

# STAT 601 Lec 2

Lectures: MW 2:26–3:45 PM, 1313 Sterling

Disc 321: F 9:30–10:45 AM, 382 Van Hise

Disc 322: F 8:00–9:15 AM, 114 Social Work

Instructor: Wei-Yin Loh

Office hours: M 4:00–5:00 PM, 1217C MSC

TA: Siyu Wang

In-person office hours: T 2:30–3:30 PM, 6173C MSC

Zoom office hours: R 11 AM–12:00 PM

Class email list: stat601-2-f23@g-groups.wisc.edu

# Topics

1. Linear regression
  - Concepts
  - Theory
  - Practice
2. Generalized linear models — logistic, Poisson, proportional hazards
3. Design and analysis of statistical experiments
  - One and two-way factorial designs
  - Randomized block designs, including Latin square designs
  - Incomplete block designs

## Recommended texts

1. Rencher (2000), *Linear Models in Statistics*, Wiley
2. Box, Hunter & Hunter (2005), *Statistics for Experimenters*, 1st or 2nd edition, Wiley
3. Mead, Gilmour & Mead (2012), *Statistical Principles for the Design of Experiments*, Cambridge University Press

# **Grading**

**Homework.** 10% of average z-score, excluding 2 lowest

**Quizzes.** 10% of average z-score

**Midterm (Wed Oct 25).** 30% of z-score

**Final (Fri Dec 15).** 50% of z-score

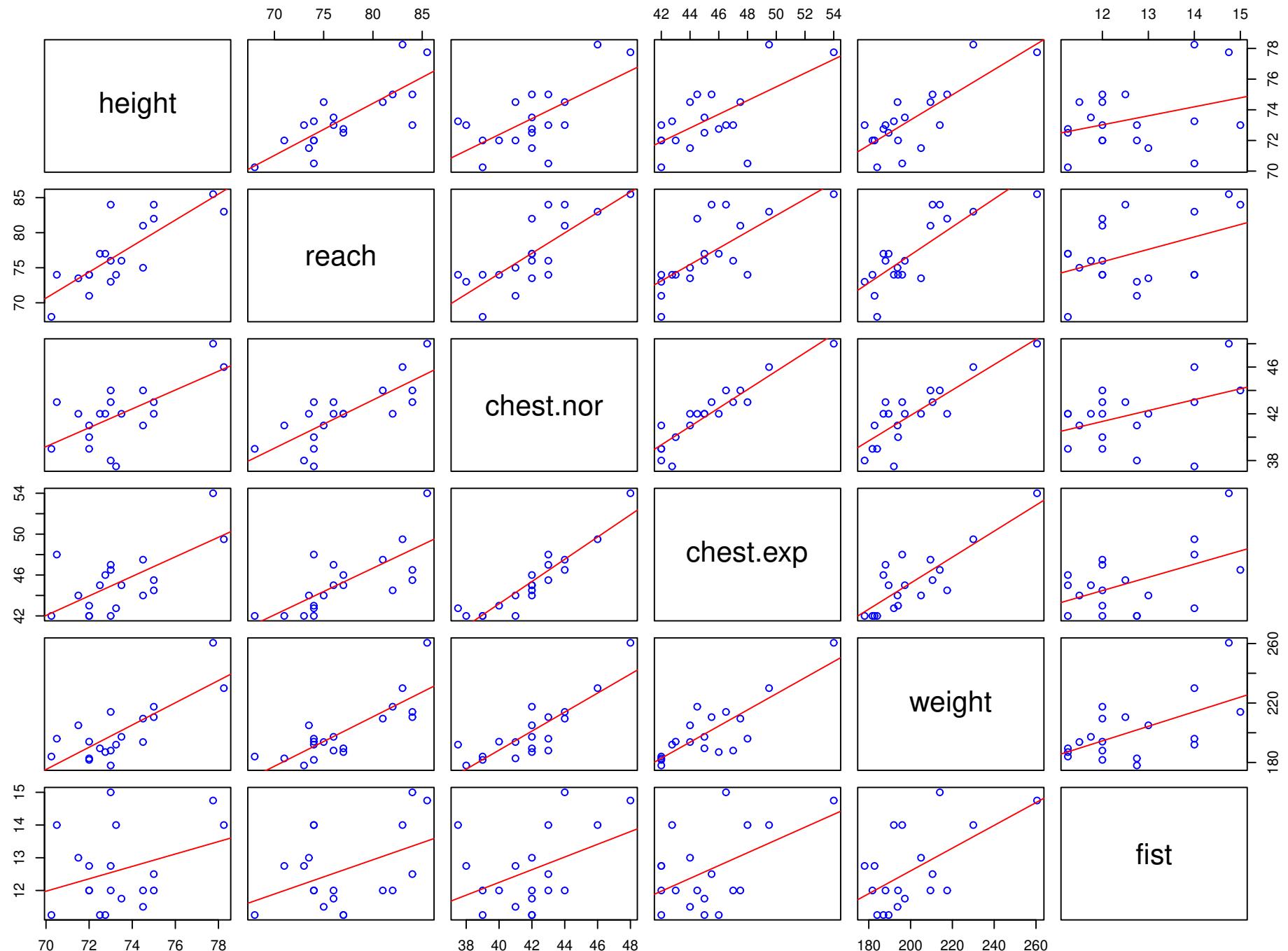
**Linear regression example:  
Height, reach and chest (normal & expanded),  
weight and fist measurements of  
19 heavyweight boxing champions**

<https://www.boxingscene.com/forums/showthread.php?t=286210>

Name	height	reach	chest.nor	chest.exp	weight	fist
Ali	75	84	43	45.5	210.5	12.5
Baer	74.5	81	44	47.5	209.5	12
Braddock	74.5	75	41	44	193.75	11.5
Carnera	77.75	85.5	48	54	260.5	14.75
Charles	72	74	39	42	181.75	12
Corbett	73	73	38	42	178	12.75
Dempsey	72.75	77	42	46	187	11.25
Foreman	75	82	42	44.5	217.5	12
Frazier	71.5	73.5	42	44	205	13
Johnson	73.25	74	37.5	42.75	192	14
Liston	73	84	44	46.5	214	15
Louis	73.5	76	42	45	197.25	11.75
Marciano	70.25	68	39	42	184	11.25
Patterson	72	71	41	42	182.75	12.75
Schmeling	73	76	43	47	188	12
Sullivan	70.5	74	43	48	196	14
Tunney	72.5	77	42	45	189.5	11.25
Walcott	72	74	40	43	194	12
Willard	78.25	83	46	49.5	230	14

## R code for pairwise plots with regression lines

```
> z <- read.csv("boxers.csv")
> names(z)
[1] "Name"    "height"   "reach"    "chest.nor" "chest.exp" "weight"
[7] "fist"
> pairs(z[,-1])  ## plain
> pairs(z[,-1],
  panel=function(x,y,...){
    points(x,y,...)
    abline(lm(y ~ x),col="red")
  },col="blue")
```



# Correlations



```
> cor(z[,-1])
```

	height	reach	chest.nor	chest.exp	weight	fist
height	1.00	0.79	0.65	0.65	0.78	0.33
reach	0.79	1.00	0.78	0.73	0.82	0.43
chest.nor	0.65	0.78	1.00	0.93	0.83	0.43
chest.exp	0.65	0.73	0.93	1.00	0.83	0.50
weight	0.78	0.82	0.83	0.83	1.00	0.58
fist	0.33	0.43	0.43	0.50	0.58	1.00

# Least squares linear regression

- Given data  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , let  $\bar{x} = n^{-1} \sum_i x_i$ ,  $\bar{y} = n^{-1} \sum_i y_i$
- Find the line  $y = \hat{\beta}_0 + \hat{\beta}_1 x$  that minimizes the sum of squared errors

sufficient statistic

$$\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

- Solution: differentiate SSE w.r.t.  $\hat{\beta}_0$  and  $\hat{\beta}_1$  and solve for them

$$\partial \text{SSE} / \partial \hat{\beta}_0 = -2 \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\partial \text{SSE} / \partial \hat{\beta}_1 = -2 \sum_i x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

- Let  $s_x^2 = (n-1)^{-1} \sum_i (x_i - \bar{x})^2$ ,  $s_y^2 = (n-1)^{-1} \sum_i (y_i - \bar{y})^2$ , and  $r = (n-1)^{-1} \sum_i (x_i - \bar{x})(y_i - \bar{y}) / (s_x s_y)$ . Least-squares estimates (LSEs) are

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_i x_i y_i - n \bar{x} \bar{y}}{\sum_i x_i^2 - n \bar{x}^2} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i x_i^2 - n \bar{x}^2} = \frac{r s_y}{s_x}$$

$$X \quad E(y|x)$$

$$\sum (y_i - ?)^2$$

minimize

X

X

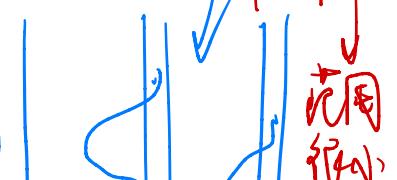
X

$$\text{min } (\cdot)^2$$

$$\text{fix } E(y|x)$$

$$E(y|x) = \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

$$= \hat{\beta}_0 + \hat{\beta}_1 x$$



为什么高度和胸围不相关  $\downarrow$  unrelated.

```
> model.nor <- lm(height ~ chest.nor, data=z)
```

```
> summary(model.nor)
```

	Estimate	Std. Error	t value	Pr(> t )
--	----------	------------	---------	----------

(Intercept) 51.5251 <sup>截距</sup> 6.2310 8.269 2.32e-07 \*\*\*

chest.nor 0.5214 0.1484 3.514 0.00266 \*\*

Residual standard error: 1.646 on 17 degrees of freedom

Multiple R-squared: 0.4208, Adjusted R-squared: 0.3867

F-statistic: 12.35 on 1 and 17 DF, p-value: 0.002661

p-value

=  $P(A \text{ statistic is as large or larger than the observed value})$

$H_0: \beta_1 = 0, \hat{\beta}_1 = 0.5214$

or larger

```
> model.exp <- lm(height ~ chest.exp, data=z)
```

```
> summary(model.exp)
```

p-value =  $P_{H_0}(|\hat{\beta}_1| \geq 0.5214)$

	Estimate	Std. Error	t value	Pr(> t )
--	----------	------------	---------	----------

(Intercept) 53.1082 5.6950 9.325 4.27e-08 \*\*\*

chest.exp 0.4478 0.1255 3.568 0.00237 \*\*

Residual standard error: 1.636 on 17 degrees of freedom

Multiple R-squared: 0.4281, Adjusted R-squared: 0.3945

F-statistic: 12.73 on 1 and 17 DF, p-value: 0.00237

Assumptions  
1. random sample  
2. linear model

## Where do the p-values come from?

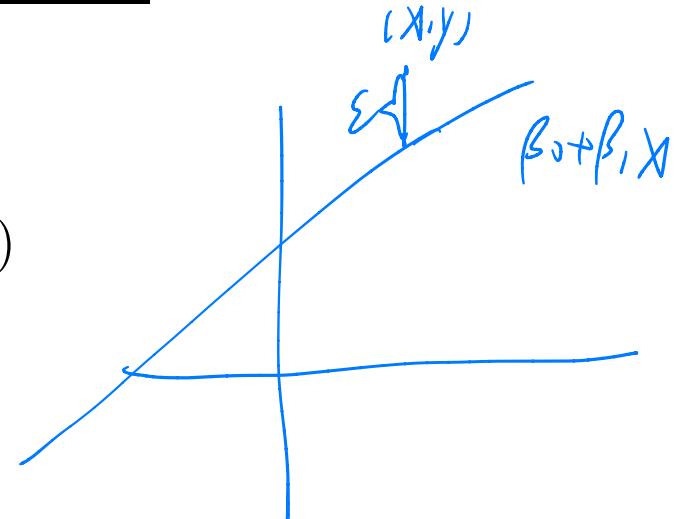
- Assume that, for each  $x$ ,

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

where  $\epsilon$  is independent, identically distributed normal (IIDN) with mean 0 and variance  $\sigma^2$

