

## INTRODUCTION

Regression modelling: rel.

$$Y = f(x_1, x_2, \dots, x_p)$$

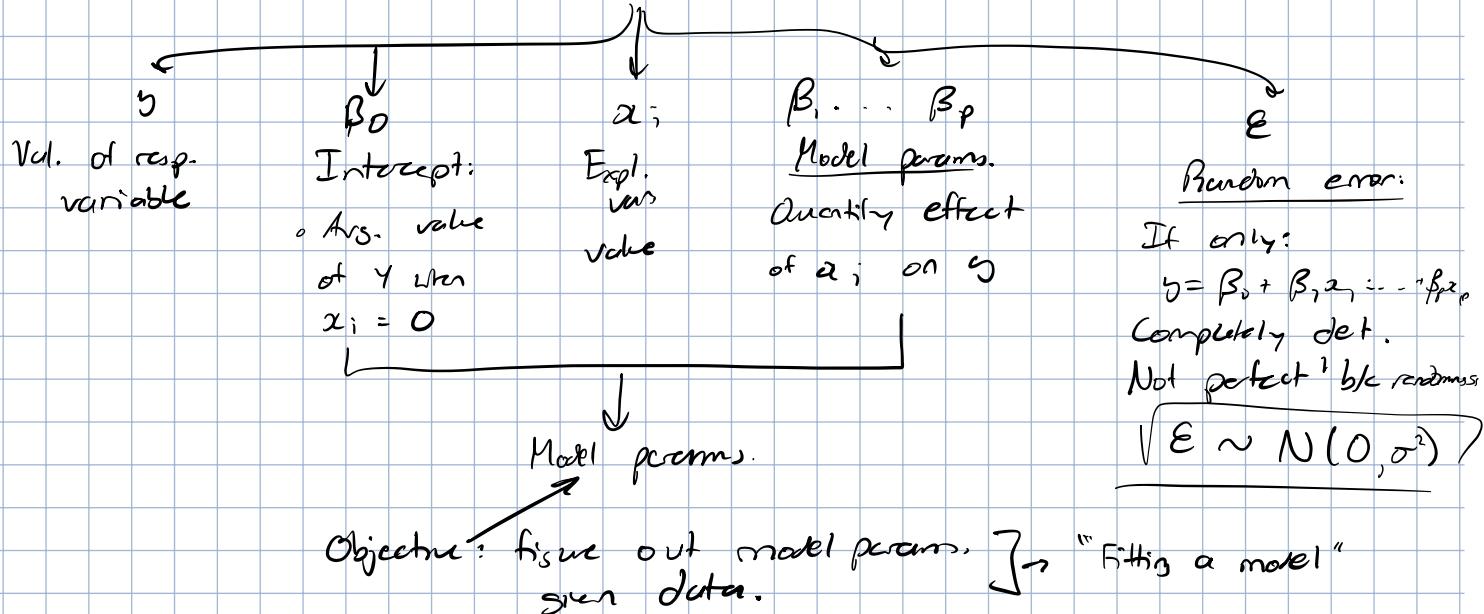
↑  
response var.  
dep. var.  
out. var.

explan. var.  
predictors  
features  
covariates.

Subset: linear RRM.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon \Rightarrow \text{Linear function}$$

Parts of LRM



Simplification: simple linear regression:

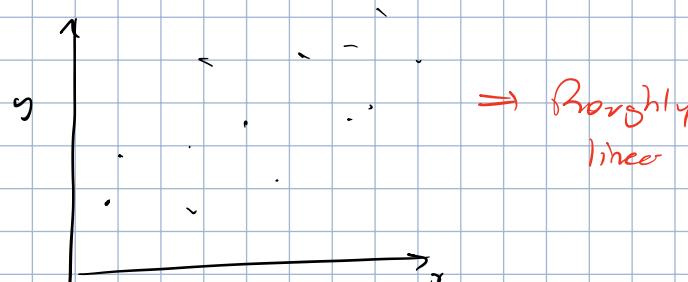
$$y = \beta_0 + \beta_1 x_1 + \epsilon \Rightarrow 1 \text{ expl. var.}$$

## SIMPLE LINEAR MODELS

Linearity

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

Check if linear rel. exists.  $\Rightarrow$  scatterplot



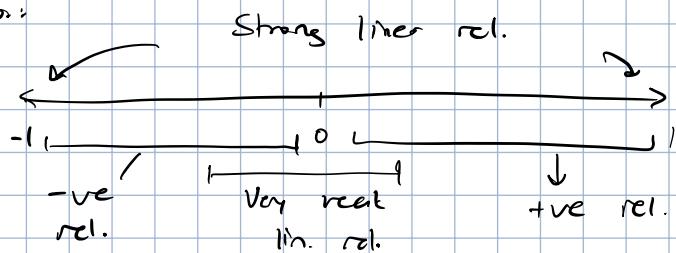
Check strength of linear rel.  $\Rightarrow$  Pearson's corr. coef.

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

Estimate  $\Rightarrow$  sample corr. coef.

$$r = \frac{S_{xy}}{\sqrt{S_{xx} \cdot S_{yy}}} \quad \begin{aligned} S_{xy} &= \sum_{i=0}^n (x_i - \bar{x})(y_i - \bar{y}) \\ S_{xx} &= \sum_{i=0}^n (x_i - \bar{x})^2 \\ S_{yy} &= \sum_{i=0}^n (y_i - \bar{y})^2 \end{aligned}$$

Interpretation:



Weaknesses:

1. Corr. coef. doesn't work in non-linear data



2. Cannot be used for prediction.

## Model Formulation

$$y = \beta_0 + \beta_1 x_i + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

Data:

$$(x_i, y_i), i = 1, \dots, n$$

Expect change / unit increase of  $x$ ,

Each value:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \Rightarrow \text{Fitted val.}$$

$$y = \beta_0 + \beta_1 x_i + \epsilon_i \Rightarrow \text{Actual val.}$$

B/c each  $y$  is  $\beta_0 + \beta_1 x_i + \epsilon_i \Rightarrow y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

$\hookrightarrow$  Proof:

$$E(y_i) = E \left[ \underbrace{\beta_0 + \beta_1 x_i}_{\text{constant}} + \epsilon_i \right]$$

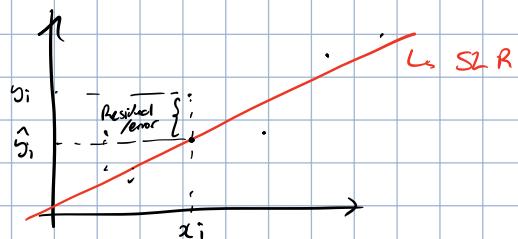
$$= \beta_0 + \beta_1 x_i + E[\epsilon_i]^0$$

$$= \beta_0 + \beta_1 x_i$$

$$\text{Var}(y_i) = \text{Var} \left[ \underbrace{\beta_0 + \beta_1 x_i}_{\text{constant}} + \epsilon_i \right]$$

$$= \sigma^2$$

## Method #1 of Estimating $\beta_0 + \beta_1$ : Least Squares Estimation



"Best" line has least amount of error

Measure of error for all values:

$$\text{Sum of squares error} = \sum_{i=0}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

Get rid of -ve + weights higher error more.

Obj: minimize SSE to get least square estimators (LSEs)  $\hat{\beta}_0, \hat{\beta}_1$

$$\underset{\beta_0, \beta_1}{\text{arg min}} \sum_{i=0}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

To minimize:

① Derivation:

$$\frac{\partial}{\partial \beta_0} \text{SSE} = -2 \sum (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial}{\partial \beta_1} \text{SSE} = -2 \sum (y_i - \beta_0 - \beta_1 x_i) (x_i)$$

② Zero:

$$a) \sum (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$b) \sum (y_i - \beta_0 - \beta_1 x_i) (x_i) = 0$$

③ Solve for  $\beta_0 + \beta_1$

Start w) a)  $\Rightarrow$  split out sum

Solve for  $\beta_0$

Sub in b)  $\Rightarrow$  solve for  $\beta_1$

$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

Remember:

$$S_{xx} = \sum (x_i - \bar{x})^2$$

$$\Rightarrow \sum (x_i - \bar{x}) = \sum x_i - n \bar{x} = \sum x_i - n \cdot \frac{\sum x_i}{n}$$

$$\begin{aligned}
 &= \sum x_i (x_i - \bar{x}) - \bar{x} \underbrace{\sum (x_i - \bar{x})}_0 \\
 &= \sum x_i^2 - \bar{x} \underbrace{\sum x_i}_n, \quad \sum x_i = n \bar{x} \\
 &= \sum x_i - n \bar{x}^2
 \end{aligned}
 \quad = \sum x_i - \sum \bar{x} = 0$$

$$\begin{aligned}
 S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \\
 &= \sum x_i y_i - \bar{x} \sum y_i - \bar{y} \underbrace{\sum (x_i - \bar{x})}_0 \\
 &= \sum x_i y_i - \bar{x} \cdot n \bar{y} \\
 &= \sum x_i y_i - \bar{x} \bar{y}
 \end{aligned}$$

$$\therefore \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

Note: Given data  $\Rightarrow$  find  $\hat{\beta}_0 + \hat{\beta}_1$

Ex. 11  $\sum x_i = a$ ,  $\sum y_i = b$ ,  $\sum x_i^2 = c$ ,  $\sum x_i y_i = d$ ,  $n$

Find SLR:

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{d - n \cdot \frac{a}{n} \cdot \frac{b}{n}}{c - n \cdot \left(\frac{a}{n}\right)^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \cdot \bar{x} = \frac{b}{n} - \hat{\beta}_1 \cdot \frac{a}{n}$$

### Residuals

$$r_i = y_i - \hat{y}_i = y_i - \beta_0 - \beta_1 x_i$$

Properties:

LSE:

- $\sum (y_i - \beta_0 - \beta_1 x_i) = 0$
- $\sum (y_i - \beta_0 - \beta_1 x_i) x_i = 0$

- $\sum r_i = 0 \Rightarrow$  From a)
- $\sum r_i x_i = 0 \Rightarrow$  From b)
- $\sum r_i \hat{y}_i = \sum r_i (\hat{\beta}_0 + \hat{\beta}_1 x_i)$   
 $= \hat{\beta}_0 \sum r_i + \hat{\beta}_1 \sum r_i x_i$   
 $= 0$

} Model checking

## Estimating $\sigma^2$ :

To estimate pop. sample size  $y_i \sim N(\mu, \sigma^2)$

$$\text{Sample var} = \frac{1}{n-1} \sum (y_i - \bar{y})^2$$

Losing 1 d.f. by using estimator  $\bar{y}$

We can do same thing for SLR:

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\begin{aligned} \therefore S^2 &= \frac{1}{n-2} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \frac{1}{n-2} \sum r_i^2 \end{aligned}$$

Unbiased estimator of  $\sigma^2$

## Method #2 of Estimating $\beta_0 + \beta_1$ : Maximum Likelihood

Theory:

① Maximum likelihood func.

$$L(\theta) = f(y; \theta) \Rightarrow \text{How likely are we seeing observed data}$$

∴ Find  $\theta$  to maximize prob.

