

$$\text{CI for } \hat{Y}_h : \hat{Y}_h \pm B \cdot \text{sd(pred)}$$

$f$ : # of observations.

$$B = t(1 - \frac{\alpha}{2g}; n-2)$$

$$\text{Scheffé procedure: } \hat{Y}_h \pm S \cdot \text{sd(pred)}$$

$$S = g \cdot F(1 - \alpha; g; n-2) \quad \text{without derivation.}$$

similar to W-H.

## Fwd for Matrix #1

---

\* (Chap 5) Matrix approach to LR.

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \quad \begin{matrix} n \times n \\ \downarrow \quad \downarrow \\ \text{row} \quad \text{col.} \end{matrix}$$

vectors:  $(n \times 1)$  matrices.

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$A^T \rightarrow$  transpose: Swap rows & cols.

$$\text{Matrix } C, EC(C) = [EC(c_{ij})], \quad y, EY = \begin{bmatrix} EY_1 \\ \vdots \\ EY_n \end{bmatrix}$$

i) dot product  $y \cdot y$  ( $\langle y, y \rangle$ )  
(scalar)

$$y \cdot y = y_1^2 + y_2^2 + \dots + y_n^2 = \|y\|^2$$

$\|y\|$ : norm / length / magnitude of  $y$ .

$$u, v \in \mathbb{R}^n, u \cdot v = u_1 v_1 + u_2 v_2 + \dots + u_n v_n.$$

•) Matrix vector product:

$$\underset{m \times n}{A} \underset{n \times 1}{u} = \begin{bmatrix} \text{Row}_1(A) \cdot u \\ \vdots \\ \text{Row}_m(A) \cdot u \end{bmatrix}$$

$$\text{e.g. } \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} u_1 + 2u_2 \\ 3u_1 + 4u_2 \end{bmatrix}$$

Matrix-matrix product:

$$\underset{m \times n}{A} \cdot \underset{n \times p}{B} = \begin{bmatrix} A \cdot B_1 & \cdots & A \cdot B_p \end{bmatrix} \quad B_i = \text{the } i\text{th col of } B.$$

Write dot product in matrix form:

$$\begin{aligned} \underline{u} \cdot \underline{v} &= \underline{u}^T \underline{v} = [u_1 \cdots u_n] \begin{bmatrix} v^1 \\ \vdots \\ v_n \end{bmatrix} \\ &= \underline{v}^T \underline{u} = [v_1 \cdots v_n] \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \Rightarrow \text{scalar.} \end{aligned}$$

$$\underline{y} \cdot \underline{y} = \|\underline{y}\|^2 = \underline{y}^T \underline{y} = \sum_i y_i^2$$

$\underline{y}^T \underline{y} \Rightarrow (n \times n)$  matrix. Outer Product.  $\neq$  dot product.

$$\underset{n \times 2}{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \underset{n \times n \ n \times n}{X^T X} = \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}$$

$$(X^T X)^T = X^T X \rightarrow \underline{A}^T = \underline{A} \Rightarrow A \text{ is symmetric.}$$

Fact: For any  $A$ ,  $A^T A$  is always a symmetric matrix.

$$I_n = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix}$$

o) Inverse of a matrix:  $A$  is invertible if  $\exists A^{-1}$  st.  $A \cdot A^{-1} = A^{-1} \cdot A = I$ .

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad A^{-1} = \frac{1}{\det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}.$$

.) Properties

$$\begin{aligned} (AB)^T &= B^T A^T \\ (A^T)^T &= A \\ (ATA)^T &= A^T A \end{aligned}$$

$$(AB)^{-1} = B^{-1} A^{-1}$$

$$(ABC)^{-1} = C^{-1} B^{-1} A^{-1}$$

$$(A^T)^{-1} = (A^{-1})^T$$

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{det}(x^T x) = n \sum x_i^2 - (\sum x_i)^2 = n \sum (x_i - \bar{x})^2$$

$$= n \sum_{i=1}^n x_i^2 - n \bar{x}^2$$

o) Back to stats.

$$\text{Var}(y) = \mathbb{E}[(y - \mathbb{E}y)^2]$$

$$\mathbb{E}y = \begin{bmatrix} \mathbb{E}y_1 \\ \vdots \\ \mathbb{E}y_n \end{bmatrix} \quad \text{Var } y = \underbrace{\mathbb{E}((y - \mathbb{E}y)(y - \mathbb{E}y)^T)}_{n \times n}.$$

Variance-covariance matrix.

$$\text{Var}(y) = \begin{bmatrix} \text{Var}(y_1) & \text{cov}(y_1, y_2) & \dots & \text{cov}(y_1, y_n) \\ \text{cov}(y_2, y_1) & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ \text{cov}(y_n, y_1) & \dots & \text{Var}(y_n) \end{bmatrix} \quad \text{Covariance matrix}$$

$$\text{Recall } w = a^T y \quad E(w) = a^T E(y), \quad \text{Var}(w) = a^T \text{Var}(y) a.$$

$$W = Ay \quad E(W) = A E(y) \neq E(y) A$$

$m \times 1 \quad m \times n \quad n \times 1$

$$\text{Var}(W) = A \text{Var}(y) A^T$$

$m \times n \quad n \times n \quad n \times m$

i) SLR Model Matrix form.

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad i=1, \dots, n. \rightarrow n \text{ eqns.}$$

$$Y_1 = \beta_0 + \beta_1 X_1 + \epsilon_1$$

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$Y_2 = \beta_0 + \beta_1 X_2 + \epsilon_2$$

$$Y_n = \beta_0 + \beta_1 X_n + \epsilon_n.$$

$$Y = X\beta + \epsilon$$

$$n \times 1 \quad n \times 2 \quad 2 \times 1 \quad n \times 1.$$

$$E(\epsilon_i) = 0.$$

$$\text{Var}(\epsilon_i) = \sigma^2$$

$$\text{cov}(\epsilon_i, \epsilon_j) = 0.$$

$$\Rightarrow E\epsilon = \vec{0}.$$

$$\text{Var}\epsilon = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \sigma^2 \end{bmatrix} = \sigma^2 I_n.$$

$$Y \sim N(\mu, \Sigma)$$

$$P_{X^1} \downarrow \text{mean vector} \rightarrow \text{covariance matrix.}$$

$$\epsilon \sim N(0, \sigma^2 I_n).$$

$$\text{PDF: } f(y) = \frac{1}{(2\pi)^{\frac{n}{2}} (\det \Sigma)^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu) \right]$$

## \* Section #7 (2023 Fall MATH223)

- Gaussian-Markov Theorem.

Normality

In-covariance (independent)

Mean-zero.

$$e_i = y_i - \hat{y}_i$$

$$e_i^* = \frac{y_i - \hat{y}_i}{\sqrt{MSE}} \quad \text{semi-studentized residual.}$$

- Fitties with the model.

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

- Possible plots:

$e_i$  vs.  $X_i$

$e_i^2$  or  $|e_i|$  vs.  $X_i$

$e_i$  vs.  $\hat{y}_i$  Covariance should be zero.

$e_i$  vs. time. Test dependent error term

$e_i$  vs. other predictors.

Normal plot of  $e_i$ 's (Q-Q plot)

1. Non-linear function

e.g.  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \propto X_i^2$   
Heteroskedasticity  
Bf test  $\rightarrow$  2. non-constant variance

RMS test  $\rightarrow$  3. not independent error term

4. outliers.

5. Error terms are not normally distributed.

6. Missing prediction.

QQ plot.

Shapiro-Wilk normality test

# Lecture 10 2023 10/7/24.

Matrix approach to LR

- SLR and Inference
- Multiple LR

SLR w/ normal error in matrix form  
 $\underline{y} = \underline{X} \cdot \underline{\beta} + \underline{\epsilon}$ ,  $\epsilon \sim N(0, \sigma^2 \mathbf{I}_n)$   
 $n \times n \quad n \times 2 \quad n \times 1 \quad n \times 1$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Least Squares Estimation of  $\beta$ .

$$Q(\beta) = Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 = \| \underline{y} - \underline{X} \underline{\beta} \|^2 = (\underline{y} - \underline{X} \underline{\beta})^T (\underline{y} - \underline{X} \underline{\beta})$$

$$\underline{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \underset{\beta_0, \beta_1}{\operatorname{arg\min}} Q(\beta_0, \beta_1) = \underset{\underline{\beta}}{\operatorname{arg\min}} Q(\underline{\beta}).$$

fitted values:  $\hat{y} = \underline{X} \cdot \underline{\beta}$       Soln.:  $\underline{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$   
 $\hat{y} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$       Hat matrix,  $H = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T$

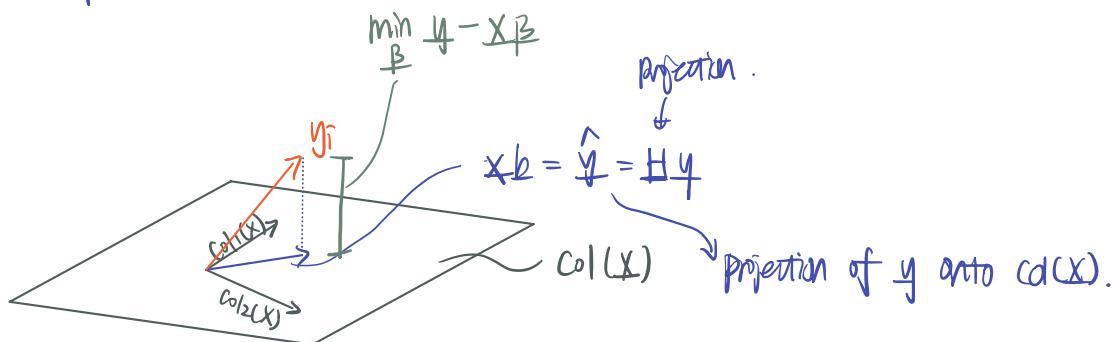
Note: all params, fitted values, residuals, are linear combi of  $y_i$ ,  $i=1, 2, \dots, n$ .

$$\begin{aligned} \underline{\beta} &= A_1 \underline{y} \\ \hat{y} &= A_2 \underline{y}, \quad A_2 = H \\ \underline{\epsilon} &= A_3 \underline{y} \end{aligned} \quad ) \quad \text{like before.} \quad SSTO = SSE + SSR.$$

Normal equations:

$$\underbrace{(\underline{X}^T \underline{X})}_{n \times n} \underline{\beta} = \underline{X}^T \underline{y} \quad \underline{A} \underline{V} = V_1 \underline{A}_1 + V_2 \underline{A}_2 + \dots + V_n \underline{A}_n$$

$$\underline{X} \underline{\beta} = \beta_0 \operatorname{Col}(1, \underline{X}) + \beta_1 \operatorname{Col}(2, \underline{X})$$