- Equivalent to writing

$$y \sim N(\beta_0 + \beta_1 x, \sigma^2)$$



# Hypothesis testing terms

**Null hypothesis**

$$H_0: \mu = \mu_0$$

**Alternative hypothesis**

$$H_1: \mu \neq \mu_0$$

**Significance level**

$$\alpha = P_{H_0} (\text{Type I error})$$

**Test statistic**

**Null distribution of test statistic**

**Rejection region**

**Critical value**

**Type I error** :  $H_0$  is true but reject  $H_0$

**Type II error**

**Power**

**P-value**

reject  $H_0$

do not  
reject  $H_0$

	True	False
True	Type I error	✓
False	✓	Type II error

# Illustration: one-sample t-test

- Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, \sigma^2)$
- Let  $\bar{X} = n^{-1} \sum_i X_i$  and  $S^2 = (n - 1)^{-1} \sum_i (X_i - \bar{X})^2$
- Then  $\bar{X} \sim N(\mu, n^{-1}\sigma^2)$ ,  $(n - 1)S^2/\sigma^2 \sim \chi_{n-1}^2$ , and  $\bar{X}$  is independent of  $S^2$

**Null hypothesis.**  $H_0 : \mu = \mu_0$

**Alternative hypothesis.**  $H_1 : \mu \neq \mu_0$

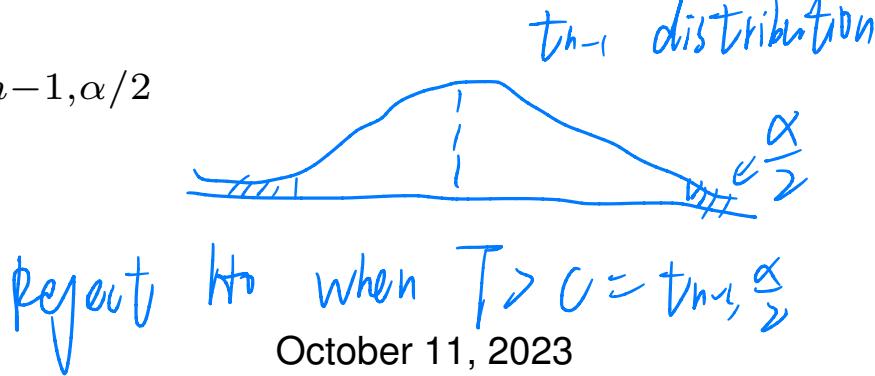
**Significance level.**  $\alpha$  (specified)

**Test statistic.**  $T = \sqrt{n}(\bar{X} - \mu_0)/S$  <sup>t(n-1)</sup> 有  $\sqrt{n}$ , 因为要变成  $t$  分布

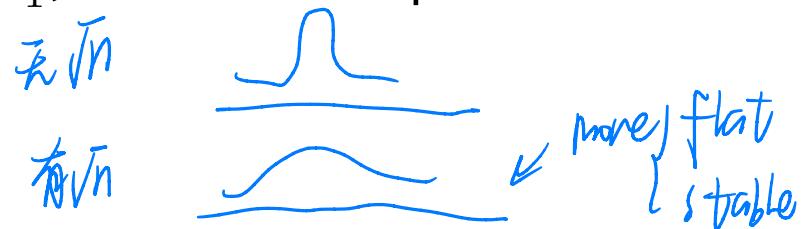
**Null distribution of statistic.**  $t_{n-1}$  ( $t$  distribution with  $n - 1$  df)

**Rejection region.**  $|T| > t_{n-1, \alpha/2}$

**Critical value.**  $t_{n-1, \alpha/2}$



$$\begin{aligned} X &\sim N(0, 1) \\ Y &\sim \mathcal{N}(y) \\ \frac{X}{\sqrt{n}} &\sim t(n) \end{aligned}$$



$$\begin{aligned} \frac{\bar{X} - \mu_0}{S/\sqrt{n}} &\sim N(0, 1) \\ SD(\bar{X}) &= \frac{\sigma}{\sqrt{n}} \\ SD(\bar{X} - \mu_0) &= \frac{\sigma}{\sqrt{n}} \\ SD\left(\frac{\bar{X} - \mu_0}{S/\sqrt{n}}\right) &= 1 \end{aligned}$$

**$P(\text{Type I error})$ .**

$$P_{\mu_0}(\text{Reject } H_0) = P_{\mu_0}(|T| > t_{n-1, \alpha/2}) = P(|t_{n-1}| > t_{n-1, \alpha/2}) = \alpha$$

**$P(\text{Type II error})$ .** For any  $\mu \neq \mu_0$  and as  $n \rightarrow \infty$ ,

$$P(\text{Type I}) + P(\text{Type II}) \neq 1$$

$$P_\mu(\text{Type II error}) = P_\mu(\text{Do not reject } H_0)$$

$$= P_\mu(|T| \leq t_{n-1, \alpha/2})$$

$$= P_\mu(-t_{n-1, \alpha/2} \leq \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \leq t_{n-1, \alpha/2})$$

$$= P_\mu\left(-t_{n-1, \alpha/2} - \frac{\mu - \mu_0}{S/\sqrt{n}} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq t_{n-1, \alpha/2} - \frac{\mu - \mu_0}{S/\sqrt{n}}\right)$$

$$\approx P\left(-t_{n-1, \alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} \leq t_{n-1} \leq t_{n-1, \alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) \rightarrow 0$$

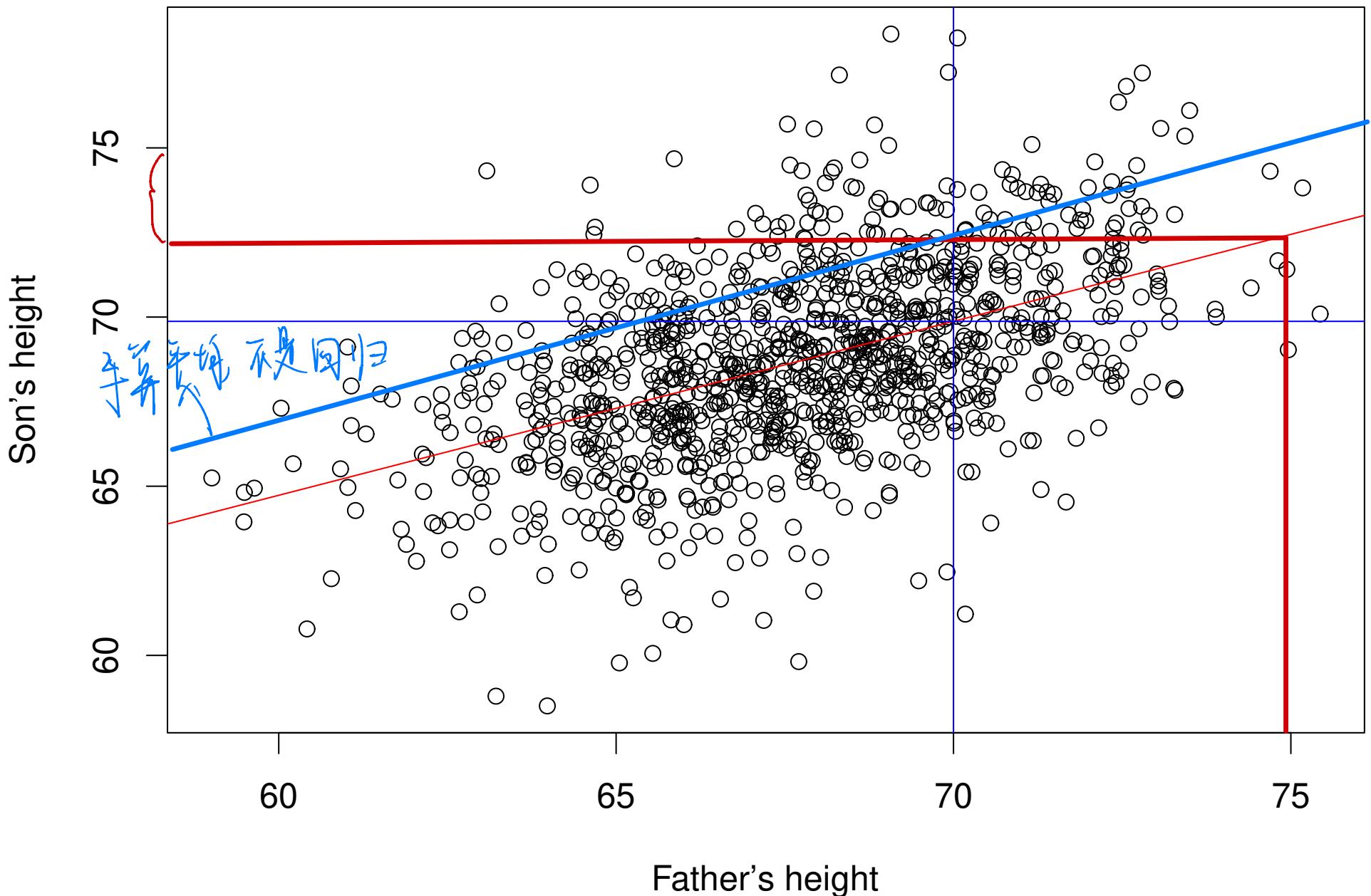
Since  $P(-t_{n-1, \alpha/2} - \delta \leq t_{n-1} \leq t_{n-1, \alpha/2} - \delta) < 1 - \alpha$  for every  $\delta \neq 0$ , it follows that  $P_\mu(\text{Type II error}) < 1 - \alpha$  for every  $\mu \neq \mu_0$

**Power.**  $P_\mu(\text{Reject } H_0) = 1 - P_\mu(\text{Type II error}) \rightarrow 1$  for  $\mu \neq \mu_0$  ( $\geq \alpha$  for all  $\mu$ )

**P-value.**  $P_{\mu_0}(|T| > T_{\text{obs}}) = P(|t_{n-1}| > \sqrt{n}(\bar{X}_{\text{obs}} - \mu_0)/S_{\text{obs}})$

# Plotting Pearson's father-son data (requires UsingR and HistData packages)

```
> library("UsingR")
> data(package="UsingR") ## print names of datasets in library
> df <- as.data.frame(father.son)
> names(df)
[1] "fheight" "sheight"
> plot(df$sheight ~ df$fheight,xlab="Father's height",
       ylab="Son's height")
> model <- lm(df$sheight ~ df$fheight)
> abline(model,col="red")
> abline(v=70,col="blue")
> mean.son.height <- model$coef [1]+model$coef [2]*70
> abline(h=mean.son.height,col="blue")
```