② Log likelihood:

$$\lambda(\theta) = \log(L(\theta))$$

③ Derivative:

$$S(\theta) = \frac{\partial}{\partial \theta} L(\theta) \Rightarrow \text{Solv func.}$$

④ Set to 0 + solve for  $\theta$ :

$$S(\theta) = 0$$

SLR:

① Likelihood function:

$$\text{Each } y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$L(\theta) = P(y_1 = \dots) \cdot P(y_0 = \dots) \dots$$

$$\theta = \{\beta_0, \beta_1, \sigma^2\} = \prod \underbrace{P(y_i = y_i)}_{\text{normal}}$$

$$= \prod (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2} \right\}$$

② Log likelihood:

$$\ell(\theta) = \sum \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2} \right]$$

③ Score function:

$$① \frac{\partial}{\partial \beta_0} = \frac{1}{\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$② \frac{\partial}{\partial \beta_1} = \frac{1}{\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i) (x_i) = 0$$

$$③ \frac{\partial}{\partial \sigma^2} = -\frac{1}{\sigma^2} + \frac{1}{2\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i)^2 = 0$$

④ Soln:

① + ②  $\Rightarrow$  sum as LSE

$$\therefore \text{MLE} = \text{LSE}$$

$$③ \sigma^2 = \frac{1}{n} \sum r_i^2$$

$\rightarrow$  This is diff.

The other is unbiased  $\Rightarrow$  This MLE is biased (don't use)

Properties of LSE

Note: 1. estimator vs. estimate

$$\begin{array}{c} \text{r.v.} \\ \tilde{\beta}_0, \tilde{\beta}_1 \xrightarrow{f(y_i)} \end{array} \xrightarrow{\text{Observed value}} N$$

2. Assumption:  $\epsilon_i \sim N(0, \sigma^2)$

Property #1: LSE are unbiased

$$a) E(\tilde{\beta}_1) = \beta_1 \leftarrow \text{true value}$$

$$\tilde{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} \rightarrow \text{contains}$$

$$E(\tilde{\beta}_1) = \frac{\sum (x_i - \bar{x}) E(y_i)}{\sum (x_i - \bar{x})^2}$$

$$b) E(\tilde{\beta}_0) = \beta_0 \quad \text{r.v.}$$

$$\begin{aligned} \tilde{\beta}_0 &= E(\bar{y} - \tilde{\beta}_1 \bar{x}) \\ &= \frac{1}{n} \sum E(y_i) - \bar{x} E(\tilde{\beta}_1) \\ &= \beta_0 + \beta_1 \bar{x} - \bar{x} \beta_1 \\ &= \beta_0 \end{aligned}$$

$$\begin{aligned}
 &= \frac{\sum (x_i - \bar{x})(\beta_0 + \beta_1 x_i)}{\sum (x_i - \bar{x})^2} \\
 &= \beta_1 \frac{\sum (x_i - \bar{x})x_i}{\sum (x_i - \bar{x})^2} \quad \text{with } \sum (x_i - \bar{x})^2 \\
 &= \beta_1
 \end{aligned}$$

Property # 2: Variances of  $\hat{\beta}_0$  and  $\hat{\beta}_1$

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

$$\begin{aligned}
 a) \quad \text{Var}(\hat{\beta}_1) &= \text{Var}\left(\frac{\sum (x_i - \bar{x})y_i}{s_{xx}}\right) \\
 &= \left(\frac{1}{s_{xx}}\right)^2 \text{Var}\left(\sum (x_i - \bar{x})y_i\right) \\
 &= \left(\frac{1}{s_{xx}}\right)^2 \underbrace{\sum (x_i - \bar{x})^2}_{s_{xx}} \cdot \text{Var}(y_i) \\
 &= \frac{\sigma^2}{s_{xx}}
 \end{aligned}$$

$$\text{Var}(\sum y_i) = \sum \text{Var}(y_i) + \sum_{i \neq j} \text{Cov}(y_i, y_j)$$

Assume indep. samples  $\Rightarrow$  0 if each sample is indep.

$$\begin{aligned}
 b) \quad \text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) \\
 &= \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2\bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1)
 \end{aligned}$$

$$\begin{aligned}
 \text{Cov}(\bar{y}, \hat{\beta}_1) &= \text{Cov}\left(\frac{1}{n} \sum y_i, \frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x})y_i\right) \\
 &= 0 \quad \text{Cov of sure variable!}
 \end{aligned}$$

$$\begin{aligned}
 \therefore \text{Var}(\hat{\beta}_0) &= \frac{\sigma^2}{n} + \bar{x}^2 \cdot \frac{\sigma^2}{s_{xx}} \\
 &= \frac{\sigma^2 \sum x_i^2}{n s_{xx}}
 \end{aligned}$$

Property # 3: Covariance of  $\hat{\beta}_0$  +  $\hat{\beta}_1$

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

## Property # 4: Distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

$y_i$  are iid  $N$  &  $\hat{\beta}_0 + \hat{\beta}_1 x_i$  or linear combo of  $y_i$

$$\therefore \hat{\beta}_0 \sim N(\beta_0, \frac{\sigma^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}), \quad \hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2})$$

Issue: all params can be calculated except  $\sigma^2$

↳ Recall:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum r_i^2 = s^2 \Rightarrow E(s^2) = \sigma^2$$

↳ MSE: mean squared error

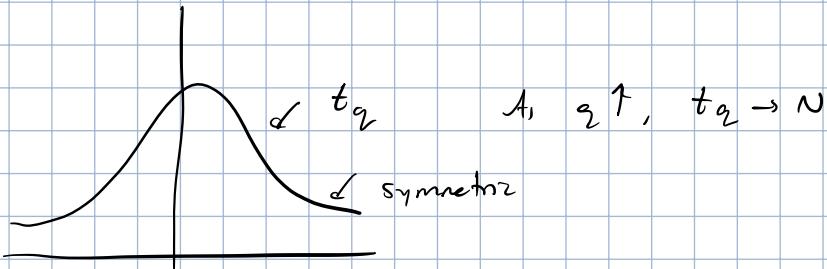
Since this is calculable  $\Rightarrow$  use this in our distn (sub  $\sigma^2$  for  $s^2$ )

Standard error:

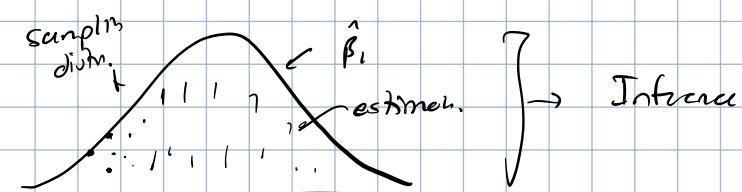
$$\begin{aligned} \text{se}(\hat{\beta}_0) &= \sqrt{\text{Var}(\hat{\beta}_0)} & \text{standard dev. using } \sigma^2 \text{ if intd} \\ &= \sqrt{\frac{s^2 \sum x_i^2}{\sum (x_i - \bar{x})^2}} & \text{se}(\hat{\beta}_1) = \sqrt{\text{Var}(\hat{\beta}_1)} \\ & & = \sqrt{\frac{s^2}{S_{xx}}} \end{aligned}$$

$$\text{If we use S.E.} \Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)} \sim t_{n-2}, \quad \frac{\hat{\beta}_0 - \beta_0}{\text{se}(\hat{\beta}_0)} \sim t_{n-2}$$

Recall:



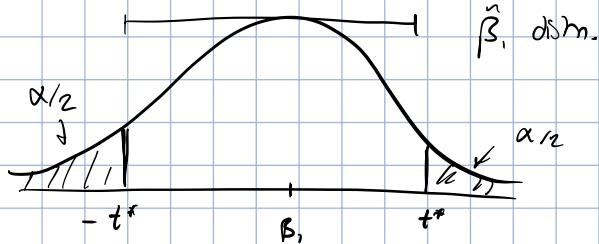
## Inference for Regression Parameters.



#1: C.I

$\alpha$  locul

$1 - \alpha$



$$\therefore P\left(-t^* < \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} < t^*\right) = 1 - \alpha$$

$$P\left(\hat{\beta}_1 - t^* se(\hat{\beta}_1) > \beta_1 > \hat{\beta}_1 + t^* se(\hat{\beta}_1)\right) = 1 - \alpha$$

↑  
true pop param.

$$C.I. = \hat{\beta}_1 \pm t^* se(\hat{\beta}_1)$$

$$\hat{\beta}_1 \pm t^* se(\hat{\beta}_1)$$

} Basis of pop param in 95%  
confidence.

}  $t$ -table via  $\alpha/2$  + d.f. ( $n-2$ )

If estimate  $\notin$  C.I.  $\Rightarrow$  can reject  $H_0$

## #2: Hypothesis Testing

Obj.: true value of param  $= \beta_1 / \in$  some estimate.