Properties:

$$H^T = H \quad (\text{symmetric})$$

$$H^2 = H H = H \quad (\text{idempotent}) \quad \text{pf: } \underline{X} (\underline{X}^T \underline{X})^{-1} (\underline{X}^T \underline{X}) (\underline{X}^T \underline{X})^{-1} \underline{X}^T = H$$

$$E\underline{y} = E[\underline{X}\underline{\beta} + \underline{\epsilon}] = \underline{X}\underline{\beta} + E[\underline{\epsilon}] = \underline{X}\underline{\beta}$$

$$\operatorname{Var} \underline{y} = \operatorname{Var} (\underline{X}\underline{\beta} + \underline{\epsilon}) = 0 + \operatorname{Var} (\underline{\epsilon}) = \sigma^2 \mathbf{I}_n$$

$$\text{distribution: } \underline{y} = \underline{\epsilon} + \underline{X}\underline{\beta} \sim N(\underline{X}\underline{\beta}, \sigma^2 \mathbf{I}_n)$$

Vector of residuals:  $\hat{e} = \hat{y} - \hat{\hat{y}} = \hat{y} - H\hat{y} = (I - H)\hat{y}$

$$0) E(\hat{e}) = E[I(I-H)\hat{y}] = (I-H)E[\hat{y}] = (I-H)(X\beta) = X\beta - H(X\beta) = 0.$$

$$\begin{aligned} 0) \text{Var}(\hat{e}) &= \text{Var}(I-H)\hat{y} \\ &= (I-H)\text{Var}(\hat{y})(I-H)^T \\ &= \sigma^2(I-H)(I-H)^T \\ &= \sigma^2(I-H)(I-H) \\ &= \sigma^2(I-HI - IH + HH) \quad \because IH = H, HH = H \\ &= \sigma^2(I-H) \end{aligned}$$

Var-Cov matrix of residual vector  $\hat{e}$ .

$$\text{Var}(e_1) = [I\sigma^2(I-H)]_{11} = \sigma^2(1-h_{11})$$

$$\text{Recall } \text{Var}(e_i) = \sigma^2(1-h_{ii}) \neq \text{Var} s_i = \sigma^2$$

### Inference in regression analysis.

$$\begin{aligned} \text{Var}(\underline{b}) &= \text{Var}((X^T X)^{-1} X^T \hat{y}) \\ &= (X^T X)^{-1} X^T \text{Var}(\hat{y}) (X^T X)^{-1} X^T \\ &= \sigma^2 (X^T X)^{-1} X^T (X^T)^T ((X^T X)^{-1})^T \\ &= \sigma^2 (X^T X)^{-1} (X^T X) ((X^T X)^{-1})^T \\ &= \sigma^2 ((X^T X)^{-1})^T \\ &= \sigma^2 (X^T X)^{-1} \\ &= \text{Var}(\underline{b}) \end{aligned}$$

$$\begin{aligned} \text{SSR} &= b_0^2 S_{xx} = \frac{S_{yy}^2}{S_{xx}} & \text{SSE} &= S_{yy} - \frac{S_{yy}}{S_{xx}} S_{xy}^2 \\ \text{SSTO} &= S_{yy} \end{aligned}$$

$$0) \text{Var}(b_0) = [\sigma^2 (X^T X)^{-1}]_{11}$$

$$S^2(b_0) = \text{MSE}((X^T X)^{-1})_{11} \quad \textcircled{1}$$

$$0) \text{Var}(b_1) = [\sigma^2 (X^T X)^{-1}]_{22}$$

$$S^2(b_1) = \text{MSE}((X^T X)^{-1})_{22} \quad \textcircled{2}$$

$$0) \text{Cov}(b_0, b_1) = [\sigma^2 (X^T X)^{-1}]_{12}$$

use in f\* (1) + (2)

CI for  $b_0, b_1$ . Before:  $S^2(b_0) = \text{MSE}(t_f + S_{xx} \bar{x}^2)$

$$S^2(b_1) = \text{MSE}(\frac{1}{S_{xx}})$$

MSE using Linear Algebra:

$$\begin{aligned} \text{MSE} &= \frac{1}{n-2} \text{SSE} = \frac{1}{n-2} \|\hat{e}\|^2 = \frac{1}{n-2} \|y - Hy\|^2 = \frac{1}{n-2} ((I-H)y)^T ((I-H)y) \\ &= \frac{1}{n-2} (y^T y, y^T Hy) \end{aligned}$$

### Inference of mean response at $X_h$ .

$$\hat{y}_h = b_0 + b_1 X_h = X_h^T \underline{b}, \quad X_h = \begin{bmatrix} 1 \\ X_h \end{bmatrix}$$

$$\text{Var}(\hat{y}_h) = \text{Var}(X_h^T \underline{b}) = X_h^T \text{Var}(\underline{b}) X_h = \sigma^2 X_h (X^T X)^{-1} X_h = \sigma^2 [t_h + S_{xx}(x_h - \bar{x})^2]$$

$$S^2(\hat{y}_h) = \text{MSE} X_h^T (X^T X)^{-1} X_h$$

use this for hypothesis testing / CI for mean response.

### Prediction of resp. variable at $X_h$ .

$$S^2(\text{pred}) = \text{MSE}[1 + X_h^T (X^T X)^{-1} X_h]$$

## Multiple Linear Regression (aka. money prediction) Chapter 6.

We have 2 predictors,  $X_1, X_2$ .

obsn's : n triplets:  $(X_{11}, X_{12}, y_1)$   
 $(X_{n1}, X_{n2}, y_n)$   
 $\vdots$

model:  $y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i, \epsilon_i \sim N(0, \sigma^2)$  i.i.d.

fitted plane:  $\hat{y}_i = b_0 + b_1 X_1 + b_2 X_2$

$\beta_0$ : mean response when predictors are all 0.

$\beta_1$ : increase in mean response when  $X_1$  increases by 1 and  $X_2$  is constant

$\beta_2$ : the other way around (vice-versa).

What if p-predictors?

In general,  $y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{p-1} X_{ip-1} + \epsilon_i, \epsilon_i \sim N(0, \sigma^2)$   
 hyperplane in  $\mathbb{R}^p$

$y = X\beta + \epsilon, \epsilon \sim N(0, \sigma^2 I_n)$

$X \in \mathbb{R}^{n \times p}$   $\beta \in \mathbb{R}^p$

## General Linear Model (GLM)

Not to be confused with generalized linear models (GLiM)

$y_i = c_0 \beta_0 + c_1 \beta_1 + \dots + c_{p-1} \beta_{p-1} + \epsilon_i$   
 functions of  $x_i$  / constants.

Note: still LR model as still linear on params  $\beta_0, \dots, \beta_{p-1}$ !

e.g. Polynomial regression.

$$y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \epsilon_i$$

$$y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i1} X_{i2} + \beta_4 X_{i1}^2 + \beta_5 X_{i2}^2 + \epsilon_i$$

2nd order term

## Qualitative Predictions (categorical var.)

gender: M or F.  $X_1 \begin{cases} 1, & \text{if F} \\ 0, & \text{M} \end{cases}$

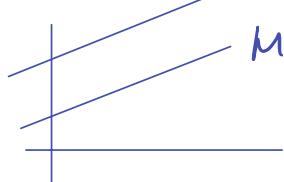
age:  $X_2 \begin{cases} 1, & \text{if M} \\ 0, & \text{F} \end{cases}$

$$\mathbb{E}y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 \begin{cases} F: (\beta_0 + \beta_1) + \beta_2 X_2 \\ M: \beta_0 + \beta_2 X_2 \end{cases}$$

C categories  $\rightarrow (C-1)$  variables

$\hookdownarrow$  "dummy var."

different intercepts.



	$X_3$	$X_4$	$X_5$
disabled	1	0	0
Partially disabled	0	1	0
not disabled	0	0	1

bad.  $X^T X$  stops being invertible.

## Lecture 11: Least Squares

\* Recall  $\underline{Y} = \underline{X}\underline{\beta} + \underline{\epsilon}$        $\underline{\epsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$   
 $n \times 1 \quad n \times p \quad p \times 1 \quad n \times 1 \quad \sim N(n \underline{0}_n, \sigma^2 \underline{I}_n)$

\* LS:  $\underline{\hat{\beta}} = \underset{\underline{\beta}}{\arg \min} Q(\underline{\beta}) = \underset{\underline{\beta}}{\arg \min} \|\underline{Y} - \underline{X}\underline{\beta}\|^2$

↳ Normal equations:  $(\underline{X}^T \underline{X}) \underline{\beta} = \underline{X}^T \underline{Y}$       p equations.

$\Rightarrow \underline{\hat{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$

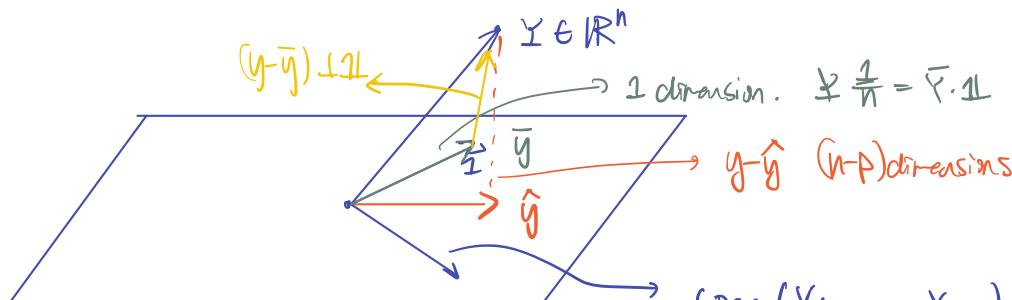
In order to calculate inverse,  $\underline{X}$  must be full rank.

\* MLE:  $L(\underline{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \left( \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 \right) \right\}$

Maximizing wrt  $\underline{\beta}, \sigma^2$ ,  $\hat{\underline{\beta}}^{\text{MLE}} = \hat{\underline{\beta}}^{\text{LS}} \leftarrow$  minimum variance, unbiased, consistent.  
 results as before for p predictors.  $\hat{\underline{\beta}}^{\text{MLE}} \xrightarrow{n \rightarrow \infty} \underline{\beta}$

Fitted values:  $\hat{\underline{y}} = \underline{H}\underline{Y} = \underline{P}\underline{Y} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$

Residuals:  $\underline{\epsilon} = (\underline{I} - \underline{H}) \underline{Y}$



Span  $(\underline{X}_1, \dots, \underline{X}_{p-1})$   
 $X_i$  means  $i$ th column of  $\underline{X}$ .  
 $(p-1)$  dimensions.

### ANOVA Table:

Source	SS	df = dim	MS.
error	SSE	$(n-p)$	$MSE = \frac{SSE}{n-p} = \ y - \hat{y}\ ^2 / (n-p)$
regression	SSR	$(p-1)$	$MSR = \frac{SSR}{p-1} = \ \hat{y} - \bar{y}\ ^2 / (p-1)$
total	SSTO	$(n-1)$	

$\min_{\underline{c}} \sum_i (y_i - c)^2, \quad \underline{c} = c \cdot \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \Rightarrow \min_{\underline{c}} \|\underline{y} - c \cdot \underline{1}\|^2, \quad c = c \cdot \underline{1}, \quad c = \bar{y}$

Like before,  $E(MSE) = \sigma^2$  (unbiased)

$E(MSR) = \sigma^2 + \underline{\beta}^T \underline{A} \underline{\beta}$

↳ Symmetric and positive definite

all eigenvalues are strictly positive.

If  $\beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$ , then  $MSE \approx MSR$ .

If at least one of  $\beta_k \neq 0$ , then  $MSE < MSR$ .

★ ★ F-test (test of overall fit / model utility test).