# Maximum likelihood estimates (MLEs)

- Let  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$
- Joint density of  $\mathbf{y}$  given  $\mathbf{x}$  is

$$f(\mathbf{y} | \mathbf{x}, \beta_0, \beta_1, \sigma) = (2\pi\sigma^2)^{-n/2} \exp[-(2\sigma^2)^{-1} \sum_i \{y_i - (\beta_0 + \beta_1 x_i)\}^2]$$

- Likelihood function is

$$\ell(\beta_0, \beta_1, \sigma | \mathbf{x}, \mathbf{y}) = (2\pi\sigma^2)^{-n/2} \exp[-(2\sigma^2)^{-1} \sum_i \{y_i - (\beta_0 + \beta_1 x_i)\}^2]$$

- Log-likelihood function is

do not depend on  $\sigma$

$$\log \ell(\beta_0, \beta_1, \sigma) = -(n/2) \log \sigma^2 - (2\sigma^2)^{-1} \sum_i \{y_i - (\beta_0 + \beta_1 x_i)\}^2 + \text{constant}$$

- For any  $\sigma^2$ ,  $\log \ell$  is maximum when  $\sum_i \{y_i - (\beta_0 + \beta_1 x_i)\}^2$  is minimum
- Hence MLEs  $(\hat{\beta}_0, \hat{\beta}_1)$  are same as LSEs

## Maximum likelihood estimate of $\sigma^2$

- MLE  $\hat{\sigma}^2$  is the solution of

$$\partial \log \ell / \partial \sigma^2 = -n/(2\sigma^2) + (2\sigma^4)^{-1} \sum_i \{y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)\}^2 = 0$$

- Therefore

$$\hat{\sigma}^2 = n^{-1} \sum_i \{y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)\}^2 = \underline{\text{SSE}/n} \quad \text{not unbiased}$$

$$s^2 = \frac{\text{SSE}}{n-2} \quad \text{unbiased}$$

$\uparrow$  2 parameter

$s \rightarrow SD(\text{sample})$   
 $\sigma \rightarrow SD(\text{population})$

$$E(s^2) = \sigma^2 \cancel{\Rightarrow} E(s) = \sigma$$



## **How to interpret regression coefficients and how to select among models**

# height vs chest.nor & height vs chest.exp

dev of  $\beta_j$  is  $\hat{\beta}_j = 0$

```
> model.nor <- lm(height ~ chest.nor, data=z)
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	51.5251	6.2310	8.269	2.32e-07 ***
chest.nor	0.5214	0.1484	3.514	0.00266 **

Residual standard error: 1.646 on 17 degrees of freedom

Multiple R-squared: 0.4208, Adjusted R-squared: 0.3867

F-statistic: 12.35 on 1 and 17 DF, p-value: 0.002661

$F = \frac{RSS}{TSS}$ , 独立变量对因变量的度量单位的多样性程度

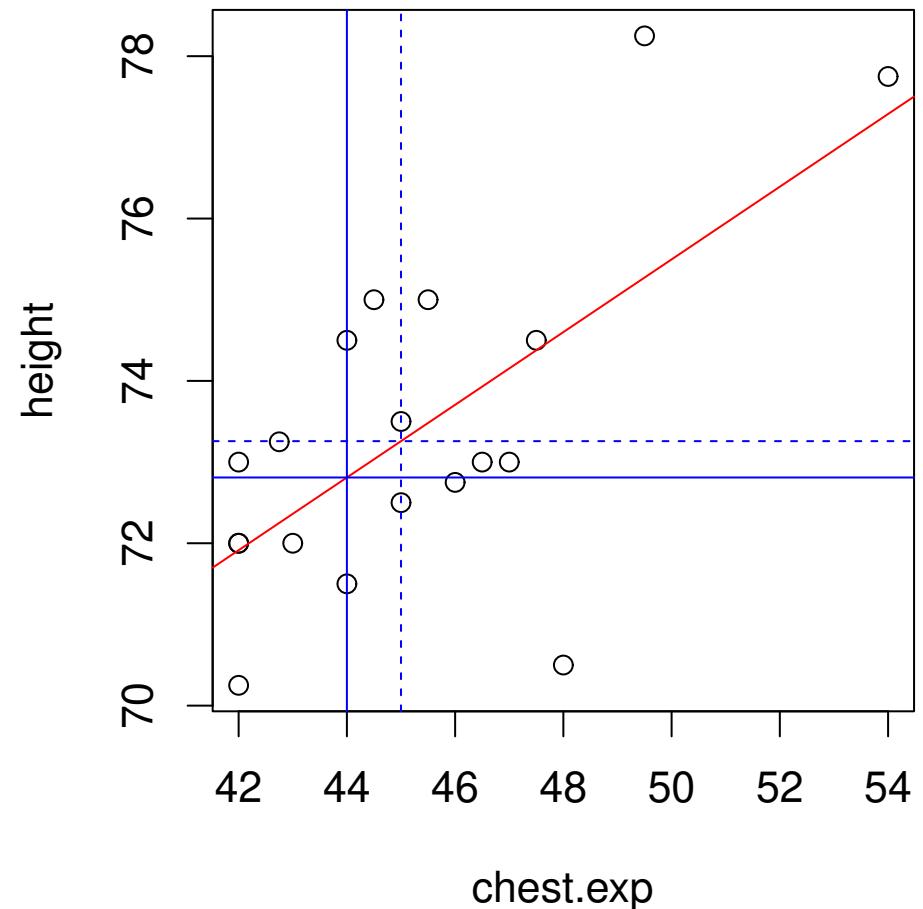
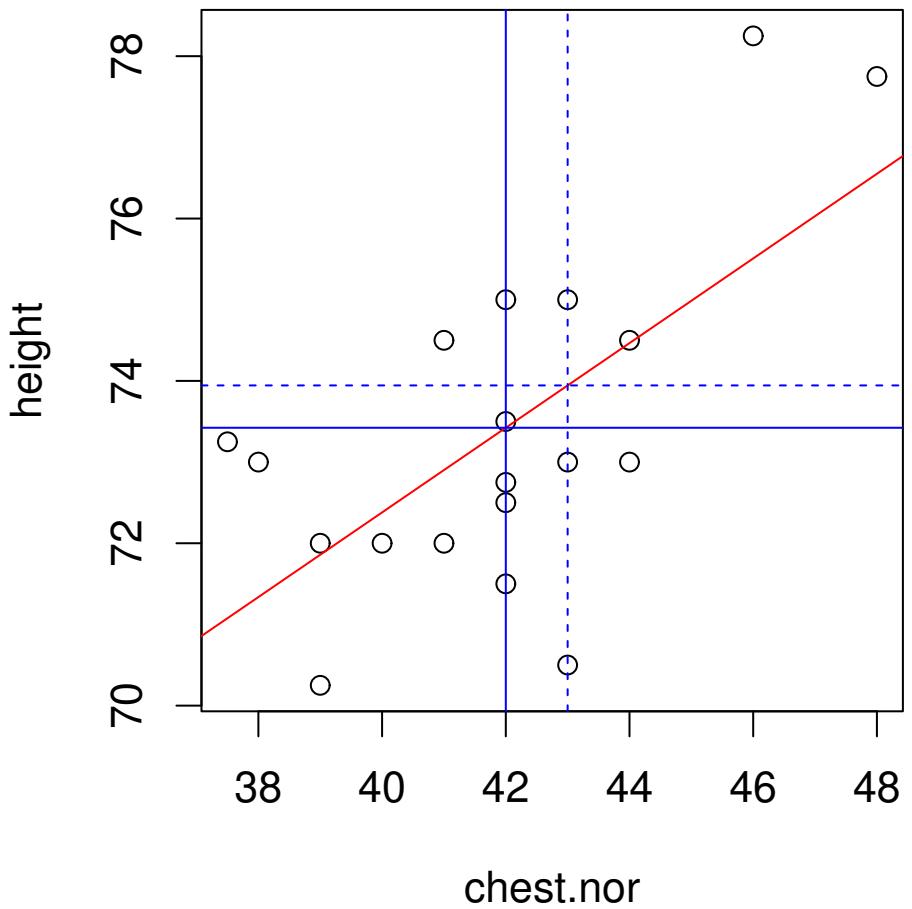
```
> model.exp <- lm(height ~ chest.exp, data=z)
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	53.1082	5.6950	9.325	4.27e-08 ***
chest.exp	0.4478	0.1255	3.568	0.00237 **

Residual standard error: 1.636 on 17 degrees of freedom

Multiple R-squared: 0.4281, Adjusted R-squared: 0.3945

F-statistic: 12.73 on 1 and 17 DF, p-value: 0.00237



检验两个模型

择取

$$Y_1 = \text{intercept} + \text{chest.nor}$$

$$Y_2 = \text{intercept} + \text{chest.exp}$$

$$Y_3 = \text{intercept} + \text{chest.exp}$$

conditional F > p 值

$$H_0: \beta_1 = 0$$

$$H_0: \beta_2 = 0$$

## height vs chest.nor & chest.exp

```
> model.both <- lm(height ~ chest.nor+chest.exp, data=z)
```

```
> summary(model.both)
```

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	51.6175	6.3120	8.178	4.17e-07	***
chest.nor	0.2407	0.4002	0.601	0.556	
chest.exp	0.2578	0.3407	0.757	0.460	