① Hypothesis:

$$\underbrace{H_0}_{\text{state quo}} \quad \text{vs.} \quad \underbrace{H_a}_{\text{Not state quo}}$$

You assume  $H_0$  is T  $\Rightarrow$  if test gets crazy result  $\Rightarrow$  reject  $H_0$

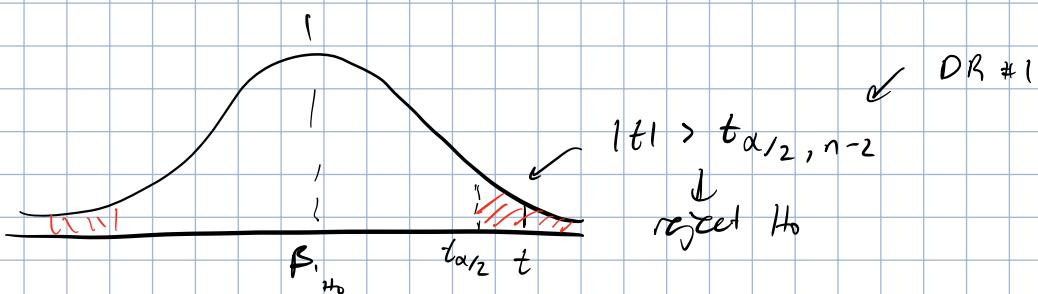
② Decision rule: what will make you reject  $H_0$

③ Test statistic  $\Rightarrow$  decision.

Test stat:  $t$ -stat

$$\frac{\hat{\beta}_1 - \beta_{1H_0}}{se(\hat{\beta}_1)} \sim t_{n-2} \quad (\text{if } H_0 \text{ is true})$$

If your data is off  $\Rightarrow$   $t$ -stat is quite far



Other decision rule: p-val

$P(|T| > t)$  = p-val of seeing data if  $H_0$  is true.

If  $p\text{-val} < \alpha \Rightarrow$  reject  $H_0$ . (DR  $\geq 2$ )

Common test:

o  $\hat{\beta}_1$  test of significance: is there a linear rel.

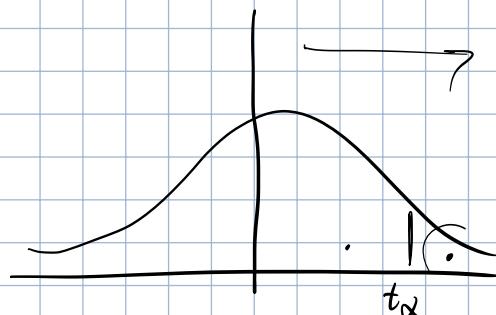
$$H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0$$

Test stat:  $\frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} \Rightarrow t\text{-ratio}$

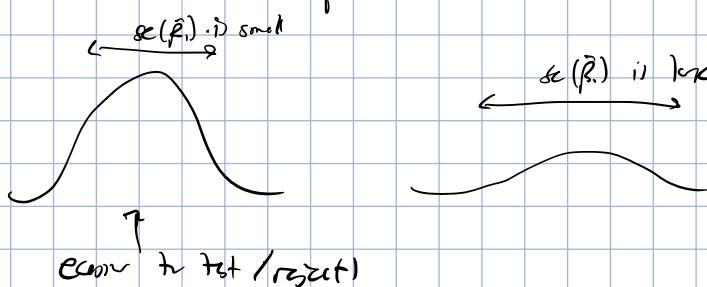
Reject in  $H_0$  & acceptance of  $H_a$ .

One-sided testing:

$$H_0: \beta_1 \geq 0, \quad H_a: \beta_1 < 0$$



This is based on  $\text{se}(\hat{\beta}_1) \propto 1/s_{\alpha/2}$



If  $s_{\alpha/2} \uparrow \Rightarrow \text{se} \downarrow \Rightarrow$   
reject / accept easier!  
Sporadically out  $\alpha$ -vals  
in LR

Ex: 1) Sign. rel. b/w prop. age & bending strength

① Hypothesis:

$$H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0$$

② Decision rule:

$$\text{If } \left| \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} \right| > t_{\alpha/2, n-2} \Rightarrow \text{reject } H_0$$

③ Math.

$$\frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} = -12.86, \quad t_{\alpha/2, n-2} = 2.101.$$

$$\text{q} t(0.975, n-2)$$

∴ Reject  $H_0$  = Interpret.

## Estimation of Mean Response

Mean response vari. vol. at  $x = x_p$

① Random var:

$$\mu = \beta_0 + \beta_1 x_p$$

② Estimate

$$\hat{\mu} = \hat{\beta}_0 + \hat{\beta}_1 x_p$$

③ Unbiased?

$$\begin{aligned} E(\hat{\mu}) &= E(\hat{\beta}_0) + x_p E(\hat{\beta}_1) \\ &= \beta_0 + x_p \beta_1 \end{aligned}$$

④ Varience?

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_p) \\ &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_p) \\ &= \text{Var}(\bar{y} + \hat{\beta}_1 (\bar{x} + x_p)) \\ &= \text{Var}(\bar{y}) + 2(\bar{x} + x_p) \text{Cov}(\bar{y}, \hat{\beta}_1) + (\bar{x} + x_p)^2 \text{Var}(\hat{\beta}_1) \end{aligned}$$

$$\begin{aligned} \text{Var}\left(\frac{1}{n} \sum y_i\right) &= \frac{1}{n^2} \cdot \sum \text{Var}(y_i) \\ &= \frac{1}{n^2} \cdot n \sigma^2 \\ &= \sigma^2 / n \end{aligned}$$

$$\begin{aligned} \text{Var}\left(\frac{1}{n} \sum y_i\right) &= \frac{1}{n^2} \sum \text{Var}(y_i) \quad \text{if } i \neq j \\ &= \frac{1}{n^2} \sum \sigma^2 \\ &= \sigma^2 / n \end{aligned}$$

⑤ Sampling distn.

$$\frac{\hat{\mu} - \mu}{\sqrt{\text{Var}(\hat{\mu})}} \sim N(0, 1)$$

↓

$$\frac{\hat{\mu} - \mu}{\text{se}(\hat{\mu})} \sim t_{n-2}$$

$$s^2 = \frac{1}{n-2} \sum r_i^2$$

Run C.I. & inference tests?

Ex: Is avg. strength of obj. at  $x = 16 < 2000$ ?

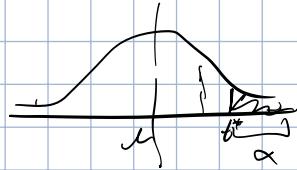
① Estimate:

$$\hat{\mu} = \hat{\beta}_0 + \hat{\beta}_1 (16) = 2031.865$$

② Hypothesis:

$$H_0: \mu \leq 2000$$

$$H_a: \mu > 2000$$



③ Decision rule:

If  $t > t_{\alpha, n-2} \Rightarrow$  reject  $H_0$

④ Maths

$$\frac{\bar{\mu} - 2000}{\text{se}(\bar{\mu})} = 1.463, \quad t_{\alpha, n-2} = 1.734$$

⑤ Conclusion:

Cannot reject  $H_0 \Rightarrow$  might be舞者.

Prediction

True value of  $y_p$  from a new subject  $x = x_p \Rightarrow$  not part of data

① R.V.

$$y_p = \beta_0 + \beta_1 x_p + \varepsilon \quad \text{future random error.}$$

② Estimation

$$\hat{y}_p = \hat{\beta}_0 + \hat{\beta}_1 x_p$$

③ Prediction:

$$y_p - \hat{y}_p \Rightarrow \text{No idea what } y_p \text{ is.}$$

Sampling distn.:

a) Mean:

$$\begin{aligned} E[y_p - \hat{y}_p] &= (\beta_0 + \beta_1 x_p) - (\hat{\beta}_0 + \hat{\beta}_1 x_p) \\ &= 0 \end{aligned}$$

b) Var:

$$\begin{aligned} \text{Var}[y_p - \hat{y}_p] &= \text{Var}[y_p - \hat{\beta}_0 - \hat{\beta}_1 x_p] \\ &= \text{Var}(\varepsilon) + \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_p) \\ &= \sigma^2 + \sigma^2 \left( V_n + \frac{(x_p - \bar{x})^2}{s_{xx}} \right) \\ &= \sigma^2 \left( 1 + V_n + \frac{(x_p - \bar{x})^2}{s_{xx}} \right) \end{aligned}$$

i) Greater than expected response.

ii) Dep. on  $(x_p - \bar{x})$

c)

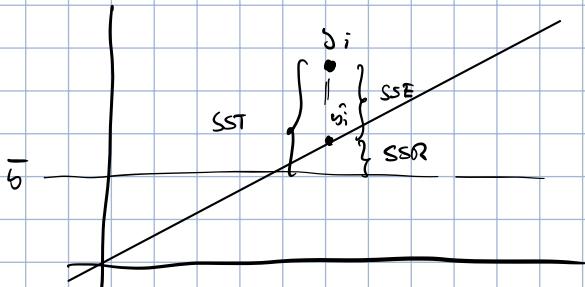
$$\frac{y_p - \hat{y}_p}{\text{se}(y_p - \hat{y}_p)} \sim t_{n-2} \Rightarrow \text{C.I. / tests.}$$

## ANOVA

$$SST = SSE + SSR$$

Total sum of squares:  $\sum (y_i - \bar{y})^2$ 
Error sum of squares:  $\sum (y_i - \hat{y}_i)^2 = \sum e_i^2$ 
Regression sum of squares:  $\sum (\hat{y}_i - \bar{y})^2$

Variance of  $y$ :



Proof:

$$\begin{aligned}
 SST &= \sum (y_i - \bar{y})^2 \\
 &= \sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\
 &= \sum (y_i - \hat{y}_i) + 2 \sum (\hat{y}_i - \bar{y})(\hat{y}_i - \bar{y}) + \sum (\hat{y}_i - \bar{y})^2 \\
 &= \text{SSE} + \underbrace{2 \left( \sum n \hat{y}_i^2 - 2 \bar{y} \sum \hat{y}_i^2 \right)}_0 + \text{SSR} \\
 &= SSE + SSR
 \end{aligned}$$

Degrees of freedom:

◦  $SST = n-1$

◦  $SSE = n-2$

◦  $SSR = 1$

Mean squares:

◦  $MSST = SST/d.f. = SST/n-1$

◦  $MSE = \frac{SSE/d.f.}{n-2} = \frac{SSE/n-2}{n-2} \rightarrow \frac{1}{n-2} \sum n_i^2 \Rightarrow \text{unbiased est. of } \sigma^2$

◦  $MSR = SSR/d.f. = SSR$

MSR:

◦  $E[MSR] = \sigma^2 + \beta_1^2 s_{xx}$

◦ Proof:

$$\begin{aligned}
 MSR &= \sum (\hat{y}_i - \bar{y})^2 \\
 &= \hat{\beta}_1^2 \sum (x_i - \bar{x})^2 \\
 &= \hat{\beta}_1^2 \cdot s_{xx}
 \end{aligned}$$

$$\hat{y}_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x})$$

$$\begin{aligned}
 E(MSR) &= E(\hat{\beta}_i^2 \cdot S_{xx}) \\
 &= S_{xx} \cdot (\text{Var}(\hat{\beta}_i) + E(\hat{\beta}_i)^2) \\
 &= \sigma^2 + \hat{\beta}_i^2 S_{xx}
 \end{aligned}$$

Inference:

①  $MSR \gg MSE \Rightarrow \beta_i \neq 0$  (error can be explained better w/ regression line)

②  $\frac{MSR}{MSE} \sim F_{1, n-2}$

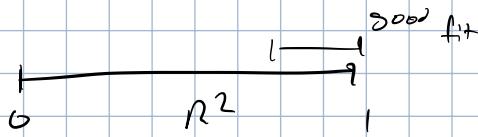
L, F-test

$$H_0: \beta_i = 0, \quad H_a: \beta_i \neq 0$$

} High  $\rightarrow$  more error expl. by SR

Coefficient of determination:

$$R^2 = \frac{SSR}{SST} \Rightarrow \text{prop. of error expl. by reg. line.}$$



## MULTIPLE LINEAR REGRESSION

Random vectors & Multivariate Normal Distribution

Random vector:

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \rightarrow \text{r.v.}$$

◦ Mean:

$$E[Y] = \begin{bmatrix} E(y_1) \\ \vdots \\ E(y_n) \end{bmatrix} = \mu$$

◦ Variance-covariance matrix:

$$\begin{aligned}
 \text{Var}(Y) &= \begin{bmatrix} \text{Var}(y_1) & \text{Cov}(y_1, y_2) & \dots & \dots \text{Cov}(y_1, y_n) \\ \vdots & \text{Var}(y_2) & & \\ \vdots & & \ddots & \\ \text{Cov}(y_n, y_1) & & & \text{Var}(y_n) \end{bmatrix} \rightarrow \text{Symmetric to lower triangle.}
 \end{aligned}$$

Positive definite

$$= \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 & \dots & \sigma_{1n}^2 \\ & \ddots & & \\ & & \ddots & \\ & & & \sigma_n^2 \end{bmatrix} \Rightarrow \Sigma_{n \times n}$$

If  $Y$  is indep.  $\Rightarrow \text{diag}\{\sigma_1^2, \dots, \sigma_n^2\}$

Properties:

① Linear combination on expectation:

$$E[AY + b] = A E[Y] + b$$

② Variance of vector + constant

$$\text{Var}[Y + b] = \text{Var}[Y]$$

③ Linear comb. on variance:

$$\text{Var}[AY + b] = A \text{Var}[Y] A' \Rightarrow \text{Square a matrix/vector: } A A'$$

Multivariate normal distn.

$$Y \sim \text{MVN}(\mu, \Sigma) \Leftrightarrow f(y_1, \dots, y_n) = (2\pi)^{-n/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu))$$

Properties:

① Linear comb.

$$Y \sim \text{MVN}(\mu, \Sigma) \Leftrightarrow AY + b \sim \text{MVN}(A\mu + b, A\Sigma A')$$

② Indiv. normality

$$Y \sim \text{MVN}(\mu, \Sigma) \Rightarrow y_i \text{ in } Y \sim N(\mu_i, \sigma_i^2)$$

③ Partitioning:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \Rightarrow \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$\hookrightarrow y_1, y_2 \sim \text{MVN}(\mu_i, \Sigma_i)$$

④ Independence of ind. vars:

$$y_1, \dots, y_n \stackrel{\text{iid}}{\sim} N(\mu_i, \sigma_i^2) \Rightarrow Y \sim \text{MVN}(\mu, \Sigma) \curvearrowright \text{diag}\{\sigma_1^2, \dots, \sigma_n^2\}$$

⑤ Variance & independence:

$$\text{Var}(Y) = \Sigma \text{ is a diag.} \Leftrightarrow \text{independence of } y_i$$

⑥ Independence of scaled MVN

$$Y \sim \text{MVN}(\mu, \Sigma), \quad V = AY, \quad W = BY$$

Not true for other distn.

$$\hookrightarrow V \perp W \Leftrightarrow A \Sigma B' = 0$$

## Multiple Linear Models

1 observation:  $y$  response,  $p$  explanatory

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

↑      ↑      ↑      ↑      ↑      ↗  
impact of  $i^{th}$  obs. of  $n^{th}$  explanatory var.

$\sim N(0, \sigma^2)$

Matrix:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \beta_0 \\ x_{11} \dots x_{1p} \\ x_{21} \dots x_{2p} \\ \vdots \\ x_{n1} \dots x_{np} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

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

↑  
design matrix

Interpretation:

$\beta_0$ : expected resp. val. if  $x_1, \dots, x_p = 0$

$\beta_i$  ( $i \neq 0$ ): change in resp. variable by 1 unit in  $x_i$  if all other vars are fixed.

Distrn:

Assumption:  $\epsilon_i \sim N(0, \sigma^2) \Rightarrow \epsilon \sim MVN(0, \Sigma) \xrightarrow{\text{diag } \{\sigma^2, \dots, \sigma^2\}}$

$\because Y$  is a liner combo of  $\epsilon \Rightarrow Y \sim MVN(X\beta, \sigma^2 I)$

## Parameter Estimation

Least squares:

$$\begin{aligned} \hat{\beta} &= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2 \\ &= \underset{\beta}{\operatorname{argmin}} \sum (Y - X\beta)' (Y - X\beta) \end{aligned}$$

Take derivatn:

$$\frac{\partial}{\partial \beta} \hat{\beta} = \frac{2}{2\beta} [YY' - Y'X\beta - X'\beta'Y + \beta'X'X\beta]$$

Aside:

$$Z = \alpha' \beta \Rightarrow \frac{\partial}{\partial \beta} Z = \alpha$$

$$Z = \beta' \alpha \Rightarrow \frac{\partial}{\partial \alpha} Z = \beta$$

$$Z = \beta' A \beta \Rightarrow \frac{\partial}{\partial \beta} Z = (A + A') \beta$$

$$\frac{\partial}{\partial \beta} \left[ Y' Y - \underbrace{Y' X \beta}_{\text{no } \beta} - \underbrace{X' \beta' Y}_{\text{no } \beta} + \beta' X' X \beta \right] = - (Y' X)' - (X' Y)' + (X' X) \beta$$

$$= - Y X' - X' Y + 2 X' X \beta$$

$$= - 2 X' Y + 2 X' X \beta$$

Set to 0:

$$- 2 X' Y + 2 X' X \hat{\beta} = 0$$

$$X' X \hat{\beta} = X' Y$$

$$\hat{\beta} = (X' Y) (X' X)^{-1}$$

Properties:

$$\textcircled{1} \quad E(\hat{\beta}) = \beta$$

$$\begin{aligned} E(\hat{\beta}) &= E[(X' Y) (X' X)^{-1}] \\ &= (X' X)^{-1} X' E[Y] \\ &= \underbrace{(X' X)^{-1}}_{+I} \underbrace{X' X \beta}_{= \beta} \\ &= \beta \end{aligned}$$

$$\textcircled{2} \quad \text{Var}(\hat{\beta}) = \sigma^2 (X' X)^{-1}$$

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \text{Var}[(X' Y) (X' X)^{-1}] \\ &= [(X' X)^{-1} X'] \text{Var}(Y) [(X' X)^{-1} X']' \\ &= ((X' X)^{-1} X') \text{Var}(Y) X (X' X)^{-1} \\ &\quad \downarrow \sigma^2 \\ &= (X' X)^{-1} \sigma^2 \end{aligned}$$

$$\textcircled{3} \quad \tilde{\beta} \sim \text{MVN}(\beta, \sigma^2 (X' X)^{-1}) \quad \text{bc } \hat{\beta} \text{ is a linear comb of } Y$$

$$\textcircled{4} \quad \hat{\beta}_j \sim N(\beta_j, \sigma^2 (X' X)^{-1}_{jj})$$

## Fitted Values & Residuals

$$SLR: \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$MLR: \hat{Y} = X \hat{\beta} = X \left[ (X'X)^{-1} X' Y \right]$$

$$= \underbrace{X (X'X)^{-1} X'}_H Y$$

$$\text{Hat matrix: } H = X (X'X)^{-1} X'$$

① Symmetric:

$$\text{Show: } H' = H$$

$$\begin{aligned} H' &= \left[ X (X'X)^{-1} X' \right]' \\ &= X \left[ X (X'X)^{-1} \right]' \\ &= X' (X'X)^{-1} X' \\ &= H \end{aligned}$$

② Idempotent:

$$\text{Show: } H H = H$$

$$\begin{aligned} H H &= \left( X \underbrace{(X'X)^{-1} X'} \right) \left( X \underbrace{(X'X)^{-1} X'} \right) \\ &= X (X'X)^{-1} X' \xrightarrow{\text{Cancel.}} \\ &= H \end{aligned}$$

③  $\hat{Y}$  is a L.T. of  $Y \Rightarrow \hat{Y} \sim MVN$

Residuals:

$$SLR: r_i = y_i - \hat{y}_i$$

$$\begin{aligned} MLR: r &= Y - \hat{Y} \\ &= Y - H Y \\ &= (I - H) Y \end{aligned}$$

①  $r$  is a L.T. of  $Y \Rightarrow r \sim MVN$

②  $E[r] = 0$

$$\text{Proof: } E[r] = E[(I - H) Y]$$

$$\begin{aligned}
&= (I - H) \times \beta \\
&= X\beta - HX\beta \\
&= X\beta - \underbrace{(X(X'X)^{-1}X')}_X X\beta \\
&= X\beta - X\beta \\
&= 0
\end{aligned}$$

$$\textcircled{3} \quad \text{Var}[r] = \sigma^2 (I - H)$$

$$\begin{aligned}
\text{Var}[r] &= \text{Var}((I - H)Y) \\
&= (I - H) \text{Var}(Y) (I - H)' \\
&= \sigma^2 (I - H)(I - H)' \quad I - H \text{ is symmetric.} \\
&= \sigma^2 (I - H) \quad \downarrow \text{Idempotent} \\
\therefore r &\sim MVN(0, \sigma^2(I - H))
\end{aligned}$$

$$\textcircled{4} \quad \sum_{i=1}^n r_i = 0$$

$$\textcircled{5} \quad \sum_{i=1}^n r_i x_{i1} = 0, \quad \sum_{i=1}^n r_i x_{i2} = 0, \quad \dots, \quad \sum_{i=1}^n r_i x_{ip} = 0$$

$$\textcircled{6} \quad \sum_{i=1}^n r_i \hat{y}_i = 0$$

Proofs of \textcircled{4}, \textcircled{5}, \textcircled{6}:

Note that  $X'r$  will give us \textcircled{4} & \textcircled{5}:

$$\begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}' \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix} = \begin{bmatrix} \sum r_i \\ \vdots \\ \sum x_{i1} r_i \end{bmatrix}$$

$$X'r = 0 \Rightarrow \text{done:}$$

$$\begin{aligned}
X'r &= X'(Y - \hat{Y}) \\
&= X'Y - X'\hat{Y} \\
&= X'Y - \underbrace{X'X(X'X)^{-1}X'Y}_X \\
&= X'Y - X'Y
\end{aligned}$$

$\approx 0$

## Estimation of $\sigma^2$ in M2R

$$\text{SLR: } \hat{MSE} = \frac{1}{n-2} \sum v_i^2$$

$$\text{MLR: } \frac{SSE}{n-(p+1)} \Rightarrow \text{unbiased estimator of } \sigma^2$$

## Inference in MLR

Recall:

- $\hat{\beta} \sim MVN(\beta, \sigma^2 (X'X)^{-1})$
- $\hat{\beta}_j \sim N(\beta_j, [\sigma^2 (X'X)^{-1}]_{jj})$

$$\therefore \frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 [\{X'X\}^{-1}]_{jj}}} \sim N(0, 1) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\sqrt{s^2 [\{X'X\}^{-1}]_{jj}}} \sim t_{n-(p+1)}$$

C.I.:

$$\hat{\beta}_j \sim t_{\alpha/2, n-(p+1)} \times se(\hat{\beta}_j)$$

Hypothesis:

$$|t| = \left| \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \right| > t_{\alpha/2, n-(p+1)}$$

Exact same mechanism.

Difference: interpretation

C.I. & test of  $\hat{\beta}_j$  is accounting for relation with other  $\hat{\beta}_j$  included in model.

## Inference for Linear Combos of Parameters

Constant vector:

$$C = (1, x_1, \dots, x_p)'$$

Estimated response:  $\hat{\mu}_c = C' \hat{\beta}$

Distrn:  $\hat{\mu}_c \sim MVN (\beta^T \beta, \sigma^2 \beta^T (\mathbf{X}' \mathbf{X})^{-1} \beta)$

$$\therefore \frac{\hat{\mu}_c - \mu_c}{\sqrt{\sigma^2 \beta^T (\mathbf{X}' \mathbf{X})^{-1} \beta}} \sim t_{n-(p+1)}$$

C.I. & t-s binned off ths.

## Prediction Inference.

True value:

$$y_p = \beta^T \beta + \epsilon_p$$

Estimate:

$$\hat{y}_p = \hat{\beta}^T \beta$$

Prediction error:  $y_p - \hat{y}_p$

Distrn?

$$\textcircled{1} \quad E[y_p - \hat{y}_p] = 0$$

$$\begin{aligned} E[\underbrace{\beta^T \beta + \epsilon_p - \hat{\beta}^T \beta}_{y_p} - \hat{\beta}^T \beta] &= \beta^T \beta + E[\hat{\epsilon}_p] - \beta^T \hat{\beta} \\ &= \beta^T \beta - \beta^T \hat{\beta} \\ &= 0 \end{aligned}$$

$$\textcircled{2} \quad \text{Var}[y_p - \hat{y}_p] = \sigma^2 (1 + \beta^T (\mathbf{X}' \mathbf{X})^{-1} \beta)$$

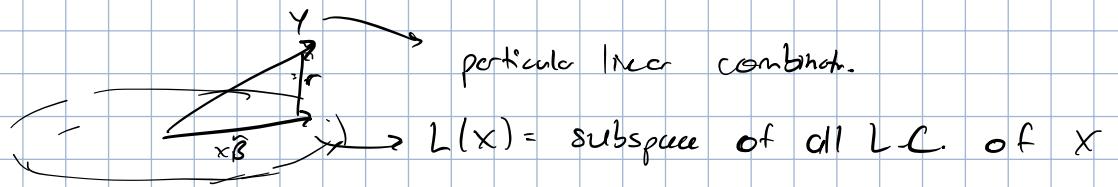
$$\therefore \frac{y_p - \hat{y}_p}{\text{se}(y_p - \hat{y}_p)} \sim t_{n-(p+1)}$$

$$\hookrightarrow \sqrt{\sigma^2 (1 + \beta^T (\mathbf{X}' \mathbf{X})^{-1} \beta)}$$

## Geometric Interpretation

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad 1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1}, \quad x_j = \begin{bmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{bmatrix}_{n \times 1}, \quad X = \underbrace{\begin{bmatrix} 1 & x_1 & \dots & x_j \end{bmatrix}}_{1 \times (p+1) \text{ row vector}}$$

All obsv. of  $j^{\text{th}}$  varcl.



To minimize  $x\hat{\beta}$  &  $r$ :

$$\textcircled{1} \text{ Define } r = y - x\hat{\beta}$$

\textcircled{2} Orthogonal to all  $x$ :

$$x'(y - x\hat{\beta}) = 0$$

$$x'y - x'x\hat{\beta} = 0$$

$$x'x\hat{\beta} = x'y$$

$$\hat{\beta} = (x'x)^{-1}x'y$$

Hat matrix:

$$Y = \hat{X}\hat{\beta}$$

$$= X \underbrace{(x'x)^{-1}x'}_{\text{H}} y$$

H: projection of  $y$  on  $L(x)$

Also:

$$x'(r) = \hat{y}'(r) = 0 \Rightarrow \text{blk orthogonal.}$$

ANOVA in MLR

Source	S.S.	d.f.	M.S.
Regression	$SSR = \sum (\hat{y}_i - \bar{y})^2$	$p$	$MSR = \frac{SSR}{p}$
Error	$SSE = \sum (y_i - \hat{y}_i)^2$	$n - (p + 1)$	$MSE = \frac{1}{n-(p+1)} \cdot SSE$
Total	$SST = \sum (y_i - \bar{y})^2$	$n - 1$	

$\underbrace{\quad}_{\text{Still adds up}}$

F-test of overall significance:

\textcircled{1} Hypothesis:

$$H_0: \beta_1 = \dots = \beta_p = 0. \quad H_a: \text{one of them is non-zero}$$

② Stmt:

$$F = \frac{MSR}{MSE} = \frac{SSR/p}{SSE/(n-(p+1))} \sim F_{p, n-(p+1)}$$

③ Conduct:

Choose a critical value & test

$$\beta^2 =$$

$$R^2 = \frac{SSR}{SST} \uparrow \propto \# \text{ of predictors}$$

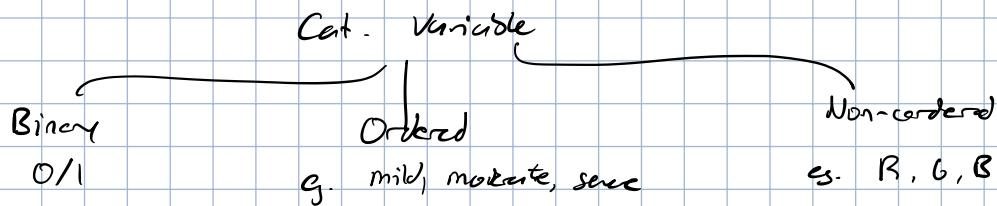
Doesn't work for MLR  $\Rightarrow$  adjust  $R^2$

$$\text{Adj. } R^2 = (1 - \frac{n-1}{n-(p+1)}) (1 - R^2)$$

Interpretation diff.

## SPECIFICATION ISSUES

## Categorical Variables



## Indicator variables:

① # of categories

② Create binary indicator for  $n-1$  categories.

$$x_{i2} = \begin{cases} 1, & \text{if "type" = t}_1 \\ 0, & \text{otherwise} \end{cases}, \quad x_{i3} = \begin{cases} 1, & \text{if "type" = t}_2 \\ 0, & \text{otherwise} \end{cases}$$

If all indicator variables = 0  $\Rightarrow$  encode info for last cat.

↳ Inclusion:  $(X'X)^{-1}$  not possible  $\rightarrow$  not L.I.  $\rightarrow$  not invertible

## Numerical method:

1 variable for all categories

$$x_{i2} = \begin{cases} 0, & \text{if } t_{i2} = t, \\ \vdots & \vdots \\ n & \vdots \end{cases}$$

→ Assures linear regn to  
2<sub>17</sub> → not un.

Model:

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$$

Categorical indicator variable

Interpretation of coef. for cat. variables?

① Take mean response for all categories type:

$$\text{type} = \text{"bc"} \rightarrow y_i = \beta_0 + \beta_1 x_1 - ①$$

$$= \text{"prot"} \rightarrow y_i = \beta_0 + \beta_1 x_1 + \beta_2 - ②$$

$$= \text{"wc"} \rightarrow y_i = \beta_0 + \beta_1 x_1 + \beta_2 + \beta_3 - ③$$

② Take diff.!

$$\text{Avg. diff. in resp. b/w "prot" \& "bc": } ② - ① = \beta_2$$

$$\text{ll} \quad \text{"wc" \& "bc": } ③ - ① = \beta_3$$

$$\text{ll} \quad \text{"wc" \& "prot": } ③ - ② = \beta_3 - \beta_2$$

Testing:

Format: is resp. stat. diff. b/w cat 1 \& cat 2:

$$H_0: \beta_{2/3} = 0 \quad \text{vs. } H_a: \beta_{-} \neq 0 \quad \left. \begin{array}{l} \\ \\ \text{cat. of indicator variable of interest} \end{array} \right\} t\text{-test}$$

Special case:

$$\text{To test } \beta_i - \beta_j = 0$$

① Estimate

② Variance:

Covariance matrix

$$\text{Var}(\beta_i - \beta_j) = \text{Var}(\beta_i) + \text{Var}(\beta_j) - 2 \text{Cov}(\beta_i, \beta_j)$$

③ Standard error:  $\sqrt{\text{Var}(\beta_i - \beta_j)}$

④ t-test

Interaction Effects

Obj: include effect of  $x_1 \leftarrow x_2$

Create interaction model:

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \varepsilon_i$$

int. effect

Interpretation: main model

Change in unit incen of  $x_1$ :  $\beta_1 + \beta_3 x_2$

II  $x_2$ :  $\beta_2 + \beta_3 x_1$

For binary variables:

Coefficient of interaction  $\rightarrow$  if slope is different b/w categories.