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0,$$

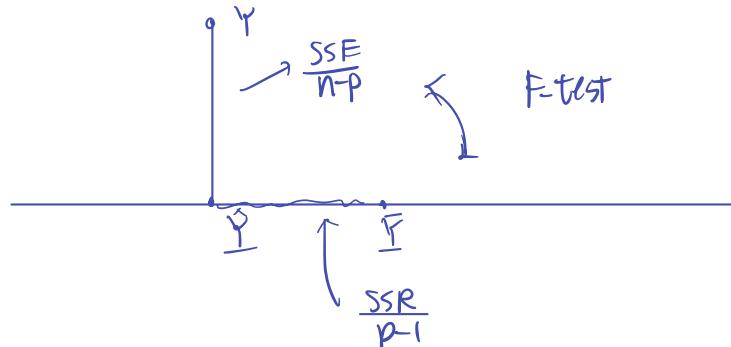
$H_a$ : at least one of them is non-zero.

$$F^* = \frac{MSR}{MSE} = \frac{SSR/(p-1)}{SSE/(n-p)}$$

If  $F^* \leq F(1-\alpha, p-1, n-p)$ , conclude  $H_0$ .

If  $F^* > F(1-\alpha, p-1, n-p)$ , conclude  $H_a$ .

F-test: Is there a useful relationship b/w  $Y$  and  $X_1, \dots, X_{p-1}$ ?



Now, F-test is no longer equivalent to t-test (one  $\beta_i$ ).

★ ★ Coefficient of multiple determination ( $R^2$ )

$$R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO} \quad \text{relative reduction in the variation of } y \text{ due to all predictors } X_1, \dots, X_{p-1}.$$

$$R^2 = 0 \Rightarrow b_1 = b_2 = \dots = b_{p-1} = 0.$$

$$y_i = b_0, \quad SSR = \hat{y}_i - \bar{y} \quad \text{No linear relationship}$$

$$R^2 = 1 \Rightarrow SSE = 0, \quad y_i - \hat{y}_i = 0. \quad \text{Perfect fit.}$$

Regression pass through all datapoint.

As we add predictors,  $R^2$  increasing "always".

$$\underset{n \times 1 \quad n \times p \quad p \times 1}{\arg \min \| \mathbf{Y} - \mathbf{X} \mathbf{B} \|^2} \quad \text{b/c you search the minimum at a larger \& larger dimension of space.}$$

★ ★ Adjusted  $R^2$ , penalize # of predictors.

$$R^2_{\text{adjusted}} = 1 - \frac{SSE/(n-p)}{SSTO/(n-1)} = 1 - \frac{MSE}{MSTO} = 1 - \frac{n-1}{n-p} \frac{SSE}{SSTO}$$

as  $p \uparrow$ ,  $\frac{n-1}{n-p} \uparrow$ ,  $R^2_{\text{adj}} \downarrow$ , (if  $SSE$  stays the same)

★ ★ Note:  $R^2 = \text{coeff. mult. det.} = \text{coeff. simple det. when } Y_i \sim \hat{Y}_i$

★ ★ Inference on  $\beta_k$ 's

$$\mathbb{E}(\mathbf{b}) = \mathbf{B}; \quad \text{Var}(\mathbf{b}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

$$\mathcal{S}(\mathbf{b}) = MSE (\mathbf{X}^T \mathbf{X})^{-1}$$

$\hookrightarrow (n-p)$  dim.

$$\hat{S}(b_k) = [S^2(b)]_{k+1, k+1}$$

$$CI_{t_d}: b_k \pm S(b_k)t(1 - \frac{\alpha}{2}; n-p)$$

$$H_0: \beta_k = 0$$

$$H_a: \beta_k \neq 0$$

$$T.S: t^* = \frac{b_k}{\hat{S}(b_k)} \sim t_{n-p} \quad \text{If } |t^*| < t(1 - \frac{\alpha}{2}; n-p), \text{ conclude } H_0.$$

\*\* If q coeffs estimated:

Berfframi procedures

$$\text{Int intervals } b_k \pm S(b_k)t(1 - \frac{\alpha}{2q}; n-p)$$

(recall joint interval for \beta\_0, \beta\_1)

## Lecture 12 2023 HT OA 10 (=)

### Inferences in MLR.

→ mean response var.

$$\hat{Y}_h = \begin{bmatrix} 1 \\ X_{h,1} \\ \vdots \\ X_{h,p+1} \end{bmatrix} \in \mathbb{R}^p$$

$\hat{Y}_h = \hat{X}_h^T \hat{b}$  = fitted value  
(estimated mean response).

$$S^2(\hat{Y}_h) = \text{MSE}(X_h^T(X^T X)^{-1} X_h) = \underbrace{\hat{X}_h^T S^2(b) \hat{X}_h}_{p \times p \quad p \times p \quad p \times 1}$$

(Scalar)

$$(CI)_{1-\alpha} = \hat{Y}_h \pm t(1 - \frac{\alpha}{2}, n-p) S(\hat{Y}_h)$$

CI for a single mean.

→ Confidence region for regression surface

$$\hat{Y}_h \pm W \cdot S(\hat{Y}_h) \quad W = \sqrt{p F_{t(\alpha/2, n-p)}}$$

"Combined" CI for all possible means. (generalization of Working-Hoteling Procedure).

→ Simultaneous CI for g means

$$\hat{Y}_h \pm t(1 - \frac{\alpha}{g}, n-p) \cdot S(\hat{Y}_h) \quad (\text{Bonferroni procedure}).$$

→ prediction of a new observations

$$(PI) \quad \hat{Y}_h \pm t(1 - \frac{\alpha}{2}, n-p) S(\text{pred})$$

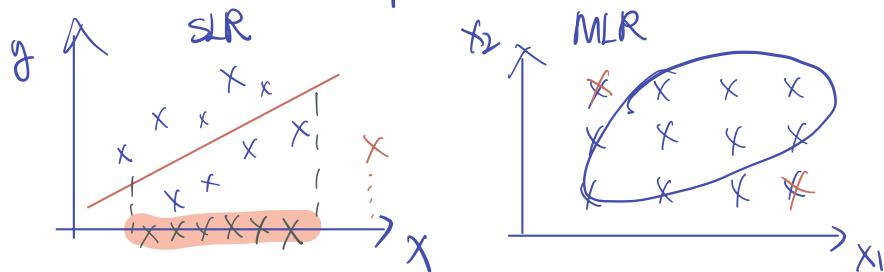
$$S^2(\text{pred}) = \text{MSE}(1 + \hat{X}_h^T (X^T X)^{-1} \hat{X}_h)$$

→ Simultaneous prediction intervals:

- Same as before in SLR:
- Bonferroni
- Scheffe.

Always pick the Narrow One.

### Issue: Hidden Extrapolations.



Cannot really detect the outlier by this graph.

2 predictors.

### Diagnostics in MLR.

① Scatter plots for all predictors against all predictors.  
Pairs (dataset).

② Correlation matrix. → niche: 3D Scatterplot ( $y_i$  vs.  $X_{i,k}$  &  $X_{i,j}$ )  
Cor (dataset)

③ Fit the model. Using lm (formula, df)  
See the p-values.

★ Diagnostics from SLR still apply

- $e_i$  vs.  $\hat{y}_i$  → nonlinearity.
- Hist. & QQ plot for  $e_i$ 's → Normality
- $e_i$  vs.  $X_{i,k}$  for all  $k$ s.
- $e_i$  vs. Omitted Predictors. → do we need more predictors?
- $|e_i|$  or  $e_i^2$  vs.  $\hat{y}_i$  → outliers.
- $|e_i|$  or  $e_i^2$  vs.  $X_{i,k}$  → Nonconstancy of variance

Interaction terms:  $X_{i,k} X_{i,j}$ , plot  $\hat{y}_i$  vs. Interaction terms to see if needed.

★ Remedial measures.

Shapiro-Wilks test for normality of error.

Brown-Forsythe test for nonconstancy of variance.

★ If Variance is NOT constant:

Transformation on  $y$ . Sometimes transform  $X$  also to fix the linearity.

★ If nonlinear,

Transformation on  $X$ , or add some predictor  $f(x)$ .

$$Y = \beta_1 X_1 + \beta_2 X_2 + \beta_0$$

$$Y = \beta_1 X_1 + \beta_2 X_2 + \beta_{12} X_1 X_2 + \beta_0.$$

★ Extra sum of squares (Ch. 7)

Model #1:  $Y \sim X_1$

Model #2:  $Y \sim X_1, X_2$

Q: Worth to add  $X_2$  to the model?

$$\begin{aligned} \text{ANOVA: } S\text{STO}_{\text{model 1}} &= S\text{SE}(X_1) + S\text{SR}(X_1) \leq M_1 \\ &= S\text{SE}(X_1 + X_2) + S\text{SR}(X_1 + X_2) \leq M_2 \end{aligned}$$

$$S\text{SE}(X_1) \geq S\text{SE}(X_1 + X_2)$$

$$S\text{SR}(X_1) \leq S\text{SR}(X_1 + X_2)$$

$$\begin{aligned} \text{Extra sum of squares: } S\text{SR}(X_2 | X_1) &= S\text{SE}(X_1) - S\text{SE}(X_1, X_2) \quad (1) \\ &= S\text{SR}(X_1, X_2) - S\text{SR}(X_1) \quad (2) \end{aligned}$$

(1) How much the error decrease when we add  $X_2$  to model  $Y \sim X$ ?

(2) How much the explained var. in  $Y$  increased when we added  $X_2$  to  $Y \sim X$ ?  
↳ Need to test if  $S\text{SR}(X_2 | X_1)$  is large enough.

→ General Linear Test (Lecture 5) :

$$(R) \text{ Reduced : } Y = \beta_0 + \beta_1 X_1 + \epsilon.$$

$$(F) \text{ Full model : } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon. \quad df: n-2 \\ df: n-3.$$

$$F^* = \frac{SSE(R) - SSE(F)}{df_R - df_F} / \frac{\text{SSF}(F)}{df_F} = \frac{SSE(X_1) - SSE(X_1, X_2)}{(n-2) - (n-3)} / \frac{\text{SSF}(X_1, X_2)}{n-3}$$

$$= \frac{SSRL(X_2 | X_1)}{1} / \frac{SSF(X_1, X_2)}{n-3} = \text{MSR}(X_2 | X_1) / \text{MSE}(X_1, X_2) \\ \sim F(1, n-3)$$



$$t\text{ test } (t^*)^2 = \frac{(\hat{\beta}_2)^2}{S^2(\hat{\beta}_2)}$$

### \* Marginal Test :

Is  $X_2$  worth adding to the model which already has  $X_1$ ?