Residual standard error: 1.667 on 16 degrees of freedom

Multiple R-squared: 0.4408, Adjusted R-squared: 0.3709

F-statistic: 6.306 on 2 and 16 DF, p-value: 0.009564

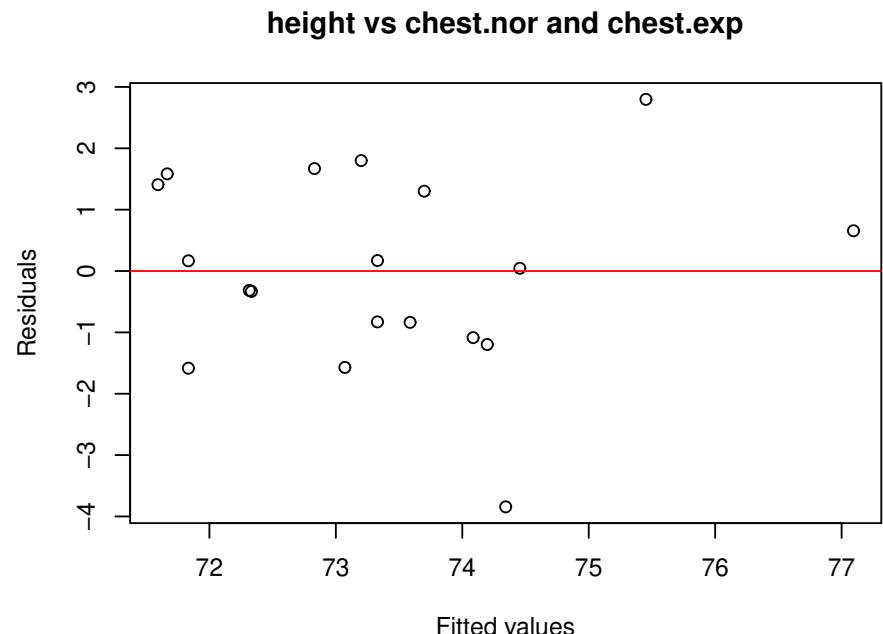
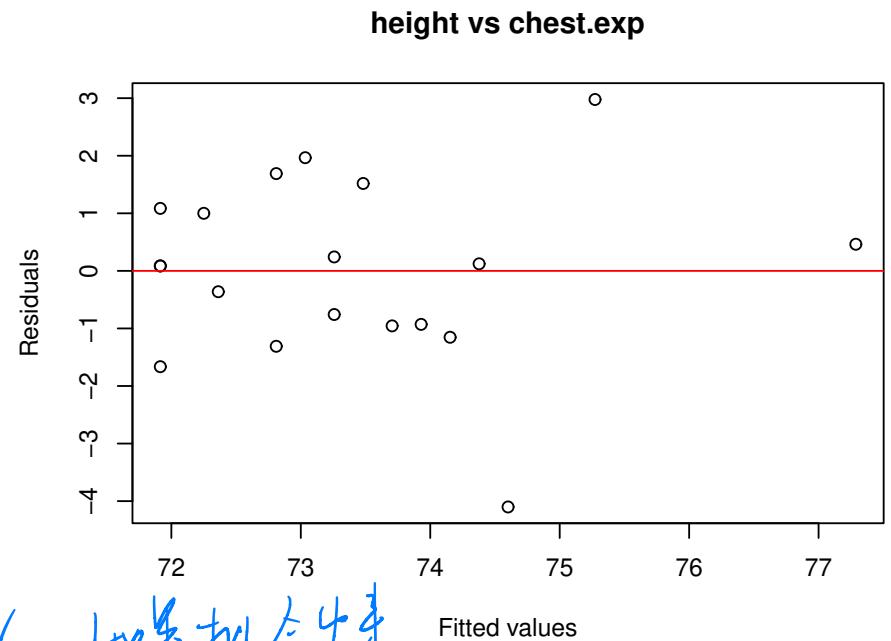
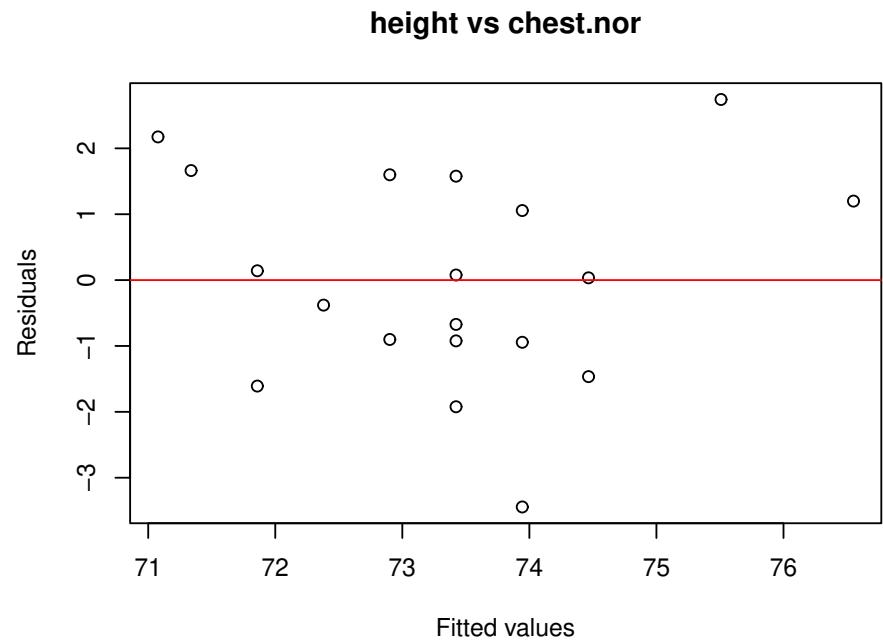
↓ 回归平方和  
残差平方和

↓  $H_0: \beta_1 = \beta_2 = \dots = 0$

違反

## Diagnostic plots to check for violations of model structure and IIDN assumptions

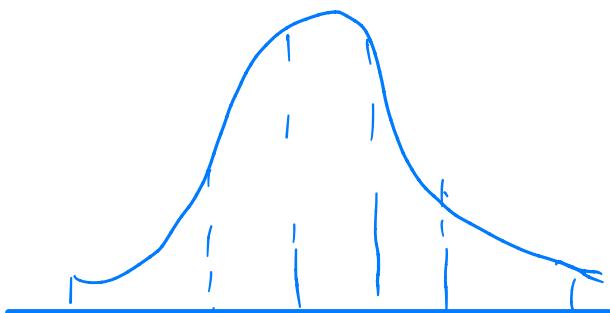
1. Plot of residuals versus fitted values—may indicate nonconstant variance
2. Normal probability plot of residuals—may indicate nonnormal  $\epsilon$  distribution



模型好,  $R^2$  大  
residual 加一起为 0。

## How is a normal probability plot constructed?

R function: `qqnorm(x)`

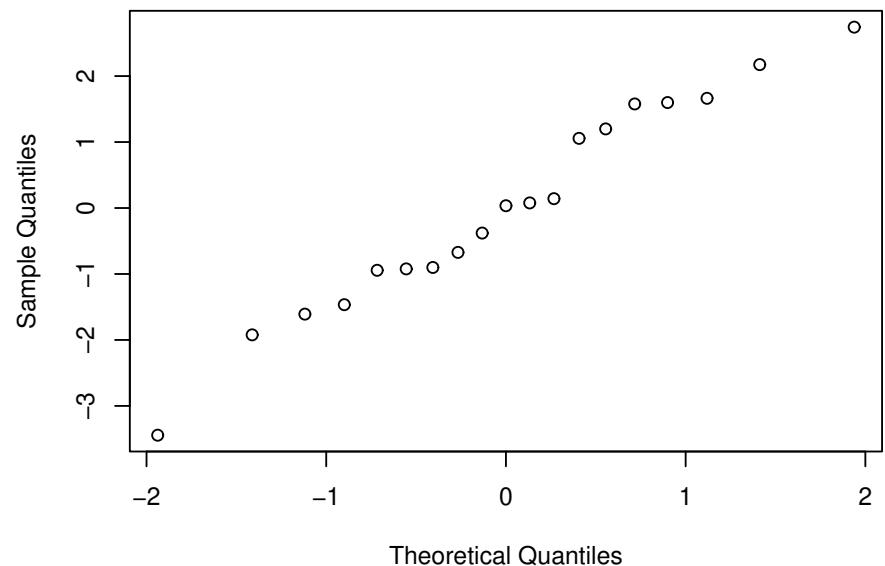


Quantile = 0.2 0.4 0.6 0.8

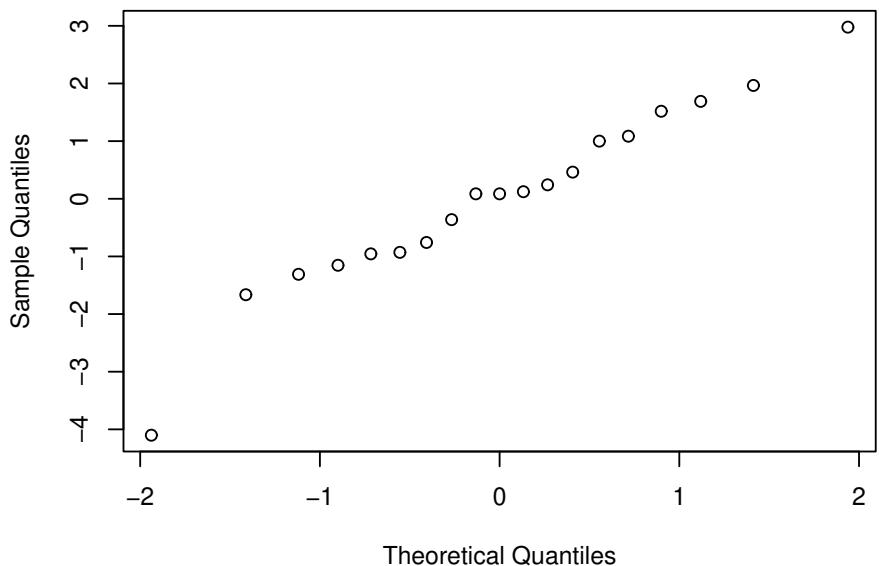
every region  
same 面积.

~ 则样本数据的分位数  
与已知标准正态分布的分位数相比较.

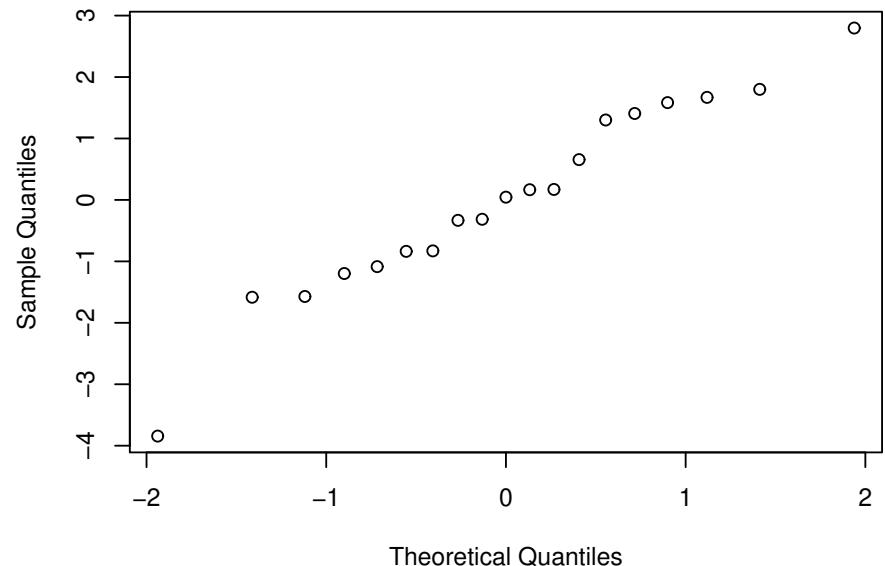
**Normal plot of residuals for height vs chest.nor**