↓

Test:

Obj: is relationship between b/w  $y$  &  $x_1, x_2 = 0$  diff conpar to rel. b/w  $y$  &  $x_1, x_2 = 1$

$H_0$ : interaction coef. = 0 vs.  $IC \neq 0$

can drop interaction effect.

## GENERAL LINEAR HYPOTHESES

Tests

Single variable  
is significant

$t$ -test

All variables  
are non-sig.

F-test for sig.

Specific comb  
of variable sign.

General linear hyp.

Test:

$H_0: A\beta = \vec{0}$   $H_a: A\beta \neq \vec{0}$

$A$ : matrix of  $d \times (p+1)$

# of constraints on  $\beta$

Ex: //

$$y_1 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$H_0: \beta_2 = \beta_3 = 0$   $H_a: \text{otherwise}$

↓

① # of linear constraints

$$d = 2$$

② Write out matrix

$$\begin{bmatrix} a & b & c & 0 \\ c & f & g & h \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} a\beta_0 + b\beta_1 + c\beta_2 + d\beta_3 \\ -f - g - h \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

∴ Put 0's in cols of variable I don't care about

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Ex://  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}$

$$H_0: \beta_1 = 0, \beta_2 = \beta_3, \text{ then } 0$$

① Make R.H.S. of const. to be 0:

$$\beta_1 = 0, \beta_2 - \beta_3 = 0$$

② Make A:

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

Q: What is actual test?

1. Reduced model: model if  $A\beta = 0$

2. Given reduced model  $\rightarrow$  find fitted value ( $\hat{Y}_A$ )

3. Find SSE of  $\hat{Y}_A \Rightarrow \sum (y_i - \hat{y}_{iA})^2 = SSE_A$

4. Test statistic:

$$\frac{\frac{(SSE_A - SSE)}{k}}{SSE / (n-p-1)} \sim F_{k, n-p-1}$$

Principle:

$SSE - SSE_A$  is "extra" error explained by including full varish

↳ we hope that  $SSE - SSE_A \approx SSE$  if reduced model is good!

↳ we hope test statistic is low  $\rightarrow$  reject H<sub>0</sub>

Ex:// Prestige example:

Full model:  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_1 x_2 + \beta_5 x_1 x_3$

eductor "Type" indicators Interaction effects

Question: is the effect of education the same for diff. types?

## ① Hypotheses

$$H_0: \beta_4 = \beta_5 = 0 \quad \text{vs.} \quad H_a: \neg H_0$$

$$A = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

## ② Create test statistic

A: Reduce model  $\Rightarrow$  assure  $H_0$  is true model?

$$y_r = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}$$

B: Fit both full & reduced model

C: Find SSE for each model

In R output  $\Rightarrow$  residual sum of sq. =  $\sqrt{MSE}$

$$\therefore SSE = (\text{res. ss})^2 \times \text{d.f.}$$

D: Create F-stat

$$\frac{(SSE_A - SSE) / l}{SSE / (n - (p + 1))} = 0.844$$

## ③ Compare against critical val.

Is  $F > F_{\alpha, l, n - (p + 1)}$ ?  $\Rightarrow$  No it wasn't, don't reject  $H_0$

ANOVA in R:

## ① anova (model\_reduced, model)

Model 1: . — —

Model 2: -- --

	Res. d.f.	$\beta_{SS}$	d.f.	sum of sq.	F	pval
1	model 1 d.f.	$SSE_{\text{reduced}}$				
2	" 2 "	$SSE_{\text{full}}$	1 d.f. $\downarrow l$	$SSE_A - SSE$	f.stat	$P(F > f\text{-val})$

## MODEL ASSUMPTIONS AND RESIDUAL ANALYSIS

Intro:

Assumptions:

- ① Linear rel. b/w  $y$  &  $x_1, \dots, x_p$
- ②  $E[\varepsilon_i] = 0$
- ③  $\text{Var}[\varepsilon_i] = \sigma^2 \quad \forall i$  (constant variance)
- ④  $\varepsilon_i \sim N$
- ⑤  $\varepsilon_i \perp \varepsilon_j \quad \forall i \neq j$   
↳ OK if random sampling

most important  
least important

If assumptions not true  $\rightarrow$  analysis is faulty

Monitors  $\varepsilon_i$  is best  $\leftarrow$  random error from theoretical model  
↳ Use residuals instead.

Claim: Residuals behave very similar to random error:

Proof:

$$\begin{aligned}
 r &= y - \hat{y} \\
 &= (I - H) y \\
 &= (I - H)(X\beta + \varepsilon) \\
 &= (X\beta - HX\beta) + (I - H)\varepsilon \quad \xrightarrow{\text{cancel } X\beta} \\
 &= (X\beta - \cancel{X\beta}) + (I - H)\varepsilon \quad \xrightarrow{\text{cancel } X\beta} \\
 &= (I - H)\varepsilon \quad \xrightarrow{\text{cancel } I - H} \varepsilon \quad \xrightarrow{\text{cancel } X\beta} \\
 &= 0
 \end{aligned}$$

linear transf. of  $\varepsilon$ !

$$\therefore r \sim MVN(0, (I - H)' \sigma^2 I (I - H))$$

$$r \sim MVN(0, \sigma^2 (I - H))$$

Analysis: If off-diag of  $I - H \neq 0 \Rightarrow r_i$  is not indep.

But, if  $H \ll I \Rightarrow r_i \approx \varepsilon_i$

## Residual Analysis

Types of residuals:

① Raw

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

② Standardized:

$$d_i = \frac{r_i}{\hat{\sigma}} \sim N(0, 1) \quad \xrightarrow{\text{estimator}} \sqrt{\text{Var}(r_i)}$$

Effectively standardizes b/c  $\text{Var}(r_i) \approx \sigma^2$  if  $H \ll I$

### ③ Studentized:

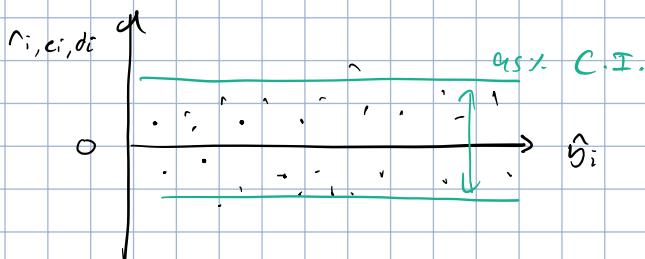
$$e_i = \frac{r_i}{\sqrt{\hat{\sigma}^2 (1 - h_{ii})}} \sim \text{Studentized residuals.}$$

Violation of assumption should reflect in residuals. We plot residuals!

Plots:

#### ① Residuals vs. fitted value

Good:

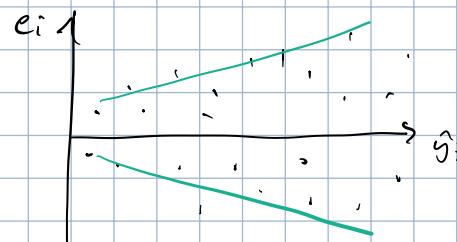


Random scatter around x-axis within horizontal band.

Why:  $\hat{y} \perp r$

Conclu: no visible model defect.

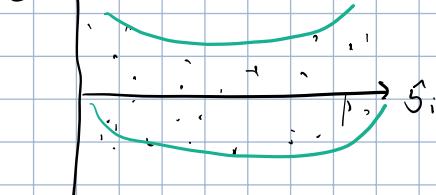
Bad:



Violation of common variance

Heteroscedasticity

Bad:



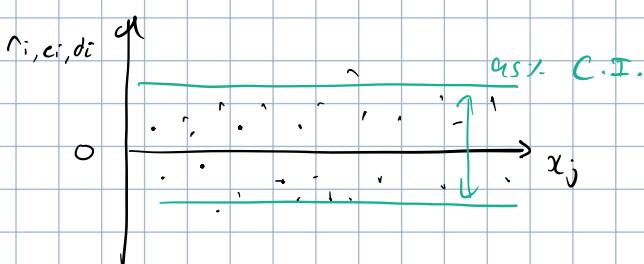
Linear model is not adequate

↳ Non-linear w.r.t. some x

↳ Missing variable

#### ② Residuals vs. explanatory

Good:

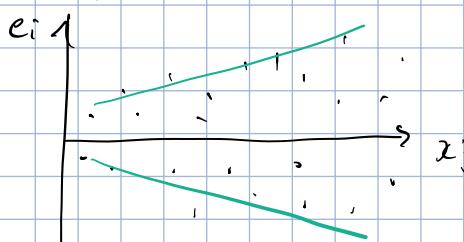


Random scatter around x-axis within horizontal band.

Why:  $x_j \perp r$

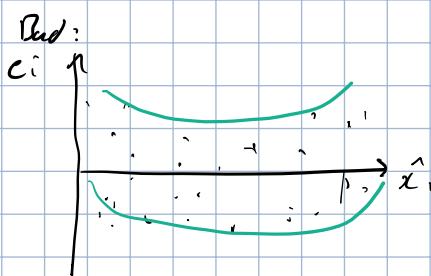
Conclu: no visible model defect.