Relationship to model utility test (F test for overall fit)

$$H_0: \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$$

$$H_a: \text{at least one } \beta_k \neq 0.$$

$$R: Y \sim 1 \quad df: n-1$$

$$F: Y \sim 1, X_1, X_2, \dots, X_{p-1} \quad df: n-p$$

$$F^* = \frac{SSE(1) - SSE(1, X_1, \dots, X_{p-1})}{(n-1) - (n-p)} / \frac{SSE(1, \dots, X_{p-1})}{(n-p)}$$

$$= \frac{SSRL(1, X_1, \dots, X_{p-1})}{(p-1)} / \frac{SSE(1, X_1, \dots, X_{p-1})}{(n-p)} = \frac{\text{MSR}}{\text{MSE}}$$

### \* General linear test

In general, we can test any subset of variables:

$$H_0: \beta_3 = \beta_5 = 0$$

$$H_a: \text{at least one } \beta_j \neq 0.$$

↳ F test only

### \* Coefficient of partial determination.

$$R^2_{y_{2|1}} = \frac{SSE(X_1) - SSE(X_1, X_2)}{SSE(X_1)} = \frac{SSR(X_2 | X_1)}{SSE(X_1)}$$

→ measures relative reduction in variation in  $y$ .

when  $X_2$  is included in model  $Y \sim X_1$ .

$$SSTO = SSR(X_1) + SSR(X_2 | X_1) + SSE(X_1, X_2)$$

Lecture 13, Oct. 16<sup>th</sup>

ASDA 913/613

Today :- general linear test : alternative form

- coefficient of partial determin'n
  - added variable plot / partial regression plot
- standardized multiple regression model
- corr'n transform
- start multicollinearity

No section on Friday: Fall break!

► General linear test (alternative form) (\* not in KNNL)

(F) General linear test

various F-tests:

(A) Model utility test / overall fit

$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \\ H_a: \text{at least one } \beta_k \neq 0 \end{cases}$$

$$F^* = \frac{\text{MSR}}{\text{MSE}} = \frac{\frac{\text{SSR}(x_1, \dots, x_{p-1})}{p-1}}{\frac{\text{SSE}(x_1, \dots, x_p)}{n-p}}$$

$\sim F(p-1, n-p) \text{ under } H_0$

(B) Individual  $\beta_k$ 's :  $H_0: \beta_k = 0$  vs.

$H_a: \beta_k \neq 0$

$$F^* = \frac{\hat{\sigma}_k^2}{\hat{\sigma}^2(\hat{\beta}_k)} \sim F(1, n-p) \text{ under } H_0$$

(C) Test various subsets of  $\beta_k$ 's, e.g.:

$$H_0: \beta_3 = \beta_5 = 0 \quad \text{vs.} \quad H_a: \beta_3 \neq 0, \beta_5 \neq 0$$

We use G.L.T. : compare reduced models vs. full model

$\Rightarrow$  We can (and will) generalize this

$$H_0: \underline{C} \underline{\beta} = \underline{0} \quad \text{vs.} \quad H_a: \underline{C} \underline{\beta} \neq \underline{0}$$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ m \times p & p \times 1 & m \times 1 \end{matrix}$

rank(C) = q ( $\neq$ )

(A) test of overall fit :  $H_0: \beta_1 = \dots = \beta_{p-1} = 0$

$$C = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & & & \ddots & 0 \\ 0 & 0 & \dots & \dots & 0 & 1 \end{bmatrix}_{p-1 \times p}, \underline{D} = \underline{0} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}_{p-1}$$

$$\underline{C} \cdot \underline{\beta} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} = \underline{\alpha} = \begin{bmatrix} 0 \\ \vdots \\ \alpha \end{bmatrix}$$

$$\text{rank}(\underline{C}) = p-1$$

$$(B) \text{ (individual)} \quad H_0: \beta_k = 0$$

$$\underline{C} = [0 \ 0 \ \dots \ 1 \ 0 \ \dots \ 0] \Rightarrow \underline{C} \cdot \underline{\beta} = \beta_k$$

$\uparrow$   
 $k+1$  position

$$\underline{\alpha} = 0$$

$$\underline{C} \cdot \underline{\beta} = \underline{\alpha} = 0 \quad \text{rank}(\underline{C}) = 1$$

Other interesting tests! E.g.:  $H_0: \beta_1 = \beta_2 \Leftrightarrow \underline{\beta}_1 - \underline{\beta}_2 = 0$

$$\underline{C} = [0 \ 1 \ -1 \ 0 \ \dots \ 0] \quad \underline{C} \cdot \underline{\beta} = \underline{\beta}_1 - \underline{\beta}_2; \underline{\alpha} = 0$$

Test statistic?

$$\text{If } \underline{x} \sim N(\mu, \sigma^2) \Rightarrow \frac{(\underline{x} - \mu)^2}{\sigma^2} \sim \chi_q^2$$

$$\text{Fact: } \underline{x} \sim N_q(\mu, \Sigma) \quad (\underline{x} - \mu)^T \Sigma^{-1} (\underline{x} - \mu) \sim \chi_q^2$$

$$\text{So: } \underline{\beta} \sim N(\underline{\beta}_f, \sigma^2 (X^\top X)^{-1}) \Rightarrow$$

$$\underline{\Sigma} \cdot \underline{\beta} \sim N(\underline{\Sigma} \cdot \underline{\beta}_f, \underbrace{\sigma^2 \underline{\Sigma} (X^\top X)^{-1} \underline{\Sigma}^\top}_{\text{var}(\underline{\Sigma} \cdot \underline{\beta})})$$

$$\text{rank}(\underline{\Sigma}) = q$$

$$(\underline{\Sigma} \cdot \underline{\beta} - \underline{\Sigma} \cdot \underline{\beta}_f)^\top (\text{var}(\underline{\Sigma} \cdot \underline{\beta}))^{-1} (\underline{\Sigma} \cdot \underline{\beta} - \underline{\Sigma} \cdot \underline{\beta}_f) \sim \chi_q^2$$

$\sigma^2$  has to be estimated as usual!

$$\frac{SSE}{\sigma^2} \sim \chi_{n-p}^2$$

$$(\underline{\Sigma} \cdot \underline{\beta} - \underline{\beta}_f)^\top [\text{var}(\underline{\Sigma} \cdot \underline{\beta})]^{-1} (\underline{\Sigma} \cdot \underline{\beta} - \underline{\beta}_f)$$

q

SSE

n-p

$$= \frac{n-p}{q} \frac{[(\underline{\Sigma} \cdot \underline{\beta} - \underline{\beta}_f)^\top (\underline{\Sigma} (X^\top X)^{-1} \underline{\Sigma}^\top)^{-1} (\underline{\Sigma} \cdot \underline{\beta} - \underline{\beta}_f)]}{SSE}$$

$$\sim F_{q, n-p}$$

$$\frac{\chi_q^2}{q}$$

$$\frac{\chi_{n-p}^2}{n-p}$$

If  $C = [0 \ 1 \ 0 \dots 0]$   $\Rightarrow$  the above is  $\frac{\hat{b}_1^2}{S^2(\hat{b}_1)}$ , i.e.

$$H_0: \beta_1 = 0$$

geometrically:

$$= \frac{\|\hat{Y}_F - \hat{Y}_R\|^2}{\|Y - \hat{Y}_F\|^2} \cdot \frac{n-p}{q} = \frac{\|Y - \hat{Y}_R\|^2}{\|Y - \hat{Y}_F\|^2} - \frac{\|Y - \hat{Y}_F\|^2}{\|Y - \hat{Y}_F\|^2} \cdot \frac{n-p}{q}$$

SSE(R) SSE(F)  
SSE(F)

$$\text{where } \hat{Y}_R = X \cdot \hat{b}_R \text{ where } \hat{b}_R = \underset{\beta}{\operatorname{arg\min}} \|Y - X\beta\|^2$$

$$\text{s.t. } C\beta = \underline{0}$$

Very versatile F-test!

We can conduct tests:

$$H_0: \beta_3 = \beta_5$$

$$H_0: \beta_3 = \beta_5 = 2$$

$$H_0: \beta_3 = 5, \beta_5 = 7$$

$$H_0: 2\beta_3 + 9\beta_5 = 12$$

e.g.:  $H_0: \beta_1 = \beta_2 = \beta_3$

G

2eqns:  $\beta_1 = \beta_2 \Rightarrow \beta_1 - \beta_2 = 0$

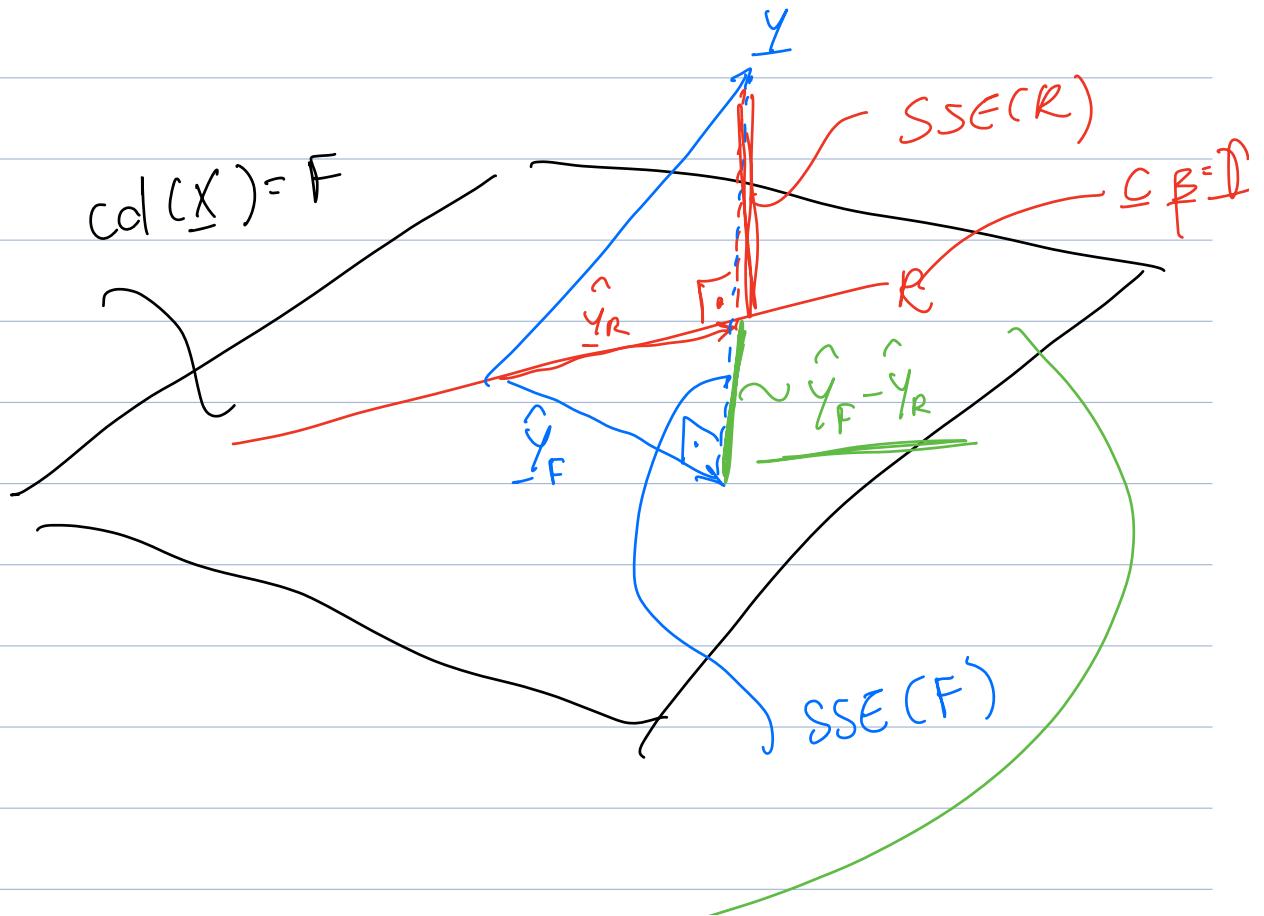
$\beta_2 = \beta_3 \Rightarrow \beta_2 - \beta_3 = 0$