**Normal plot of residuals for height vs chest.exp**



**Normal plot of residuals for height vs chest.nor and chest.exp**



测试正态分布，不能判断是否显著

```

> model.reach <- lm(height ~ reach, data=z)
> summary(model.reach)

              Estimate Std. Error t value Pr(>|t|)
(Intercept) 47.25588    4.86946   9.705 2.40e-08 ***
reach        0.33953    0.06316   5.376 5.03e-05 ***
Residual standard error: 1.316 on 17 degrees of freedom
Multiple R-squared:  0.6296,      Adjusted R-squared:  0.6078
F-statistic: 28.9 on 1 and 17 DF,  p-value: 5.033e-05

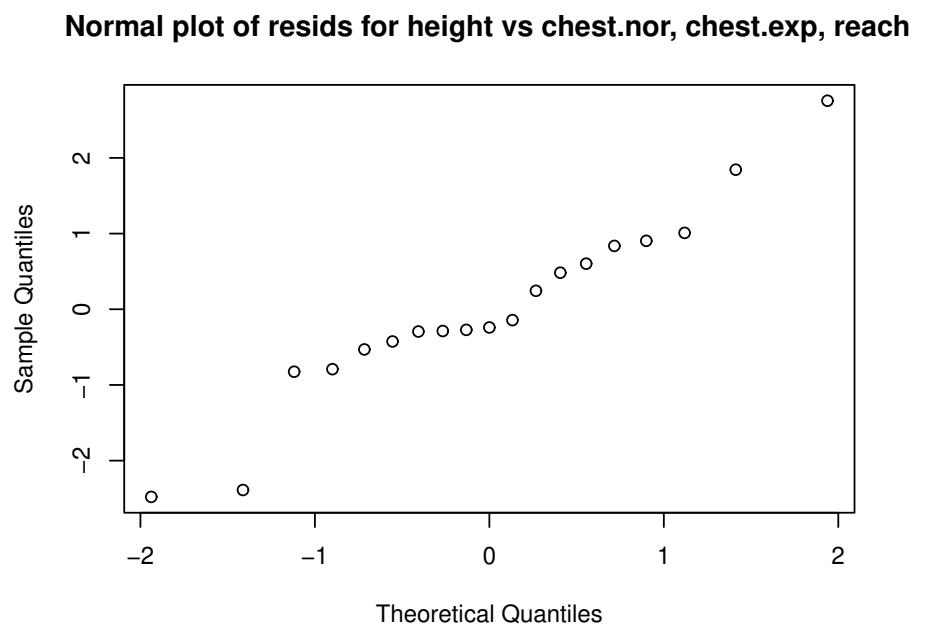
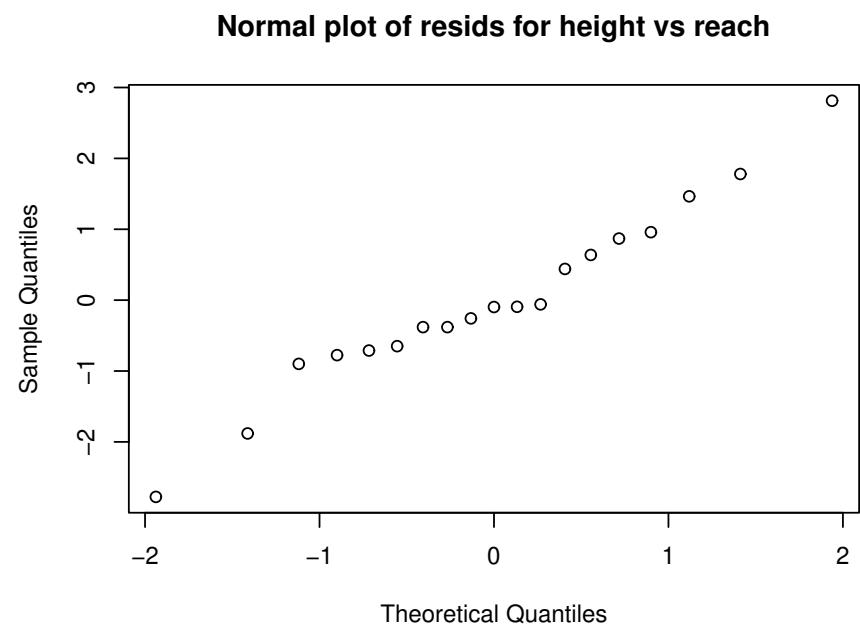
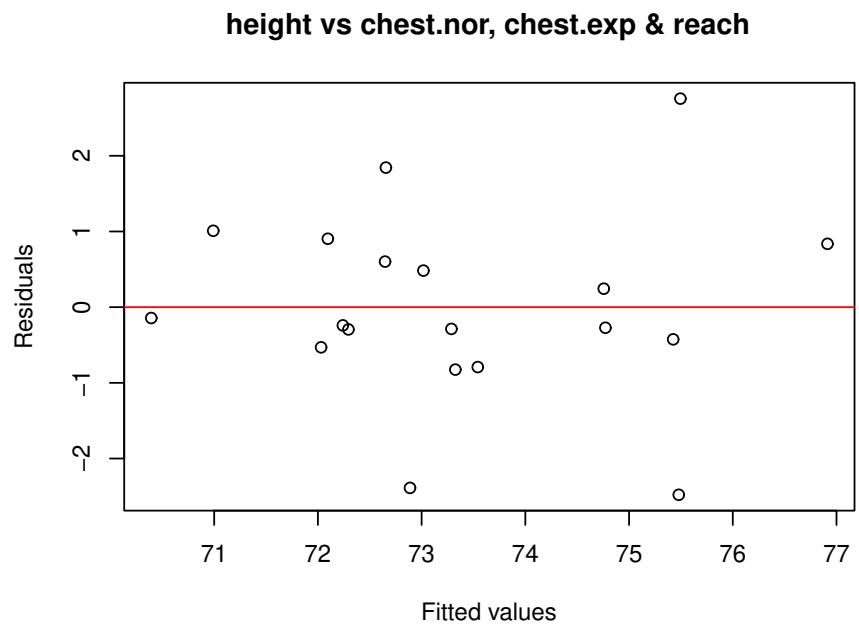
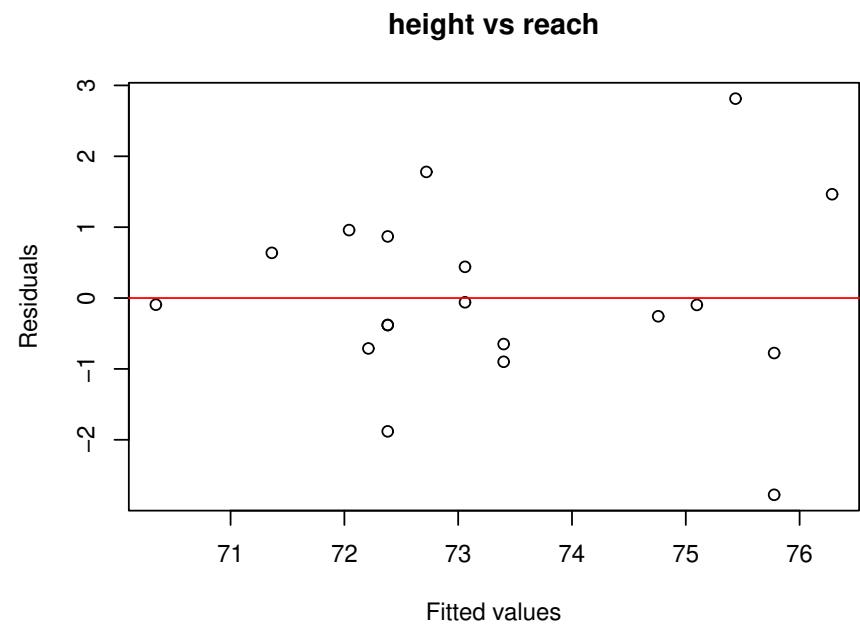
```

```

model.3 <- lm(height ~ chest.nor+chest.exp+reach, data=z)
> summary(model.3)

              Estimate Std. Error t value Pr(>|t|)
(Intercept) 46.7189    5.4423   8.584 3.58e-07 ***
chest.nor   -0.1630    0.3559  -0.458  0.65352
chest.exp    0.2166    0.2802   0.773  0.45153
reach        0.3079    0.1042   2.954  0.00986 **
Residual standard error: 1.369 on 15 degrees of freedom
Multiple R-squared:  0.6464,      Adjusted R-squared:  0.5757
F-statistic: 9.142 on 3 and 15 DF,  p-value: 0.001104

```



偏残差

## Partial residuals

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon_j$$

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

$$= \beta_0' + \beta_1' x_1 + \beta_2' x_2 + \epsilon_j$$

```
> lm(height ~ chest.nor+chest.exp+reach, data=z)
  (Intercept)    chest.nor      chest.exp       reach
        46.7189     -0.1630       0.2166       0.3079
```

```
> model1 <- lm(height ~ chest.exp+reach, data=z)
```

$$y \sim x_2 + x_3$$

```
> model2 <- lm(chest.nor ~ chest.exp+reach, data=z)
```

$$x_1 \sim x_2 + x_3$$

```
> lm(model1$res ~ model2$res)
```

```
(Intercept)  model2$res
```

```
2.339e-17 -1.630e-01
```

# Homework 1

## (due 11:59 PM Tue Sep 26 in Canvas)

Compute the leave-one-out cross-validation estimates of root mean squared prediction error (RMSPE) for the model

```
height ~ chest.nor + chest.exp + reach + weight + fist
```

and all its submodels

Which model has the lowest RMSPE?

Note: Let  $\mathcal{D}$  denote the full data set. Let  $(x_i, y_i)$  denote the  $i$ th observation and  $\mathcal{D}_{-i}$  the data subset without  $(x_i, y_i)$ . Let  $\mathcal{M}_{-i}$  denote the model fitted to  $\mathcal{D}_{-i}$  and  $\hat{y}_i^{(-i)}$  be the predicted value of  $y_i$  based on  $\mathcal{M}_{-i}$ . The leave-one-out RMSPE is

