p=3$
model 2	$C_p: 5.1$	$p=6$

$C_p \leq p$

Data Example

$\text{adj-}R^2$	model 8	M Ed Pol M.F (U1) U2 Ineq Prob
AIC	model 8	
BIC	model 6	M Ed Pol — — U2 Ineq Prob.
$C_p$	model 6	

stepwise

# Multicollinearity Issue in model selection

- step
- ① Automatic selection (p-value criteria)
  - ② goodness-of-fit measure search