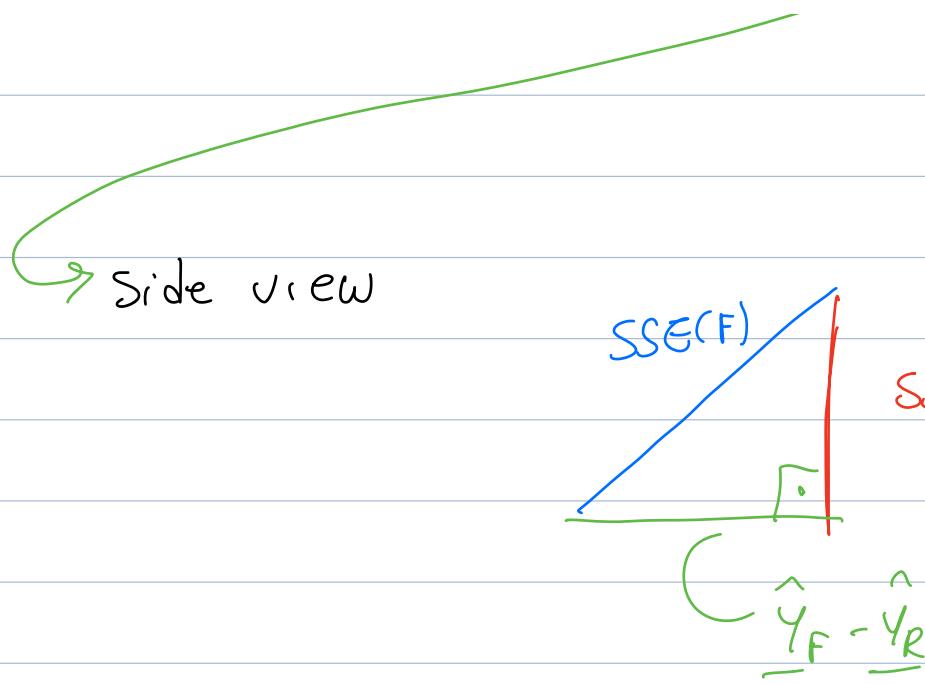
$$C = \begin{bmatrix} 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

C is not unique!

but outcome of the test will be the same!

$$\text{col } (\underline{X}) = F$$





► Coefficient of partial determination (Ch. 7.4)

→ coeff. simple def.:  $R^2 = SCR$

→ coeff. multiple def.:  $R^2 = MLR$

→ coeff. of partial def.:  $R_{Y_{Z11}}^2$

$$R_{Y_{Z11}}^2 = \frac{\overbrace{SSE(x_1)} - \overbrace{SSE(x_1, x_2)}}{SSE(x_1)}$$

= proportionate reduction in variation of  $Y$  when  $x_1$  is included in model and we add  $x_2$

meaning:

$$\textcircled{1} \quad Y \sim X_1 \Rightarrow e_i(Y_i | X_1) = Y_i - \hat{Y}_i(X_1)$$

$$\textcircled{2} \quad X_2 \sim X_1 \Rightarrow e_i(X_{i2} | X_1) = X_{i2} - \hat{X}_{i2}(X_1)$$

Regress  $e_i(Y_i | X_2) \sim e_i(X_{i2} | X_1) \Rightarrow$   
 $R^2 = \text{coeff. of simple def'n:}$   
=

coeff. of partial def'n  $R^2_{Y_{i1} | X_2}$

→ measures the strength of linear relationship  
btw  $Y$  &  $X_2$  when  $X_1$  taken into account

→ Interesting fact:  $\hat{Y}_i = b_0 + b_1 X_1 + b_2 X_2$  are the same!!

$$\hat{e}_i(Y_i | X_1) = \hat{b}_0 + b_2 e_i(X_{i2} | X_1)$$

→ check R-script!

► Additional diagnostic tool

↳ added variable plot / partial reg. plot

plot of  $e_i(Y|X_1)$  vs  $e_i(X_2|X_1)$

↳ check constancy of variance, non linear relationship

► Standardized multiple linear regression model

- to make  $\beta_k$ 's comparable
- for numerical stability in the matrix  $(X^T X)^{-1}$

standardization:  $\frac{y_i - \bar{y}}{s_y}, \quad \frac{x_{ik} - \bar{x}_k}{s_k}$

↳ sample sd( $y$ )      ↳ sample sd( $x_k$ )

► Correlation transform

$$y_i^* = \frac{1}{\sqrt{n-1}} \frac{y_i - \bar{y}}{s_y}$$

$$x_{ik}^* = \frac{1}{\sqrt{n-1}} \frac{x_{ik} - \bar{x}_k}{s_k}$$

In these standardized models, there is no intercept!

$$Y_i^* = \beta_0^* X_{i1}^* + \beta_1^* X_{i2}^* + \dots + \beta_{p-1}^* X_{ip-1}^* + \epsilon_i^*$$

all sample means are 0! so model has to pass through origin  $\rightarrow$  no  $\beta_0^*$  (intercept)

$$\underline{X} = \begin{bmatrix} X_{11}^* & \dots & X_{1,p-1}^* \\ \vdots & & \vdots \\ \vdots & & \vdots \\ X_{n-1}^* & & X_{n,p-1}^* \end{bmatrix}$$

Fact:  $\underline{X}^\top \underline{X} = r_{XX}$

correlation matrix  
btw  $X_i$  &  $X_j$ !

$$r_{ij} = \frac{1}{n-1} \sum_k \frac{(X_{ki} - \bar{X}_i)(X_{kj} - \bar{X}_j)}{S_i S_j}$$

$$|r_{ij}| \leq 1$$

Similarly,  $\underline{X}^\top \underline{Y} = r_{YX} = \begin{bmatrix} r_{Y,1} \\ r_{Y,2} \\ \vdots \\ r_{Y,p-1} \end{bmatrix} \rightsquigarrow r_{Y,j} : \text{sample corr'n btwn } Y \text{ & } X_j$

$$\underline{b^*} = \underline{r_{xx}^{-1}} \cdot \underline{r_{xy}} = \begin{bmatrix} b_1^* \\ \vdots \\ b_{p-1}^* \end{bmatrix} \rightsquigarrow \begin{array}{l} \text{can recover } b_k^* \text{'s} \\ \text{(unstandardized)} \\ \text{as} \end{array}$$

$$b_k = \frac{s_y}{s_k} \cdot b_k^*$$

$$\begin{aligned} b_0 &= \bar{y} - b_1^* \bar{x}_1 - b_2^* \bar{x}_2 - \dots \\ &\quad - b_{p-1}^* \bar{x}_{p-1} \end{aligned}$$

## ► Multicollinearity (A 7.6)

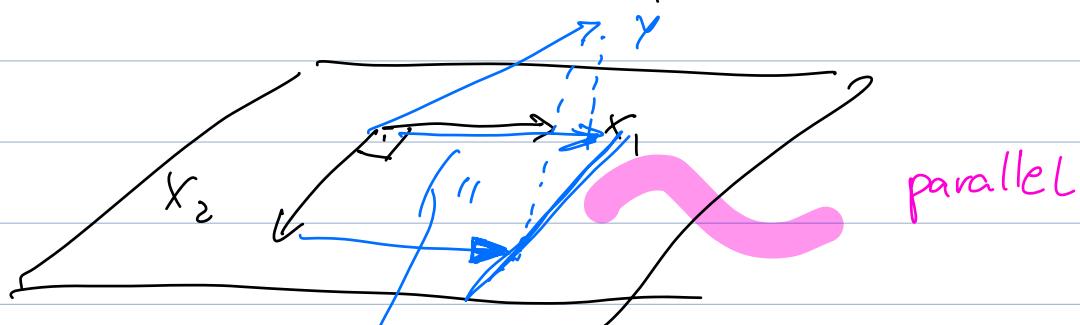
wrong def'n : predictors are correlated

||  
multicollinearity / intercorrelation

o) when  $\text{corr}(X_1, X_2) = 0$ ,

$$SSR(X_2 | X_1) = SSR(X_2)$$

$$SSR(X_1 | X_2) = SSR(X_1)$$



SSR is the same

If we fit 3 models:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$Y = \tilde{\beta}_0 + \beta_1 X_1$$