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

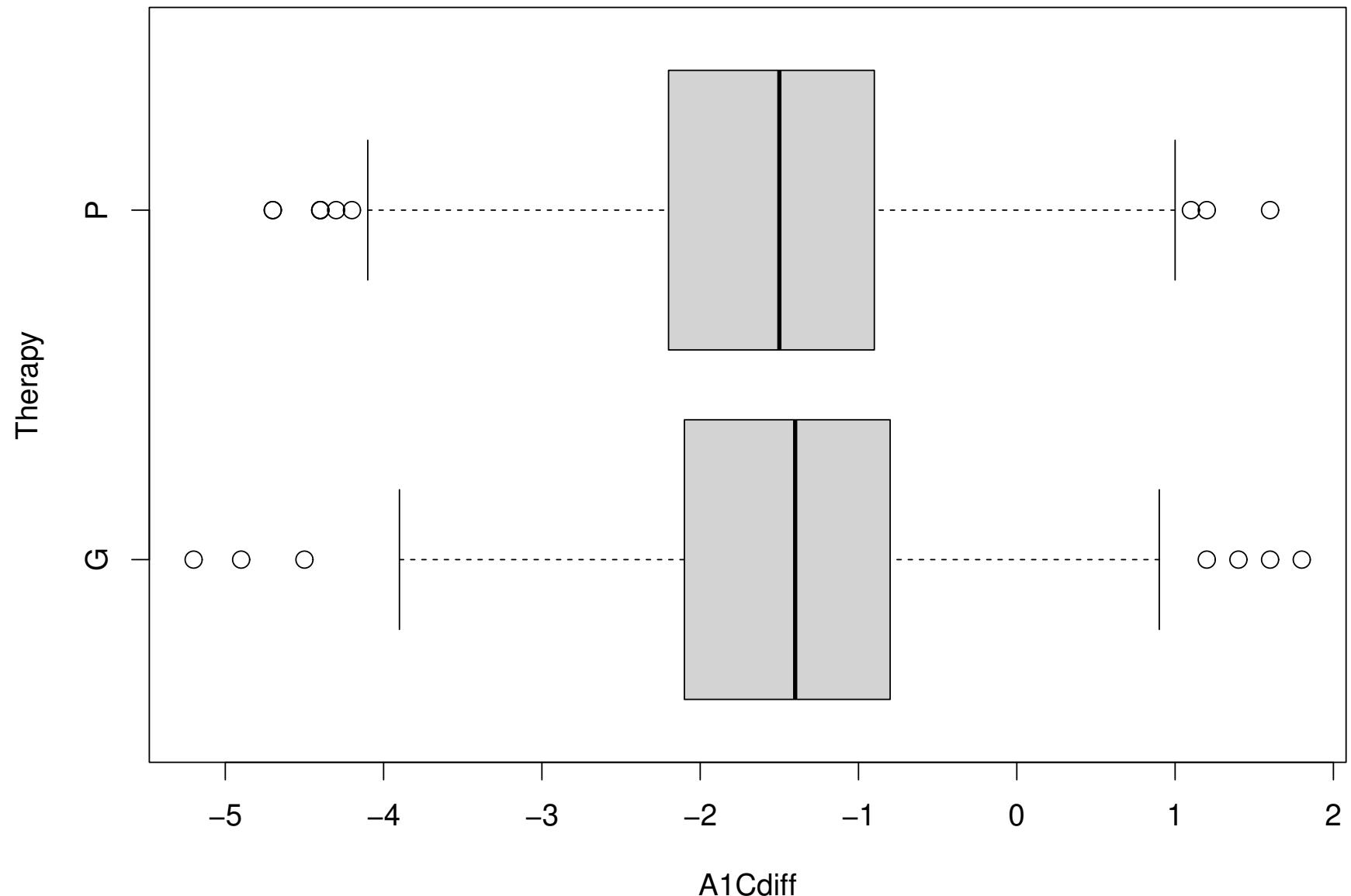
# Diabetes data

- 802 subjects from a multi-center, randomized double-blind clinical trial
- Subjects are 35–75 years old with HbA1c (blood sugar level) between 7.5% and 11.0%, for whom diet was prescribed for at least 3 months (non-diabetics have HbA1c below 5.7%; prediabetics 5.7–6.4%)
- Subjects randomized to a 52-week treatment period of drug G (*Gliclazide*) or P (*Pioglitazone*)
- *Gliclazide* increases amount of insulin produced by the pancreas
- *Pioglitazone* improves how body uses insulin (“insulin sensitizer”)
- 23 baseline (time 0) variables measured for each subject as well as their HbA1c at 52 weeks
- Interest is in whether there is a statistically significant difference between the effects of the two drugs on reduction of HbA1c

# Diabetes variables

Variable	Meaning	Variable	Meaning
HDL	high-density cholesterol	Age	age in years
LDL	low-density cholesterol	Weight	weight in kg
TotalChol	total cholesterol	BMI	body mass index
Triglycerides	type of fat (lipid)	Waist	waist in cm
Creatinine	kidney function test	A1CBase	HbA1C at week 0
GGT	gamma-glutamyl transferase	HomaS	measure of insulin sensitivity
ALT	alanine aminotransferase	HomalR	measure of insulin resistance
AST	aspartate aminotransferase	HomaB	measure of beta-cell function
FastInsulin	fasting insulin	Diastolic	Diastolic blood pressure
C-peptide	connecting peptide	Systolic	Systolic blood pressure
DDuration	Diabetes duration	Pulse	pulse rate
FastBG	Fasting blood glucose	A1C52	HbA1C at week 52

# Boxplots of A1Cdiff = A1C52 - A1CBase



## Proof

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, s_{\hat{\beta}} = \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$t = \frac{\bar{y}_A - \bar{y}_B}{\hat{\sigma} \cdot \sqrt{\frac{1}{n_A} + \frac{1}{n_B}}}$$

## Two-sample t-test

## Interpretation

Two tests are testing the same hypothesis:  $H_0 : \mu_A = \mu_B$ .

$$\hat{\beta} / s_{\hat{\beta}} = t$$

```
> z <- read.csv("diabetes.csv")
> names(z)
[1] "HDL"          "LDL"          "TotalChol"     "Triglycerides"
[5] "Creatinine"   "FastInsulin"  "ALT"          "AST"
[9] "GGT"          "DDuration"   "Age"          "Weight"
[13] "BMI"          "Waist"        "A1CBase"      "FastBG"
[17] "HomaS"        "HomaIR"       "HomaB"        "Diastolic"
[21] "Systolic"     "Pulse"        "Therapy"      "A1C52"
[25] "A1Cdiff"

> t.test(A1Cdiff ~ Therapy, data=z, var.equal=TRUE, alternative="two.sided")
t = 1.4089, df = 800, p-value = 0.1592
 无法拒绝原假设
    因为是0-1分布

> summary(lm(A1Cdiff ~ Therapy, data=z))
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -1.46445   0.05317 -27.542  <2e-16 ***
TherapyP     -0.10465   0.07428  -1.409    0.159
```

$\hat{Y} = \beta_0 + \beta_1 \cdot \text{Treat} + \beta_2 \cdot \text{Salary}$   
 $\downarrow$   $\beta_1$  表示  $\text{Treat}$   
 $\beta_2$  表示  $\text{Salary}$   
 $\text{Treat}$  和  $\text{Salary}$  的系数

## Im(A1Cdiff ~ . - A1CBase - A1C52, data=z)

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	1.3845825	0.7412652	1.868	0.06216	.
HDL	7.0921082	7.4252863	0.955	0.33981	
LDL	6.7457429	7.4299154	0.908	0.36420	
TotalChol	-6.6762442	7.4283109	-0.899	0.36906	
Triglycerides	3.1720212	3.4028299	0.932	0.35154	
Creatinine	-0.0014036	0.0025733	-0.545	0.58559	
FastInsulin	-0.0042363	0.0022309	-1.899	0.05794	.
ALT	-0.0002606	0.0038294	-0.068	0.94576	
AST	-0.0038380	0.0056152	-0.684	0.49449	
GGT	0.0002968	0.0010794	0.275	0.78339	
DDuration	0.0134485	0.0086608	1.553	0.12088	
Age	0.0022026	0.0046120	0.478	0.63308	
Weight	0.0128551	0.0049309	2.607	0.00931	**
BMI	-0.0355874	0.0136901	-2.599	0.00951	**
Waist	-0.0111150	0.0056152	-1.979	0.04812	*

放进 weight  
的系数 变量

## **Im(A1Cdiff ~ . - A1CBase - A1C52, data=z)**

FastBG	<b>-0.1690230</b>	0.0237942	-7.104	2.75e-12	***
HomaS	-0.0032221	0.0014140	-2.279	0.02295	*
HomaIR	0.2318614	0.0806527	2.875	0.00415	**
HomaB	-0.0039261	0.0028481	-1.379	0.16844	
Diastolic	0.0001408	0.0051229	0.027	0.97809	
Systolic	0.0005319	0.0027948	0.190	0.84912	
Pulse	-0.0080700	0.0038454	-2.099	0.03617	*
TherapyP	<b>-0.0887768</b>	0.0716958	-1.238	<b>0.21600</b>	

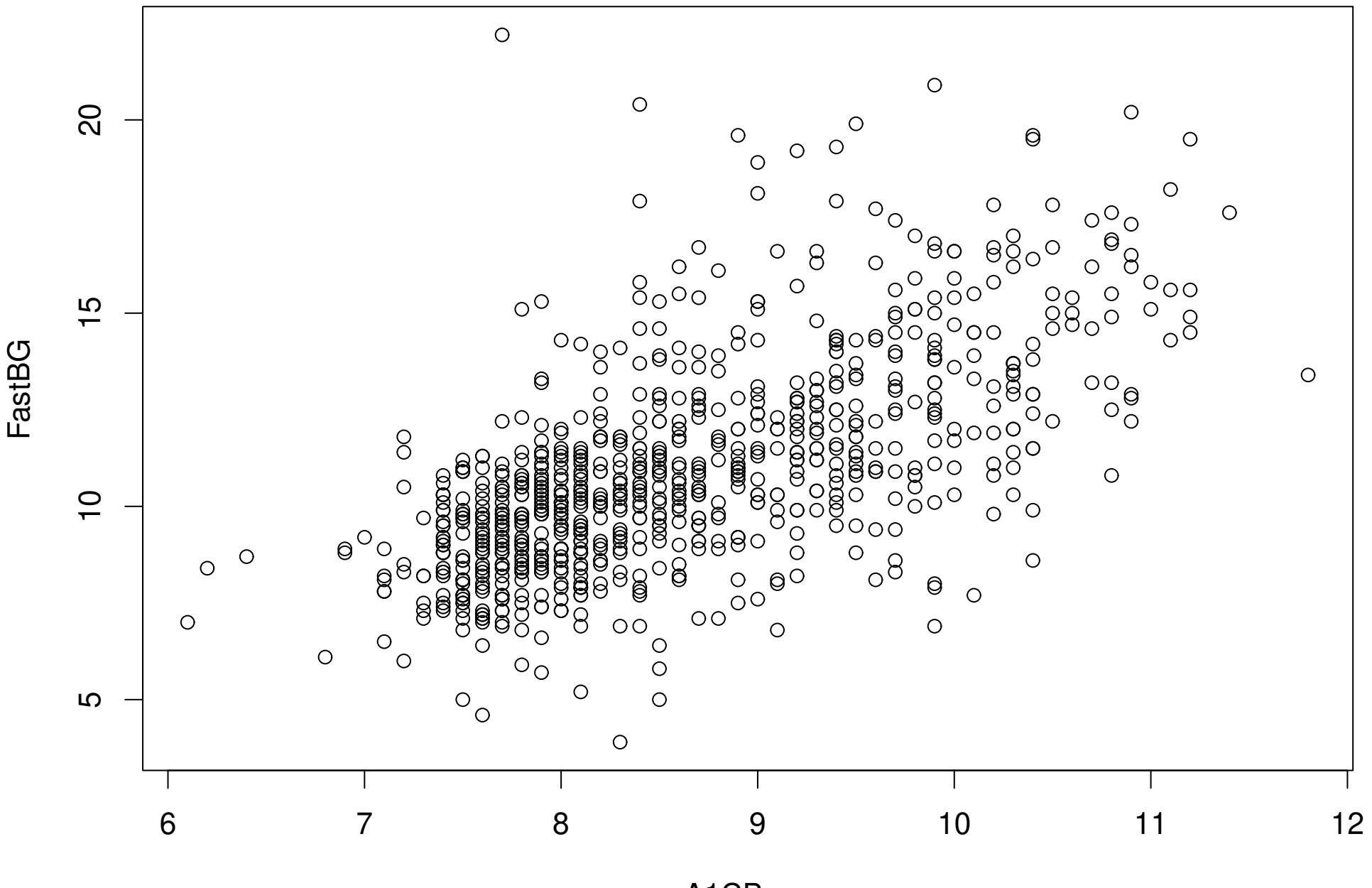
## Im(A1Cdiff ~ . - A1C52, data=z)

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	4.950e+00	6.685e-01	7.404	3.44e-13	***
HDL	7.042e+00	6.355e+00	1.108	0.2681	
LDL	6.887e+00	6.359e+00	1.083	0.2791	
TotalChol	-6.829e+00	6.357e+00	-1.074	0.2831	
Triglycerides	3.217e+00	2.912e+00	1.105	0.2697	
Creatinine	-1.312e-03	2.202e-03	-0.596	0.5517	
FastInsulin	6.488e-04	1.931e-03	0.336	0.7370	
ALT	1.283e-03	3.279e-03	0.391	0.6957	
AST	-6.220e-03	4.808e-03	-1.294	0.1961	
GGT	1.847e-05	9.239e-04	0.020	0.9841	
DDuration	6.068e-03	7.425e-03	0.817	0.4140	
Age	9.895e-04	3.948e-03	0.251	0.8021	
Weight	4.659e-03	4.248e-03	1.097	0.2731	
BMI	-1.508e-02	1.178e-02	-1.280	0.2009	
Waist	-6.533e-03	4.813e-03	-1.357	0.1751	