Bad:



Violation of common variance

Heteroscedasticity

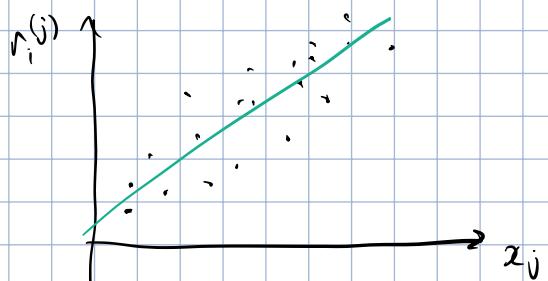


Linear model is not adequate  
 ↳ Non-linear not  $x_j$  (or  $x_j^2$ ?)  
 ↳ Missing variable

### ③ Partial residual plot

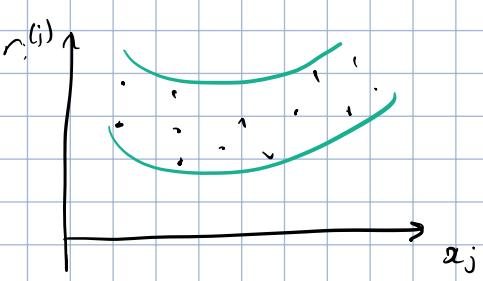
Partial residual of  $x_j$ :  $r_i^{(j)} = \underbrace{e_i}_{\text{full model residuals}} + \underbrace{\hat{\beta}_j x_{ij}}_{\text{effect of } x_j}$

Good:



$x_j$  is linear

Bad:



Non-linear  $\rightarrow$  higher order term  
 might help.

### ④ Q-Q plots

Check normality of  $r_i / d_i / e_i$ :

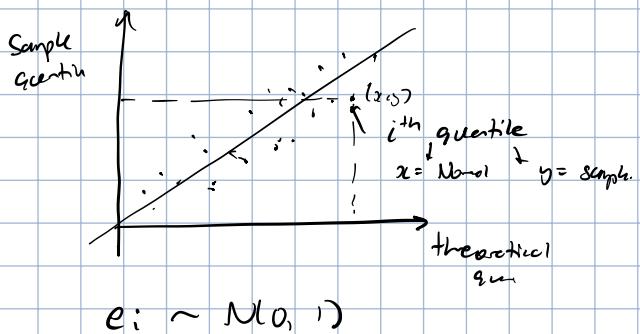
Steps:

① Order standardize residuals  $d_{(1)} < d_{(2)} < \dots < d_{(n)}$

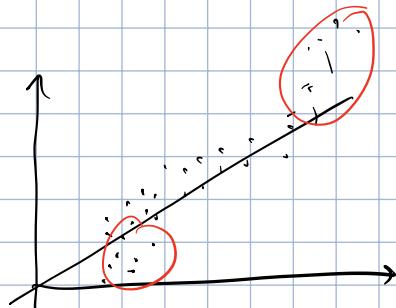
② Find quantiles of data

③ Plot against quantiles of  $N(0, 1)$

Good:

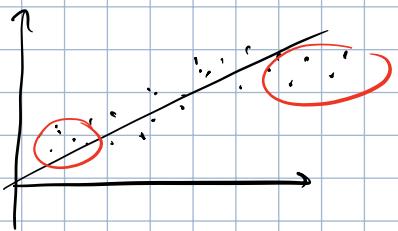


$e_i \sim N(0, 1)$

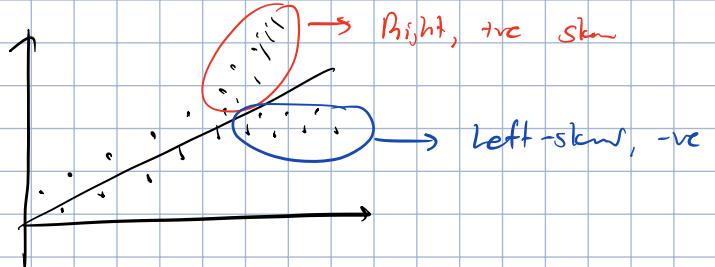


Heavy-tailed  $\Rightarrow$  outliers.

How to fix

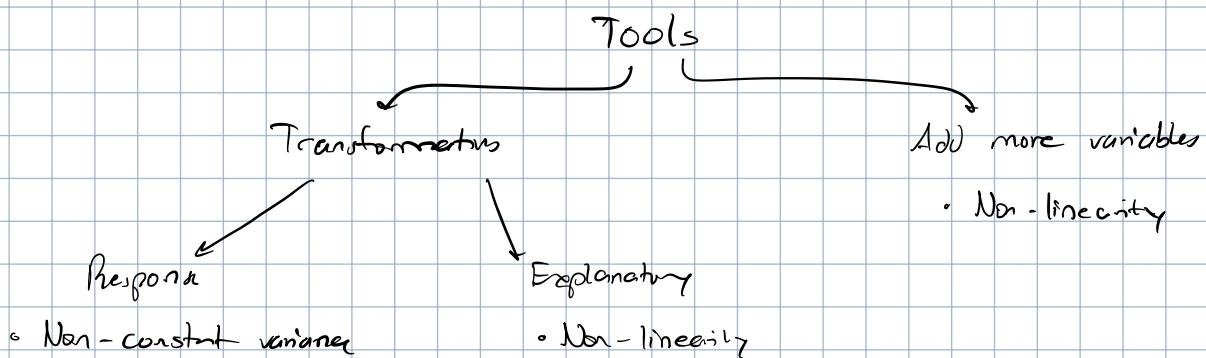


Light-tailed  
Tric



Skew  
Transform → make normal

## Addressing Model Assumption Problems



## Transforming Response Variables

Obj: constant variance  $\rightarrow$  apply transform  $g(y)$   $\rightarrow$  fit model  $g(y) = \dots$

Assumption: variance of  $y$  is func. of mean  $\mu_i$

$$\begin{aligned} \text{Var}(y_i) &= \text{Var}(\varepsilon_i) \\ &= h(\mu_i) \sigma^2 \quad \text{Assum, } h > 0 \end{aligned}$$

Back off sol:

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

Taylor expansion at  $\mu_i$ :

$$g(y_i) \doteq g(\mu_i) + g'(\mu_i)(y_i - \mu_i)$$

$$\text{Var}(g(y_i)) \doteq g'(\mu_i)^2 \text{Var}(y_i)$$

$$= g'(\mu_i)^2 \cdot h(\mu_i) \sigma^2$$

$$\sigma^2 = g'(\mu_i)^2 \cdot h(\mu_i) \sigma^2$$

$$\boxed{g'(\mu_i)^2 \propto \frac{1}{h(\mu_i)}}$$

In practice: make assumption of  $h(\cdot)$   $\rightarrow$  use formula to find  $g(\cdot)$

Ex:// a)  $h(\mu_i) = \mu_i \Rightarrow \text{Var}(y_i) = \sigma^2 \mu_i$

$$s^2(\mu_i)^2 \propto \frac{1}{\mu_i}$$

$$s^2(\mu_i) \propto \frac{1}{\sqrt{\mu_i}}$$

$$s(\mu_i) \propto \sqrt{\frac{1}{\mu_i}} = \sqrt{\mu_i} + c$$

$\therefore s$  should be sqrt transform.

Box-Cox Transform

$$s(y_i) = \begin{cases} \frac{y_i^{\lambda-1}}{\lambda}, & \lambda \neq 0 \\ \log y_i, & \lambda = 0 \end{cases}$$

Q: How to choose  $\lambda$ ?

Try many  $\lambda \rightarrow$  choose one w/ highest log likelihood.

Special case:

$$\lambda = \frac{1}{2} \rightarrow \text{sqrt}$$

$$\lambda = 0 \Rightarrow \log$$

$$\lambda = 1 \rightarrow \text{no transform}$$

$$\lambda = -1 \rightarrow \text{reciprocal transform.}$$

Interpretation changes w/ transform  $\Rightarrow \hat{\beta}_j$  is change in  $s(y_i)$  if  $x_j \uparrow$

Ex://  $y_i$  is log transform.

$$\log y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

Original:

$$y = c^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}$$

$\therefore$  1 unit chng in  $\hat{\beta}_j \Rightarrow 100(e^{\hat{\beta}_j} - 1) \%$  chng in  $y$

Transforming explanatory vars. / add more vars.

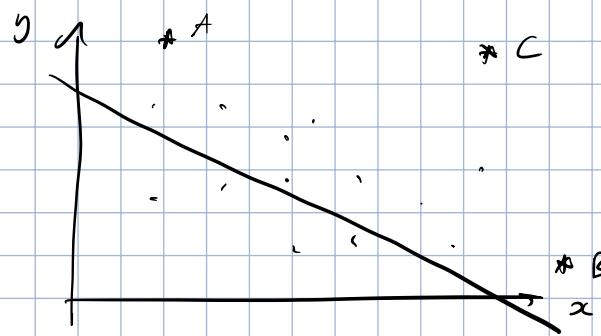
(1) Power transform of variable

(2) Polynomial term addition:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \epsilon_i$$

## EFFECTS OF INDIVIDUAL OBSERVATIONS

### Intro



outliers:

- A: outlier in response
- B: outlier in expl.
- C: outlier in response & expl.

Outliers have big impact on regression

Q: Why do outliers occur?

1. Human error
2. Sampling issue (sample from wrong pop.)
3. Natural variation

Q: What do we do w/ outliers?

Human error → correct remov. O.W. → don't!

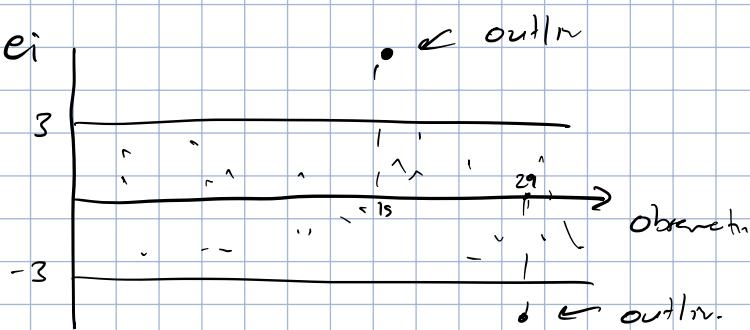
### Studentized Residuals - Outliers in Response

Studentized residual:

$$e_i = \frac{r_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}} \sim N(0, 1)$$

If  $|e_i| > 3 \Rightarrow$  outlier in response (v. extreme in standard normal distn.)