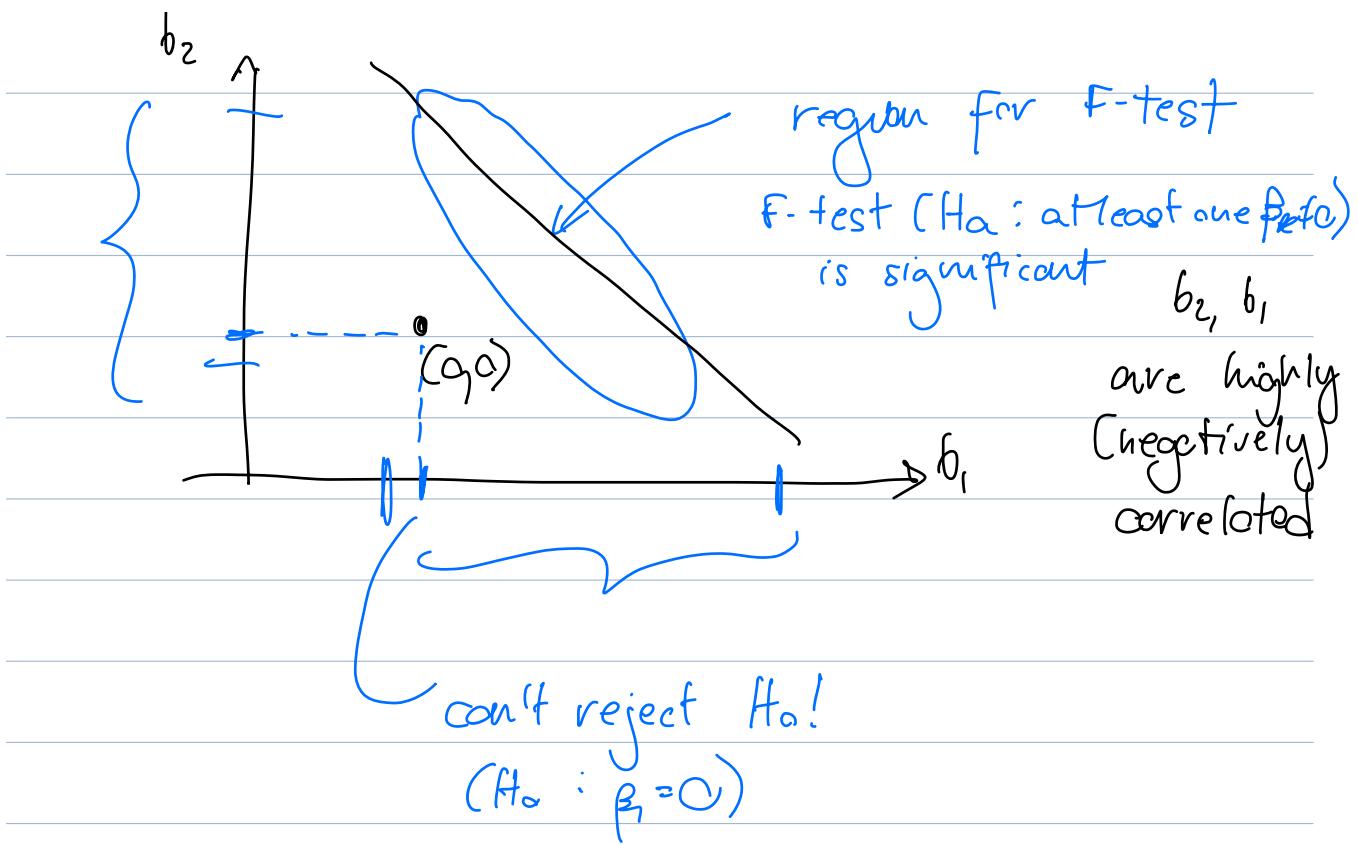
$$Y = \hat{\beta}_0 + \beta_2 X_2$$

• when  $\text{corr}(X_1, X_2) = 1$ , e.g.  $X_1 = X_2$

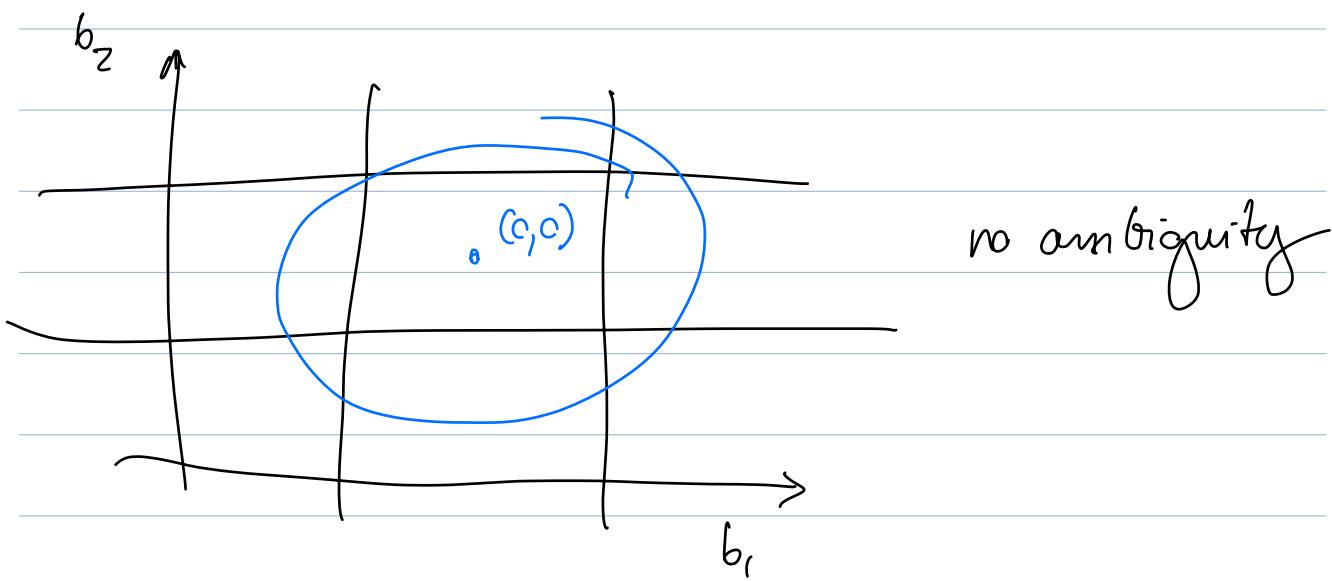
$$\beta_0 = 2, \beta_1 = 3, \beta_2 = 5$$

$$\begin{aligned} \hat{Y} &= 2 + 3X_2 + 5X_2 \\ &= 2 + 4X_1 + 4X_2 \\ &= 2 + (-10)X_1 + 18X_2 \end{aligned} \quad \left. \begin{array}{l} \text{SD } (\beta_k) \text{ are very} \\ \text{large} \end{array} \right\}$$

c)  $\text{corr}(\beta_2, \beta_1)$  is high!



$$o) \text{ corr}(b_2, b_1) \approx 0$$

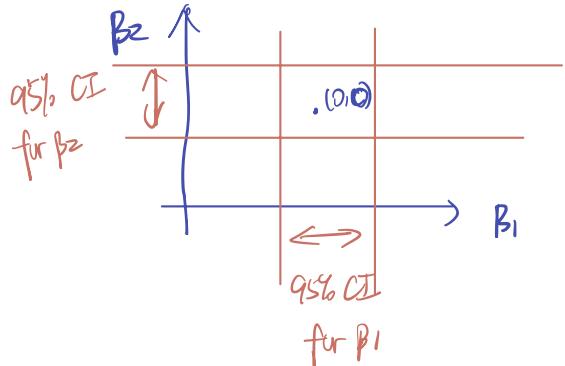


## Lecture 14 2023 LoFi (BA) . Chp 10.

\* Collinearity Significant vs. non-significant.

$\beta_k$  is significant  $\Leftrightarrow H_0$  nonzero value of  $\beta_k$  (H.T.)  
 $\Leftrightarrow p\text{ value} < \alpha$ .  
 $\Leftrightarrow 0 \notin CI_{1-\alpha}(\beta_k)$

\* Confidence region for marginal tests.



)  $H_0: \beta_1 = 0$ . Do NOT reject

)  $H_0: \beta_2 = 0$ . Do NOT reject.

)  $H_0: \beta_1 = 0$  (do not reject)

$\Leftrightarrow$  comparing 2 models  $y \sim X_2$  vs.  $y \sim X_1 + X_2$ .

\* Confidence region for F-test.

$$H_0: \beta_1 = \beta_2 = 0.$$

$$\text{Recall: } F^* = \frac{(\underline{c} \underline{b})^T (\underline{c} (\underline{X}^T \underline{X})^{-1} \underline{c})^{-1} (\underline{c} \underline{b})}{SSE} \quad \frac{n-p}{q} = \text{rank}(\underline{c}). \quad \underline{\chi} = 0$$

$$F^* = \underline{b}^T \underline{\Sigma} \underline{b} \quad \text{PSD matrix. all eigenvalues of } \underline{\Sigma} \text{ are } \geq 0.$$

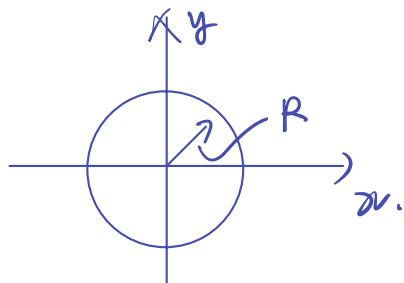
$$\underline{\Sigma} \underline{v} = \pi_i \underline{v} + \underline{u}_i, \pi_i \geq \pi_i \geq 0.$$

$$F^* = \underline{b}^T \underline{\Sigma} \underline{b} \leq \text{critical value} \rightarrow \text{non-rejection region}$$

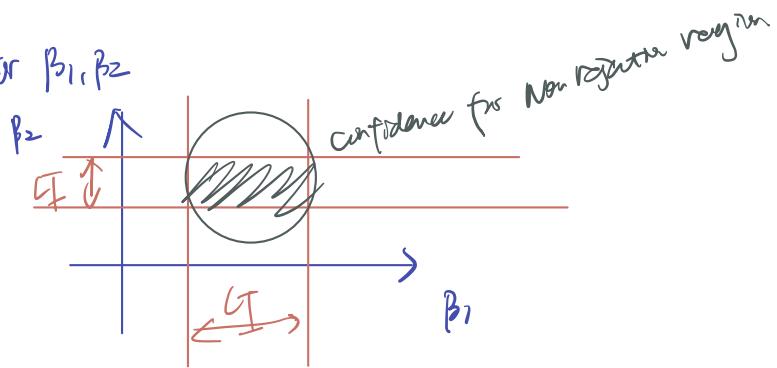
quadratic form in  $\underline{b}$

Assume: If  $\underline{\Sigma} = \underline{I}$ ,  $\underline{b}^T \underline{I} \underline{b} = \underline{b}^T \underline{b} = \sum_{i=1}^p b_i^2 \leq \text{c.v.}$

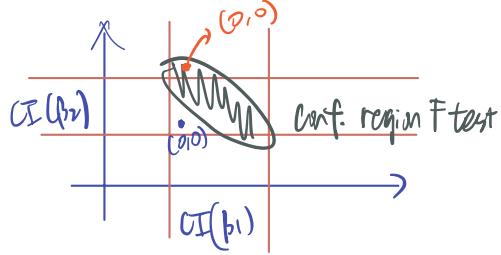
$$x^2 + y^2 \leq R^2$$



For  $\beta_1, \beta_2$



+ If  $\Sigma \neq I \Rightarrow b^T \Sigma^{-1} b$  defines an ellipsoid.



$H_0: \beta_1 = \beta_2 = 0$ , Reject  $H_0$ .

- (1)  $H_0: \beta_1 = 0$  do not reject
- (2)  $H_0: \beta_2 = 0$  do not reject
- (3)  $H_0: \beta_1 = \beta_2 = 0$  reject.

different scenarios.

$H_0:$

- (1)  $\beta_1 = 0$  do not reject
- (2)  $\beta_2 = 0$  reject
- (3)  $\beta_1 = \beta_2 = 0$  do not reject.

In general, multicollinearity is when the corr ( $x_i, \sum_{j \neq i} x_j$ ) is high.  $\Rightarrow \det(x^T x) \leq 0$ .

$\Rightarrow$  large standard deviation for b's.  
Non significant results.

Note: Correlation matrix C or pairwise correlations  $r_{ij}$   
do NOT measure this.  $\rightarrow$  we need more diagnostics.

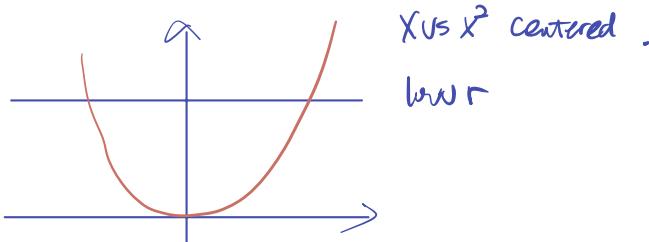
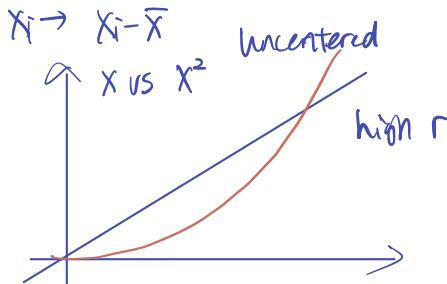
## Regression Models with Qualitative and Quantitative Predictors (Chap. 8),

$\rightarrow$  One pred, 2nd order model.

$$Y_i = \beta_0 + \beta_1 x_i + \beta_{11} x_i^2 + \epsilon_i, \text{ s.t. } N(0, \sigma^2)$$

$\rightarrow$  One pred, 3rd order model.

$$Y_i = \beta_0 + \beta_1 x_i + \beta_{11} x_i^2 + \beta_{111} x_i^3 + \epsilon_i$$



Caution: always center  $X_i$ 's for polynomial regression.