## Im(A1Cdiff ~ . - A1C52, data=z)

A1CBase	<b>-7.005e-01</b>	4.145e-02	-16.899	< 2e-16	***
FastBG	<b>5.209e-02</b>	2.420e-02	2.152	0.0317	*
HomaS	-1.227e-03	1.216e-03	-1.009	0.3132	
HomaIR	-3.275e-02	7.078e-02	-0.463	0.6437	
HomaB	-3.039e-04	2.447e-03	-0.124	0.9012	
Diastolic	1.569e-04	4.384e-03	0.036	0.9715	
Systolic	-3.867e-04	2.392e-03	-0.162	0.8716	
Pulse	-7.553e-03	3.291e-03	-2.295	0.0220	*
TherapyP	<b>-7.457e-02</b>	6.136e-02	-1.215	<b>0.2247</b>	

# FastBG vs A1CBase



# Linear regression with $p$ predictors

- Denote the sample by  $\{(y_i, x_{i1}, x_{i2}, \dots, x_{ip}), i = 1, 2, \dots, n\}$
- The linear model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, 2, \dots, n$$

may be written in matrix form as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

# Vector differentiation

**Theorem 1** Let  $\partial/\partial\beta$  denote a column vector with  $i$ th component  $\partial/\partial\beta_i$ . For any constant vector  $\mathbf{a}$  and symmetric matrix  $\mathbf{A}$ ,

$$\frac{\partial(\beta' \mathbf{a})}{\partial\beta} = \mathbf{a}$$
$$\frac{\partial(\beta' \mathbf{A}\beta)}{\partial\beta} = 2\mathbf{A}\beta$$

**Proof.** 1st equation is due to  $\partial(\beta' \mathbf{a})/\partial\beta_i = \underbrace{\partial(\sum_j \beta_j a_j)/\partial\beta_i}_{= a_i}$ . For the 2nd,

$$\begin{aligned}\frac{\partial(\beta' \mathbf{A}\beta)}{\partial\beta_i} &= \frac{\partial}{\partial\beta_i} \left( \sum_j \sum_k a_{jk} \beta_j \beta_k \right) \\ &= 2a_{ii}\beta_i + 2 \sum_{i \neq j} a_{ij}\beta_j \\ &= 2 \underbrace{\sum_j a_{ij}\beta_j}_{= (\mathbf{A}\beta)_i} \\ &= \underbrace{2(\mathbf{A}\beta)_i}_{= 2\mathbf{A}\beta_i}\end{aligned}$$

數學二教材

## Least squares theory

Why we don't need normal distribution here?  
Cause it only uses mean and covariance,  $\hat{\beta} = \bar{y} - \bar{X}\bar{\beta}$   
(moment)

**Theorem 2** Let  $y = X\beta + \epsilon$ , where  $y$  is an  $n \times 1$ ,  $X$  is a  $n \times k$  matrix of rank  $k$ , and  $\epsilon$  is an  $n \times 1$  vector of independent errors with mean 0 and constant variance  $\sigma^2$ . Let  $\hat{\beta}$  minimize the sum of squares

$$\|y - X\beta\|^2 = (y - X\beta)'(y - X\beta) = y'y - 2\beta'X'y + \beta'X'X\beta$$

$$= \|\hat{\epsilon}\|^2$$

and let  $SSE = \|y - X\hat{\beta}\|^2$  be the sum of squared errors. Then

$$A = (X'X)^{-1}X'y \quad A' = X(X'X)^{-1}, \quad \text{Var}(y) = \text{Var}(X\hat{\beta} + \epsilon) \\ = \text{Var}(\hat{\epsilon}) = \sigma^2 I$$

$$\text{Cov}(\hat{\beta}) = \text{Cov}(A'y) = \text{Cov}(A\hat{\epsilon}) = A \text{Cov}(\hat{\epsilon}) A' \\ = AA' \sigma^2 I = \sigma^2 (X'X)^{-1}$$

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

$$2. E(\hat{\beta}) = \beta \quad \sum_{i=1}^n c_i y_i = c'y$$

$$3. \text{Cov}(\hat{\beta}) = \sigma^2 (X'X)^{-1} \quad \sum_{i=1}^n \epsilon_i$$

$$4. s^2 = SSE/(n - k) \quad \frac{n-k}{n-k}$$

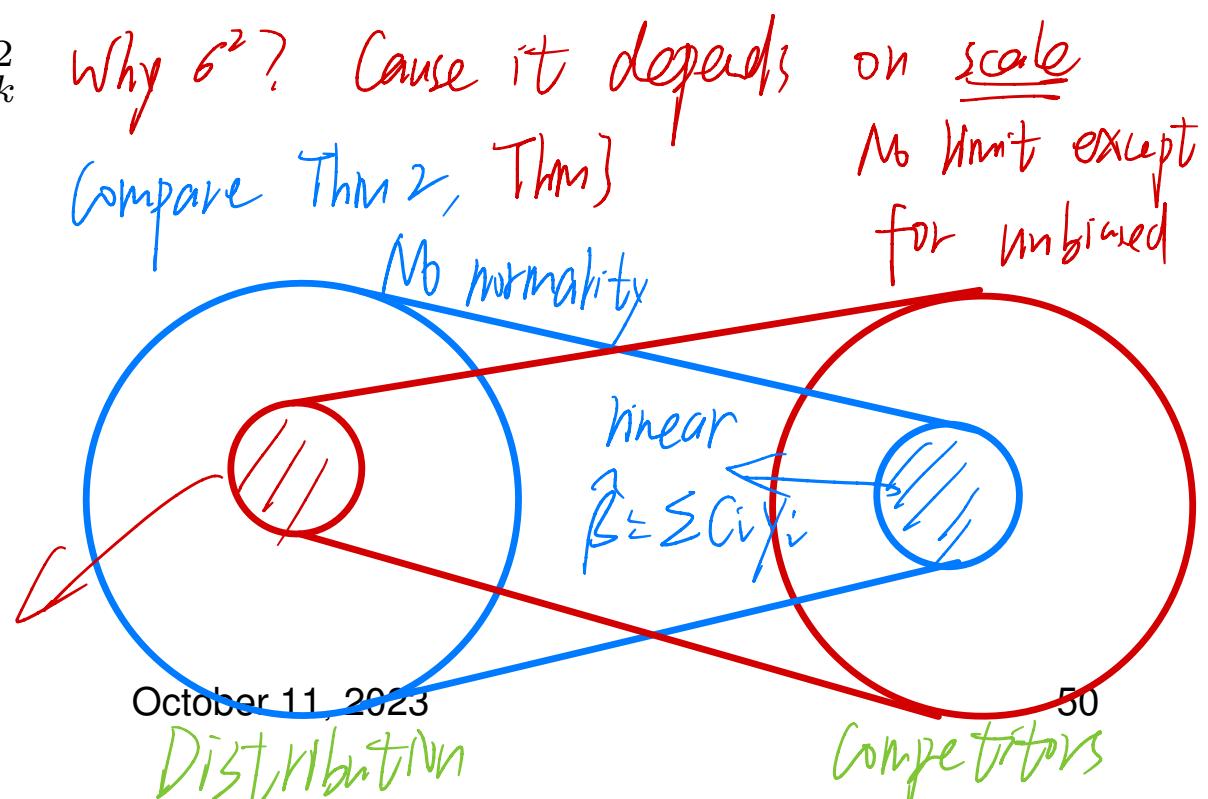
5. For any vector  $c$ ,  $c'\hat{\beta}$  is the best linear unbiased estimate (BLUE) of  $c'\beta$

**Theorem 3** Suppose in addition that  $\epsilon_1, \epsilon_2, \dots, \epsilon_n$  are normally distributed. Then

1.  $\hat{\beta}$  is the uniformly minimum variance unbiased (UMVU) estimate of  $\beta$ , i.e., it minimizes the variance among all unbiased estimates for all values of  $\beta$
2.  $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$
3.  $(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) / \sigma^2 \sim \chi_k^2$
4.  $\hat{\beta}$  is independent of  $s^2$
5.  $(n - k)s^2 / \sigma^2 \sim \chi_{n-k}^2$

$$\frac{SSE}{\sigma^2} = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sigma^2}$$

normality  
of  $\Sigma$



## Model in centered form

- Write  $\beta = (\beta_0, \beta_1, \dots, \beta_p)$
- The model

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

can be written in terms of centered  $x$ 's as

$$y_i = \alpha + \beta_1(x_{i1} - \bar{x}_1) + \beta_2(x_{i2} - \bar{x}_2) + \dots + \beta_p(x_{ip} - \bar{x}_p) + \epsilon_i$$

where  $\bar{x}_j = \underbrace{n^{-1} \sum_{i=1}^n x_{ij}}$ ,  $j = 1, 2, \dots, p$ , and

$$\alpha = \beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 + \dots + \beta_p \bar{x}_p$$

- Let  $\beta_1 = (\beta_1, \beta_2, \dots, \beta_p)'$  and define

$$\cancel{X} = \left( \begin{array}{c|c} & X_1 \end{array} \right) \quad \mathbf{X}_1 = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \quad \cancel{\neq X}$$

- Let  $\mathbf{J}$  a  $n \times n$  matrix of 1's and  $\mathbf{I}$  be the  $n \times n$  identity matrix
- Define the *centered design matrix*

$$\mathbf{X}_c = (\mathbf{I} - n^{-1}\mathbf{J})\mathbf{X}_1 = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \dots & x_{2p} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \dots & x_{np} - \bar{x}_p \end{pmatrix}$$

↑  
Centred

- Let  $\mathbf{j} = (1, 1, \dots, 1)'$  be a  $n \times 1$  vector of 1's
- The centered model can be written as

$$\mathbf{y} = (\mathbf{j}, \mathbf{X}_c) \begin{pmatrix} \alpha \\ \boldsymbol{\beta}_1 \end{pmatrix} + \boldsymbol{\epsilon}$$

with least-squares estimates (LSEs)

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\boldsymbol{\beta}}_1 \end{pmatrix} = [(\mathbf{j}, \mathbf{X}_c)'(\mathbf{j}, \mathbf{X}_c)]^{-1}(\mathbf{j}, \mathbf{X}_c)'\mathbf{y} = \begin{pmatrix} \bar{y} \\ (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{y} \end{pmatrix}$$

- For the *simple linear regression model*  $y = \alpha + \beta(x - \bar{x}) + \epsilon$ , define  $s_x^2 = (n-1)^{-1} \sum_{i=1}^n (x_i - \bar{x})^2$ ,  $s_y^2 = (n-1)^{-1} \sum_{i=1}^n (y_i - \bar{y})^2$ , and  $r = (n-1)^{-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) / (s_x s_y)$ . Then

$$\hat{\alpha} = \bar{y} \quad \hat{\beta} = r s_y / s_x$$

$$\frac{L_{xy}}{L_{xx}}$$

# Sums of squares

↓ Total

- Define the corrected total sum of squares  $SST = \sum_i (y_i - \bar{y})^2$  and the regression sum of squares

$$SSR = \|\mathbf{X}_c \hat{\beta}_1\|^2 = (\mathbf{X}_c \hat{\beta}_1)' (\mathbf{X}_c \hat{\beta}_1) = \hat{\beta}_1' \mathbf{X}_c' \mathbf{X}_c \hat{\beta}_1 = \hat{\beta}_1' \mathbf{X}_c' \mathbf{y}$$