Plot:  $e_i$



### Leverage - Outliers in Explanatory

Theorem:  $h_{ii}$  indicates rel. b/w  $\hat{y}_i$  &  $y_i$ :

$$\text{L, } \hat{y} = X \hat{\beta}$$

$$\hat{y} = H Y$$

Consider  $\hat{y}_i$ :

$$\hat{y}_i = [h_{i1} \ h_{i2} \ \dots \ h_{ip}] \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$= h_{ii} y_i + \sum_{i \neq j} h_{ij} y_j$$

Properties of  $h_{ii}$ :

①  $0 < h_{ii} \leq 1$

Proof:

From idempotency ( $HH = H$ )

$$\begin{aligned} h_{ii} &= [h_{i1} \ \dots \ h_{ip}] \begin{bmatrix} h_{ii} \\ \vdots \\ h_{ii} \end{bmatrix} \\ &= \sum_j h_{ij} h_{ji} \end{aligned}$$

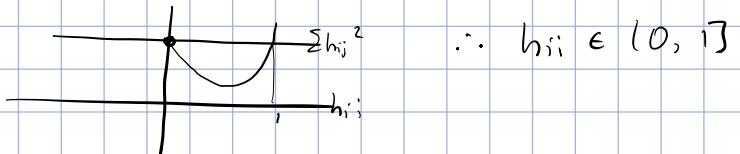
From symmetry:  $h_{ij} = h_{ji} \ \forall i, j$

$$\therefore \sum h_{ij} h_{ji} = \sum h_{ij}^2$$

$$= h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$$

$$\therefore h_{ii} = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$$

$$h_{ii}(1-h_{ii}) = \underbrace{\sum h_{ij}^2}_{\geq 0} \rightarrow \geq 0$$



② As  $h_{ii} \rightarrow 1$ , off diagonal  $\rightarrow 0$

Part:

$$\begin{aligned} h_{ii}(1-h_{ii}) &= \sum h_{ij}^2 \\ &\rightarrow 0 \text{ if } h_{ii} \rightarrow 1 \quad \rightarrow 0 \end{aligned}$$

Impl:

$$\hat{y}_i = y_i h_{ii} + \sum h_{ij} y_j \quad 0 \text{ if } h_{ii} \rightarrow 1$$

$$\therefore \hat{y}_i \approx y_i$$

③ If  $h_{ii} \gg h_{ij} \Rightarrow \hat{y}_i$  is det. by  $y_i$

$\therefore$  If  $i$  was an outlier  $\rightarrow$  fitter value in model will be skewed.

$\therefore i$  has high leverage!

④ If, obs.  $(x_{i1}, \dots, x_{ip})$  is close to center of exp.  $(\bar{x}_1, \dots, \bar{x}_p)$

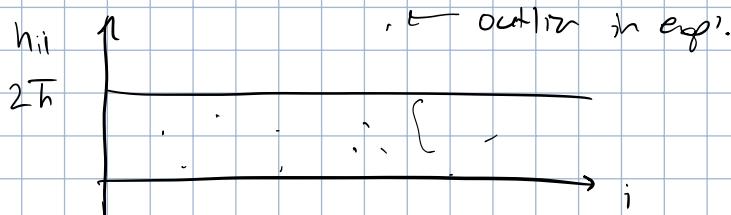
$\Leftrightarrow h_{ii}$  is small.

Rule of thumb:

$$h_{ii} > 2\bar{h} = \frac{2(p+1)}{n} \Rightarrow i^{\text{th}} \text{ obs. has high leverage.}$$

Corollary: From prop. 4, outlier in explanatory.

Plot:



### Cook's Distance - Influential Data points

Defn of infl. d. p.: removing point causes big change in  $\hat{\beta}$

Note:  $\hat{\beta}^{(i)}$  is fit to  $\hat{\beta}$  after removing  $i^{\text{th}}$  obs.

To find distance b/w  $\hat{\beta}$  &  $\hat{\beta}^{(i)}$   $\Rightarrow$  Cook's distance.

$$D_i = \frac{(\hat{\beta} - \hat{\beta}^{(i)})' (X'X) (\hat{\beta} - \hat{\beta}^{(i)})}{\hat{\sigma}^2 (p+1)}$$

$$= e_i^2 \cdot \frac{h_{ii}}{1-h_{ii}} \cdot \frac{1}{p+1}$$

$\therefore D_i$  is high if obs. is influential in both expl. & response.

Just b/c outlier  $\rightarrow$  not always influential

Rule of thumb:

$$D_i = \begin{cases} > 0.5 \Rightarrow \text{might be influential} \\ > 1 / \text{sig. diff from other data point} \Rightarrow \text{influential} \end{cases}$$

## MODEL SELECTION

Goal: find optimal # of predictors

### Selection Criteria

①  $R^2_{\text{adj}}$

$$R_{\text{adj}}^2 = 1 - \frac{SSE / (n - k - 1)}{SST / (n - 1)} \quad \left. \quad \right\} \begin{matrix} k = \# \text{ of \\ predictors} \end{matrix}$$

$$= 1 - \frac{n - 1}{n - k - 1} \cdot (1 - R^2)$$

Brill paralitic if  $k \uparrow$

If we just use  $R^2$ , model with more predictors will have bigger  $R^2$

## ② Alsaike Informator Criteria (AIC)

$$\text{AIC} = 2g - 2 \ln(L(\hat{\theta})) - \frac{n}{2} \log\left(\frac{\text{SSE}}{n}\right) + \text{constant}$$

$\uparrow$  # of model param =  $k+2$

The lower AIC, better the model

AIC penalizes high predictor models

### ③ Bayesian Inf. Criteria (BIC)

$$BIC = \frac{1}{2} \ln(n) - 2 \ln(L(\hat{\theta}))$$

Results high predictor model and & incl sample size

None of these interpretable.

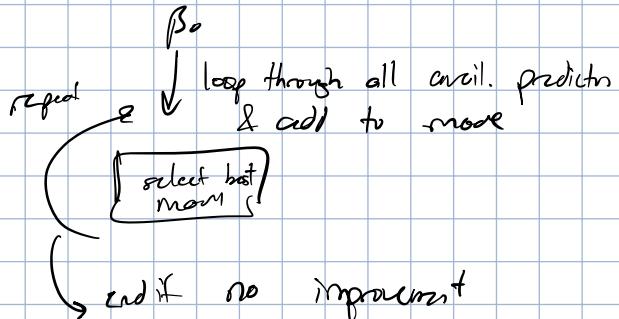
## Search Strategy

## ① All subset regression

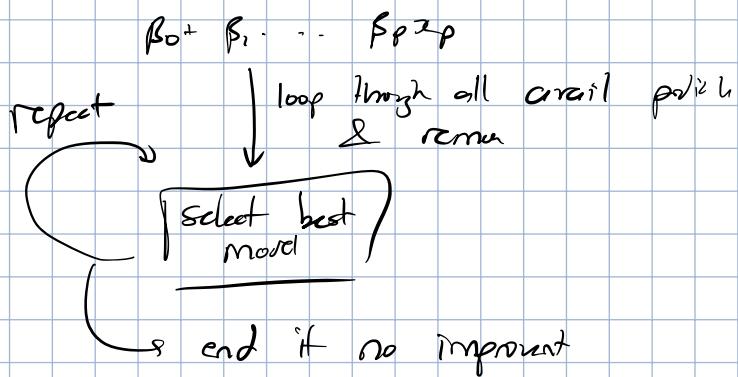
Find all subsets of predictors  $\rightarrow$  fit  $\rightarrow$  choose best

Very expensive (2<sup>o</sup> months)

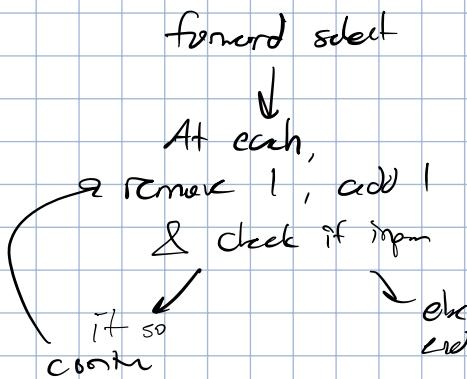
② Formed selection



### ③ Backward selection:



#### ④ Forward-backward stepwise:



### BUILDING PREDICTIVE MODELS

Metrics to evaluate model on new data (on testing/validation data)

#### ① Mean squared prediction error

$$MSPE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

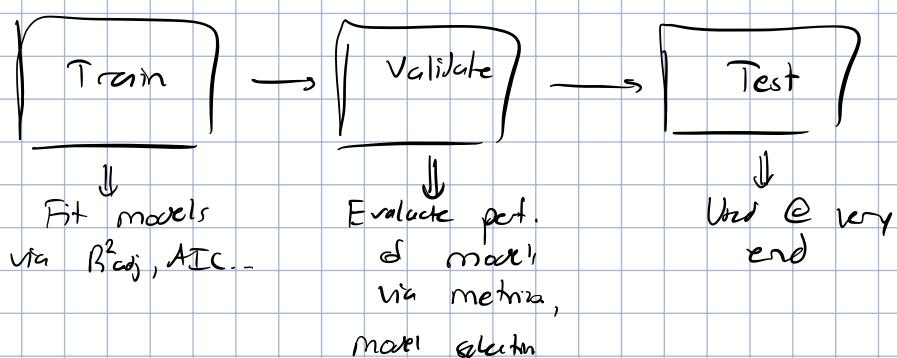
#### ② Root mean squared error

$$RMSE = \sqrt{MSPE}$$

#### ③ Mean absolute error:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Process of building:



How to split data?

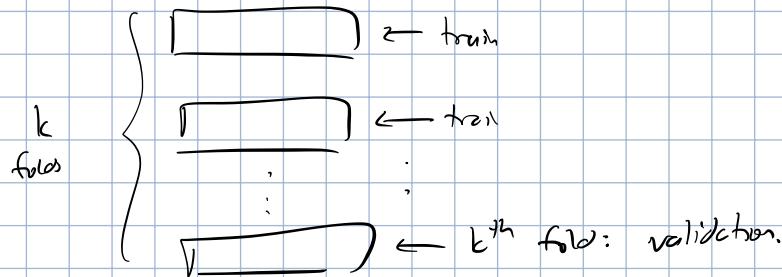
① Randomly

+ Simple

- Variation in diff. split  $\leftarrow$  randomness

② K-Fold C.V.

1. Divide data into k folds



2. Train model on  $k-1$  fold, eval on  $k^{\text{th}}$  fold

3. Repeat so each fold becomes eval fold once

4. Take avg. of error

$$\frac{1}{k} \sum_{k=1}^K \text{MSPE}_k$$