$\rightarrow$  Two predictors, 2nd order:

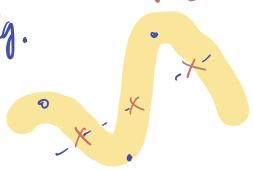
$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_{11} x_{i1}^2 + \beta_{22} x_{i2}^2 + \beta_{12} x_{i1} x_{i2} + \epsilon_i$$

$\underbrace{\quad}_{\text{Interaction term}}$   $\underbrace{\quad}_{\text{Interaction effect coeff.}}$

effect of  $X_1$  on  $y$  depends on  $X_2$

Caution for polynomial model: Problem with interpolating & extrapolating.

Eg.



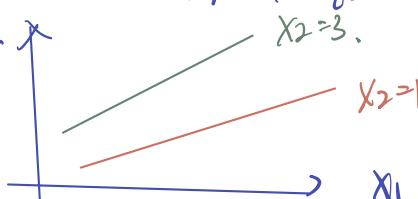
$x \rightarrow \pm \infty$  hyper order polys. oscillate & diverge.

$$E(y) = 10 + 2x_1 + 5x_2 + 0.5x_1 x_2$$

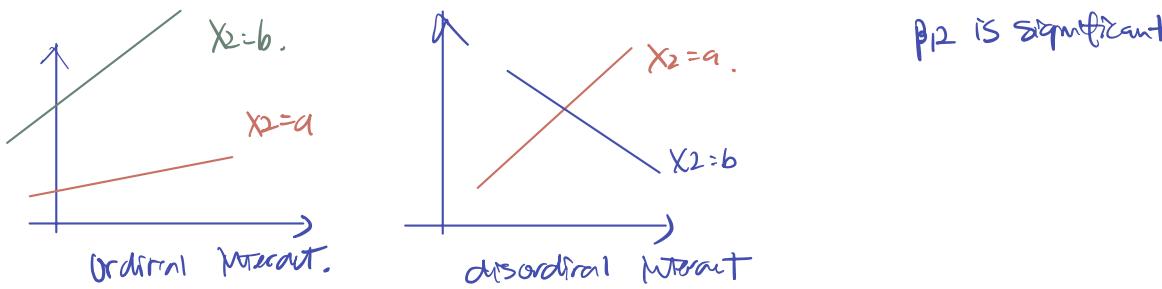
$$x_2 = 1, E(y) = 10 + 2x_1 + 5 + 0.5x_1 = 15 + 2.5x_1$$

$$x_2 = 3, E(y) = 25 + 3.5x_1$$

Interaction plot:  $y$  vs  $x_1$



If parallel,  $\beta_{12}$  is NOT significant.



quantitative predictors : numerical

qualitative " " : categorical.

C categories  $\rightarrow$  (-1) var. ("dummy")

cat c:  $X_1 = X_2 = \dots = X_{c-1} = 0$ .

$\hookrightarrow$  baseline category .

Why don't we use C variables, e.g.:  $X_C = \begin{cases} 1 & \text{if baseline} \\ 0 & \text{o.w.} \end{cases}$

	$X_1$	$X_2$	$\dots$	$X_C$	" $X_0$ "	
Cat 1	1	0	$\dots$	0	1	$\Rightarrow$ Can't invert $\underline{X^T X}$ .
Cat 2	0	1	$\dots$	0	1	
Cat 3	0	0	$\dots$	0	1	
:	:	:	$\vdots$	:	$\vdots$	
Cat C	0	0	$\dots$	1	1	

$$\sum_{1 \leq i \leq c} X_i$$

- .) If all Var's are Qualitative  $\Rightarrow$  ANOVA (Analysis of var. Chp. 1b)
- .) If Var's are both qual. & quant.  $\Rightarrow$  ANCOVA (analysis of cov.)  
 $\hookrightarrow$  math (e.g. dosage of drug).

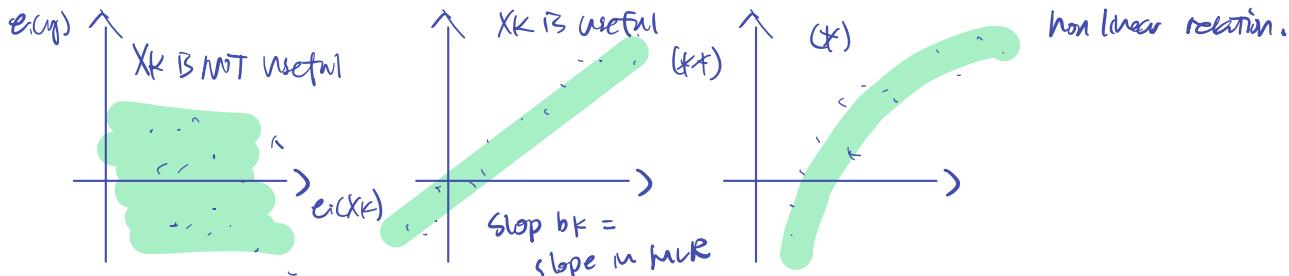
## ¶¶ Diagnosestus

- .) Added variable / partial reg. plot

$y \sim$  all pred's but  $X_k \rightarrow$  residual  $e(y)$

$X_k \sim$  all " " " " "  $\rightarrow$  residual  $e(X_k)$   $\nearrow$  contain unexplained variation in  $y$  &  $X_k$ .

Plot  $e(y)$  vs.  $e(X_k)$

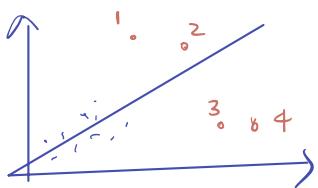


The plot should show no patterns normally.

$$y_i = b_0 + b_1 X_{i1} + \dots + b_k X_{ik} + \dots + b_{p-1} X_{ip}, p-1$$

- Check:
- improper function form ( $f(x)$ )
  - whether  $X_k$  is useful ( $f(x_k)$ ).
  - non constant variance. (like the cone shape)

### 2) Outlying observations

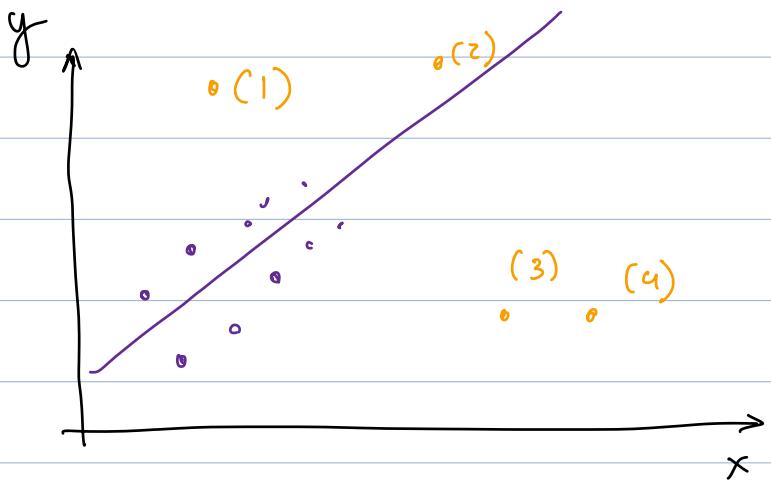


Outliers?	Outlying wrt X	2, 3, 4.
→	" Y	1, , 3, 4.
	" $X+Y$ .	3, 4.
	Compare with fitted line.	

## ► Lecture 15, ASDA 413/613

Today: → outlying obsn's / outliers (10.2)  
 → hidden extrapolations (10.3)

o) outlying obsn's



→ outlying wrt x (2), (3), (4)

→ outlying wrt y (1), (3), (4) → compare w/

→ outlying wrt x & y (3), (4) fitted values

$$e_i = y_i - \hat{y}_i \rightarrow e_i^* = \frac{e_i}{\sqrt{MSE}}$$

semi-studentized residuals

$$\underline{e} = \underline{y} - \hat{\underline{y}} = (\underline{I} - \underline{H}) \underline{y}$$

$$\text{var}(\underline{\epsilon}) = \text{var}((\underline{I} - \underline{H}) \underline{Y}) = (\underline{I} - \underline{H}) \text{var}(\underline{Y}) (\underline{I} - \underline{H})^T$$

↓ symm.

$$= (\underline{I} - \underline{H}) \text{var}(\underline{Y}) (\underline{I} - \underline{H}) = \sigma^2 (\underline{I} - \underline{H}) (\underline{I} - \underline{H}) =$$

↗ ↓ idempotent

$$\sigma^2 = \sigma^2 (\underline{I} - \underline{H})$$

$$\text{var}(e_i) = \sigma^2(f - \underbrace{h_{ii}}_0) \quad \text{cov}(e_i, e_j) = -\sigma^2 h_{ij}$$

$a_{ii}$  = i<sup>th</sup> diagonal term of A

$$h_{ii} = \left[ \underline{X} \left( \underline{X}^T \underline{X} \right)^{-1} \underline{X}^T \right]_{ii} = \underline{x_i} \left( \underline{X}^T \underline{X} \right)^{-1} \underline{x_i}^T$$

## Exercise: verify

With obs'n  
vector

$\text{def'n} \equiv h_{ii} = \text{leverage obs'n } i$

$$h_{ij} = \underline{x}_i (\underline{X}^T \underline{X})^{-1} \underline{x}_j^T \rightarrow \text{verify}!!$$

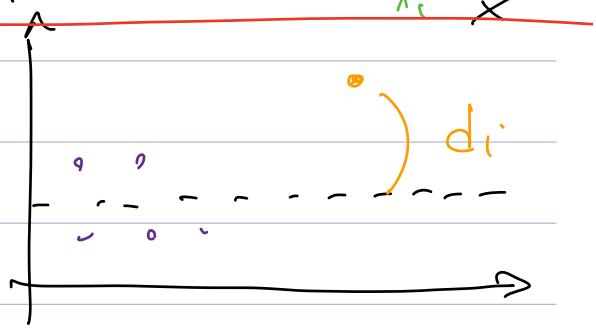
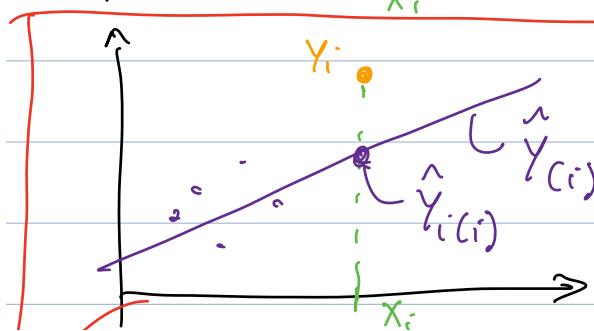
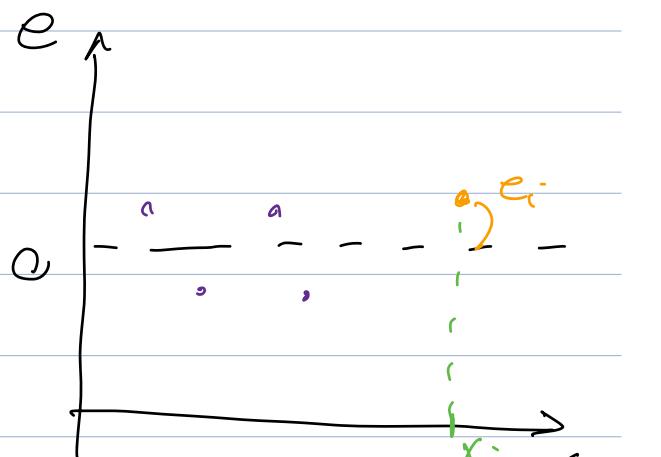
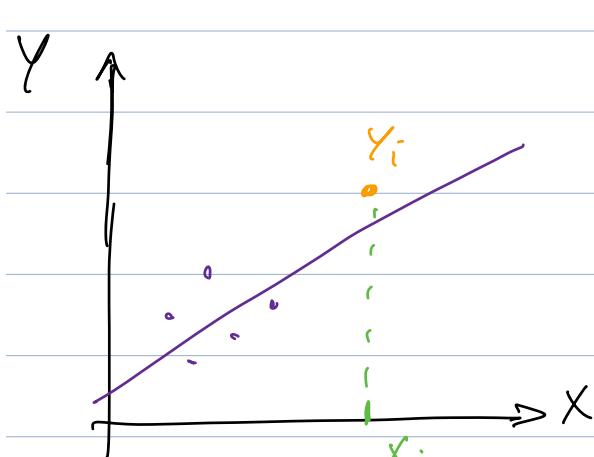
Note:  $x_i$  here is  $i$ th row of  $X$ , in row form! i.e.  $x_i \in \mathbb{R}^{d \times p}$

- To correctly studentize residuals we compute:

$$r_i = \frac{e_i}{s(e_i)} = \frac{e_i}{\sqrt{MSE(1-h_{ii})}}$$

"internally" studentized residuals

WARNING: no consistent naming of standardized, studentized residuals, etc, especially in R! (always check the documentation)



↳

Idea: let's forget about  $y_i$ !

$d_i$  = deleted residual

$$d_i = y_i - \hat{y}_{i(i)} \quad \text{removed obs'n } (i)$$

↳ pred'n at point  $x_i$

Fact:  $d_i = \frac{e_i}{1 - h_{ii}}$  (C.A. trick!)

! can compute this w/o fitting n different regressions

Proof sketch:

$$\text{Fact: } (\underline{\underline{x}}_{(i)})^\top (\underline{\underline{x}}_{(i)})^{-1} = \underbrace{(\underline{\underline{x}}^\top \underline{\underline{x}})^{-1}}_{p \times p} + \underbrace{(\underline{\underline{x}}^\top \underline{\underline{x}})^{-1} \underline{\underline{x}}_{(i)}^\top \underline{\underline{x}}_{(i)} (\underline{\underline{x}}^\top \underline{\underline{x}})^{-1}}_{p \times p \quad p \times 1 \quad 1 \times p \quad p \times p}$$

↳ design matrix  
w/  $i^{\text{th}}$  row removed

From this we get:

$$b_{(i)} = \underline{b} - (\underline{x}^T \underline{x})^{-1} \frac{\underline{x}_i^T \underline{e}_i}{1 - h_{ii}} \quad (\text{no refitting!})$$

$$d_i = \underline{y}_i - \hat{\underline{y}}_{(i)} = \underline{y}_i - \underline{x}_i \underline{b}_{(i)} =$$

$$\underline{y}_i - \underline{x}_i \underline{b} + \underline{x}_i (\underline{x}^T \underline{x})^{-1} \frac{\underline{x}_i^T \underline{e}_i}{1 - h_{ii}} =$$

$$\underline{e}_i + \frac{h_{ii} \underline{e}_i}{1 - h_{ii}} = \frac{\underline{e}_i (1 - h_{ii}) + h_{ii} \underline{e}_i}{1 - h_{ii}} =$$

$$\frac{\underline{e}_i}{1 - h_{ii}} \Rightarrow d_i = \frac{\underline{e}_i}{1 - h_{ii}}$$

when is  $d_i$  large?

both

- $\underline{e}_i$  is large
- $h_{ii} \approx 1$
- $\underline{e}_i$  is large,  $h_{ii} \approx 0$

$$S(d_i) = \underbrace{MSE_{(i)}}_{\text{i-th case removed}} \left[ 1 + \underbrace{\underline{x}_i}_{\text{i-th case}} \underbrace{(\underline{X}_{(i)})^\top (\underline{X}_{(i)})^{-1} \underline{x}_i}_{\substack{\text{design} \\ \text{matrix w/} \\ \text{i-th case removed}}} \right]$$

$$= \frac{MSE_{(i)}}{1 - h_{ii}}$$

Fact:  $\underbrace{(n-p)MSE}_{\text{SSE}} = (n-p-1) MSE_{(i)} + \frac{e_i^2}{1-h_{ii}}$

$$\Rightarrow SSE = (n-p-1) MSE_{(i)} + \frac{e_i^2}{1-h_{ii}} \Rightarrow$$

$$\frac{SSE (1-h_{ii}) - e_i^2}{n-p-1} = \boxed{MSE_{(i)} (1-h_{ii})}$$

studentize  $d_i$ :

$$t_i = \frac{d_i}{s(d_i)} = \frac{\overline{e_i}_{\bar{x} - h_{ii}}}{\sqrt{\frac{MSE_{(i)}}{1 - h_{ii}}}} = \frac{e_i}{\sqrt{MSE_{(i)}(1 - h_{ii})}}$$

→ w/o refitting the regression line!

$$\rightarrow e_i = e_i \sqrt{\frac{n-p-1}{SSE_{(\bar{x} - h_{ii})} - e_i^2}}$$

$t_i$  = externally studentized residuals

IF  $MSE_{(i)} \approx MSE$  then  $r_i \approx t_i$   
ISR ESR

$t_i \sim t(n-1-p)$  # predictors  
# obsn's

→ Test for outliers:

$$\text{Outlier } |t_{ij}| > t(1 - \frac{\alpha}{2n}, n-p-1)$$

↳ Bonferroni correction

= simultaneous test correction  
(conservative!)

but it works!

¶ Outlying obsn's wrt X:

e.g. (2)

related to Mahalanobis distance

$\underline{H}$  = distance measure in the space of pred's

Properties:

•  $0 \leq h_{ii} \leq 1$  (w/o proof)

•  $\sum_{j=1}^n h_{jj} = p$

$$\sum_{j=1} h_{jj} = \text{trace}(\underline{H}) = \text{tr}(\underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T) =$$

$$\text{tr}(A\underline{B}) = \text{tr}(\underline{B} A)$$

$$= \text{tr}((\underline{X}^T \underline{X}) (\underline{X}^T \underline{X})^{-1}) = \text{tr}(I_p) = p$$

$h_{ii} \approx$  distance  $\underline{x}_i$  to the center of the data in predictor space

$$\hat{Y}_i = (\underline{H}\underline{Y})_i = h_{i1} Y_1 + h_{i2} Y_2 + \dots + h_{ir} Y_r + \dots + h_{in} Y_n$$

If  $h_{ii} = 1 \Rightarrow \text{var}(e_i) = 0 \Rightarrow e_i = 0 \Rightarrow \hat{Y}_i = \hat{y}_i$   
 ↳ highest possible fit

$h_{ii}$  = leverage, "influence" of  $i^{th}$  obs'n on fitted values

Rule of thumb:  
 (to detect outliers  
 wrt  $X$ )

$$h_{ii} > \frac{2p}{n}$$

"2x the avg leverage"

$$\frac{\sum h_{ii}}{n} = \frac{p}{n}$$

$h_{ii} > 0.5$  = "high leverage"

NOT

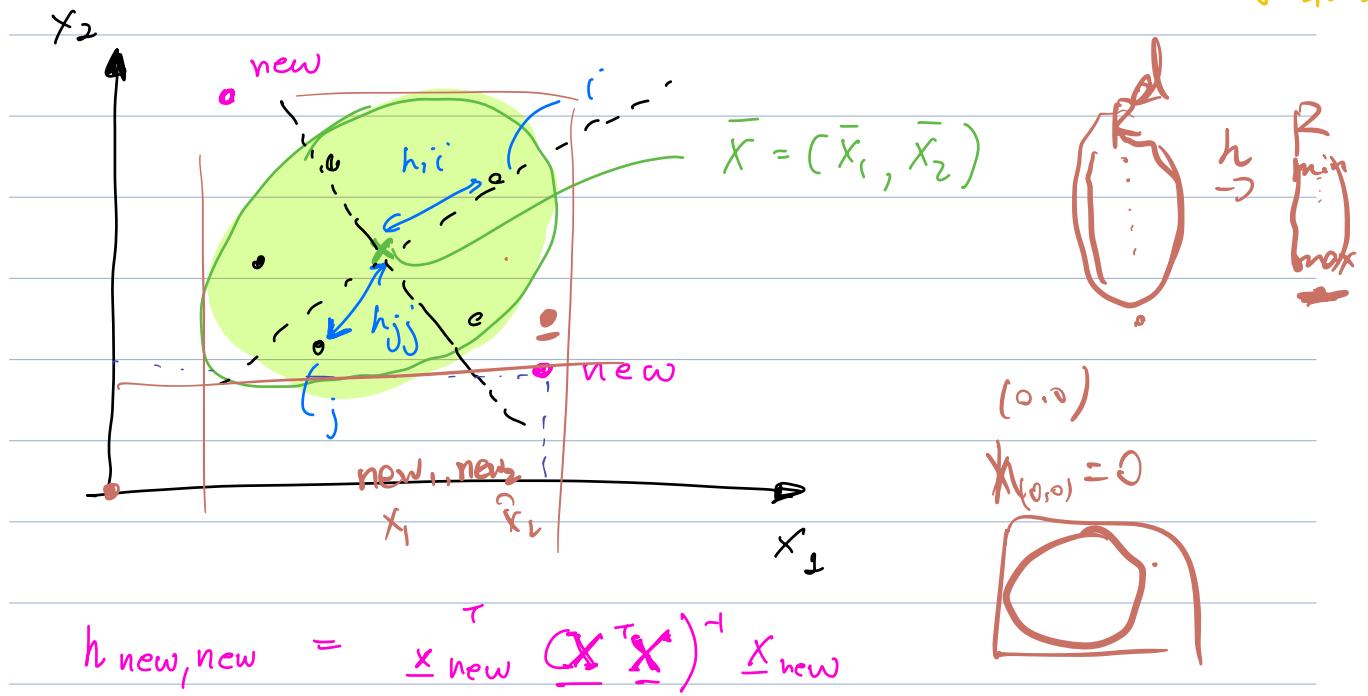
$0.2 < h_{ii} < 0.5$  = "moderate leverage"

TESTED!

→ hidden extrapolations

$$B\text{-有1d 的时候 } h_{ii} = \frac{1}{n} + \frac{1}{S^2} (x_i - \bar{x})^2$$

可以看距离。



$\min(h_{ii}) \leq h_{new,new} \leq \max(h_{ii}) \Rightarrow$  no hidden extrapolation!

MT2 stops here!

① Q3(b)

② Q2(ii).

③ hidden extrapolation vs. outliers.