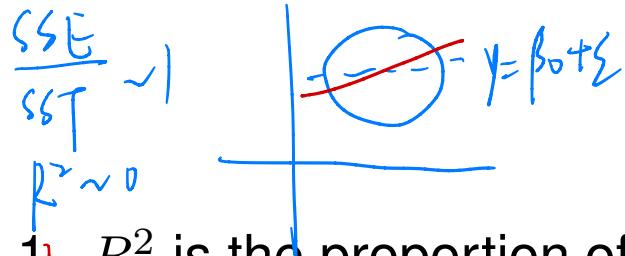
(recall that  $\hat{\beta}_1 = (\mathbf{X}_c' \mathbf{X}_c)^{-1} \mathbf{X}_c' \mathbf{y}$ )

- The error sum of squared errors can be expressed as

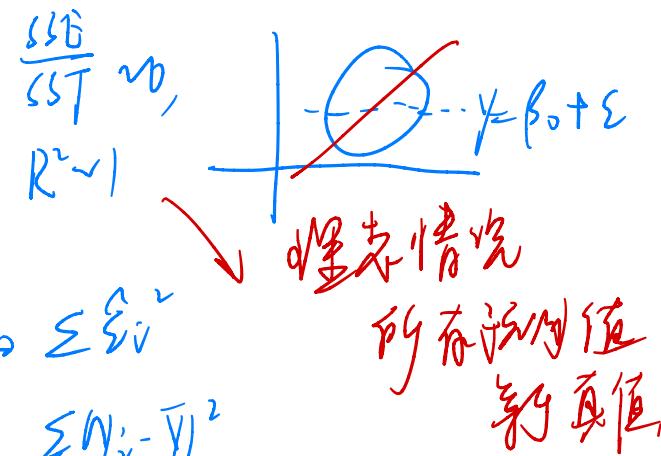
$$\begin{aligned} SSE &= \|(\mathbf{y} - \mathbf{X}\hat{\beta})\|^2 = \|(\mathbf{y} - \bar{y}\mathbf{j} - \mathbf{X}_c \hat{\beta}_1)\|^2 \\ &= \sum_i (y_i - \bar{y})^2 + \|\mathbf{X}_c \hat{\beta}_1\|^2 - 2(\mathbf{y} - \bar{y}\mathbf{j})' \mathbf{X}_c \hat{\beta}_1 \\ &= SST + SSR - 2(\mathbf{y}' \mathbf{X}_c \hat{\beta}_1 + \bar{y}\mathbf{j}' \mathbf{X}_c \hat{\beta}_1) \\ &= SST + SSR - 2\hat{\beta}_1' \mathbf{X}_c' \mathbf{y} = SST - SSR \end{aligned}$$

because  $SSR = \hat{\beta}_1' \mathbf{X}_c' \mathbf{y}$  and  $\mathbf{j}' \mathbf{X}_c = \mathbf{0}'$

- Therefore  $SST = SSR + SSE$



## Properties of $R^2$



- $R^2$  is the proportion of SST due to regression

这直线表示简单的线性回归只有  
y值的平均值。

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \rightarrow \frac{\sum \hat{e}_i^2}{\sum (y_i - \bar{y})^2}$$

- Adjusted  $R^2$  is  $R_{adj}^2 = 1 - \frac{SSE/(n-1-p)}{SST/(n-1)}$

- $0 \leq R^2 \leq 1$  and  $R^2 = 0$  if  $\hat{\beta}_1 = \hat{\beta}_2 = \dots = \hat{\beta}_p = 0$  ( $\hat{\beta}_0$  need not be 0)

- Let  $\bar{y} = n^{-1} \sum_i \hat{y}_i$ . The positive square root  $R = \sqrt{R^2}$  is called the *multiple correlation coefficient* and equals the correlation between the  $y_i$  and  $\hat{y}_i$

$$R = \frac{\sum_i (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_i (y_i - \bar{y})^2 \sum_i (\hat{y}_i - \bar{\hat{y}})^2}}$$

7. makes residuals be less i.e.  $SSE \downarrow \Rightarrow R^2 = 1 - \frac{SSE}{SST} \uparrow$

- $R^2$  increases if a variable is added to the model but  $R_{adj}^2$  can decrease

- $E(R^2) = p(n-1)^{-1}$  if  $\beta_1 = \beta_2 = \dots = \beta_p = 0$

# Expectation and covariance operators

## Definition 1

- Let  $\mathbf{Z}$  be an  $m \times n$  matrix of random variables  $Z_{ij}$
- Then  $E(\mathbf{Z})$  is a  $m \times n$  matrix of expected values  $E(Z_{ij})$

**Theorem 4** If  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  are  $l \times m$ ,  $n \times p$ , and  $l \times p$  matrices of constants,

$$E(\mathbf{A}\mathbf{Z}\mathbf{B} + \mathbf{C}) = \mathbf{A}E(\mathbf{Z})\mathbf{B} + \mathbf{C}$$

**Theorem 5** If  $\mathbf{A}$  and  $\mathbf{B}$  are  $m \times n$  matrices of constants and  $\mathbf{X}$  and  $\mathbf{Y}$  are  $n \times 1$  random variables, then

$$E(\mathbf{A}\mathbf{X} + \mathbf{B}\mathbf{Y}) = \mathbf{A}E(\mathbf{X}) + \mathbf{B}E(\mathbf{Y})$$

## Definition 2

1. If  $\mathbf{X}$  and  $\mathbf{Y}$  are  $m \times 1$  and  $n \times 1$  random vectors,  $\text{Cov}(\mathbf{X}, \mathbf{Y})$  is a  $m \times n$  matrix with  $(i, j)$ th element  $\text{Cov}(X_i, Y_j) = E\{(X_i - EX_i)(Y_j - EY_j)\}$
2.  $\text{Var}(\mathbf{X}) = \text{Cov}(\mathbf{X}, \mathbf{X})$  with  $(i, j)$ th element  $E\{(X_i - EX_i)(X_j - EX_j)\}$

**Theorem 6**  $\text{Cov}(\mathbf{X}, \mathbf{Y}) = E\{(\mathbf{X} - E\mathbf{X})(\mathbf{Y} - E\mathbf{Y})'\}$

**Theorem 7** If  $\mathbf{X}$  and  $\mathbf{Y}$  are  $m \times 1$  and  $n \times 1$  random vectors and  $\mathbf{A}$  and  $\mathbf{B}$  are  $l \times m$  and  $p \times n$  matrices of constants,

$$\text{Cov}(\mathbf{AX}, \mathbf{BY}) = \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{Y}) \mathbf{B}'$$

**Theorem 8** If  $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $\mathbf{U}$ , and  $\mathbf{V}$  are random vectors of the same dimension, then for all real numbers  $a, b, c$ , and  $d$

$$\begin{aligned}\text{Cov}(a\mathbf{X} + b\mathbf{Y}, c\mathbf{U} + d\mathbf{V}) &= ac\text{Cov}(\mathbf{X}, \mathbf{U}) + ad\text{Cov}(\mathbf{X}, \mathbf{V}) \\ &\quad + bc\text{Cov}(\mathbf{Y}, \mathbf{U}) + bd\text{Cov}(\mathbf{Y}, \mathbf{V})\end{aligned}$$

**Theorem 9** Let  $\mathbf{X}$  be an  $n \times 1$  random vector and let  $\mathbf{A}$  be an  $n \times n$  symmetric matrix of constants. Then

$$E(\mathbf{X}' \mathbf{A} \mathbf{X}) = \text{tr}(\mathbf{A} \Sigma) + \theta' \mathbf{A} \theta$$

where  $\theta = \underline{\mathbf{E}(\mathbf{X})}$  and  $\Sigma = \underline{\text{Var}(\mathbf{X})}$

~~difference~~

$$\hat{\sigma}_1^2 = \frac{\sum (X_i - \bar{X})^2}{n-1} \quad \begin{matrix} \bar{X} \text{ always known} \\ \Leftarrow \text{unbiased} \end{matrix}$$

$$\hat{\sigma}_2^2 = \frac{\sum (X_i - \mu)^2}{n} \quad , \quad E(\hat{\sigma}_2^2) = \frac{1}{n} \sum E(X_i - \mu)^2 = \frac{1}{n} \sum \text{Var}(X_i), \text{ still unbiased}$$

## Noncentral chi-square distribution

**Definition 3** Let  $Z_i$  be independent  $N(\theta_i, 1)$ ,  $i = 1, 2, \dots, p$ . Then

1.  $X = \sum_i Z_i^2$  has a noncentral  $\chi_p^2(\lambda)$  with  $p$  df and noncentrality parameter  $\lambda = \sum_i \theta_i^2$
2. The mean and variance of  $X$  are  $(p + \lambda)$  and  $2(p + 2\lambda)$ , respectively
3. The characteristic function of  $X$  is

$$Ee^{itX} = \frac{\exp\{i\lambda t/(1 - 2it)\}}{(1 - 2it)^{\nu/2}}$$

If  $U_1 \sim \chi^2_{p_1}$ ,  $U_2 \sim \chi^2_{p_2}(\lambda)$ .

$$\frac{U_2/p_2}{U_1/p_1} \sim F_{p_2, p_1}(\lambda), \quad \frac{U_1/p_1}{U_2/p_2} \sim \frac{1}{F_{p_2, p_1}(\lambda)}$$

## Noncentral F distribution

← non centred

**Definition 4** Let  $U_1 \sim \chi^2_{p_1}(\lambda)$  be independent of  $U_2 \sim \chi^2_{p_2}$ . Then

$$\frac{U_1/p_1}{U_2/p_2} \sim F_{p_1, p_2}(\lambda) \quad \uparrow \text{centred } \chi^2$$

a noncentral F distribution with  $p_1, p_2$  degrees of freedom and noncentrality parameter  $\lambda$

**Theorem 10** Suppose  $\mathbf{y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ , where  $\mathbf{X}$  is  $n \times (p + 1)$  of rank  $p + 1 < n$  and let

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

and  $\lambda = \boldsymbol{\beta}'_1 \mathbf{X}'_c \mathbf{X}_c \boldsymbol{\beta}_1 / \sigma^2$ . Then

1.  $SSR/\sigma^2 \sim \chi_p^2(\lambda)$
2.  $SSE/\sigma^2 \sim \chi_{n-p-1}^2$
3.  $SSR$  and  $SSE$  are independent
4.  $F \sim F_{p, n-p-1}(\lambda)$  (noncentral  $F$  distribution)

△ 62-75

## Testing a subset of $\beta$ s

Reject  $H_0$  if

$$F > F_{n, n-p_1, \alpha}$$

- Consider testing  $H_0 : \beta_2 = \mathbf{0}$  vs  $H_1 : \beta_2 \neq \mathbf{0}$  in the model

$$\mathbf{y} = \mathbf{X}\beta + \epsilon = (\mathbf{X}_1, \mathbf{X}_2) \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \epsilon = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \epsilon$$

where  $\beta_2$  is  $h \times 1$  and the intercept term  $\beta_0$  is included in  $\beta_1$

- Let  $SSE_0$  and  $SSE_1$  be the residual sums of squares under  $H_0$  and  $H_1$ , respectively

$$v_1 = h - (p+1)$$

- Then

$$v_0 = n - (p+1) + h$$

more degree

$$F = \frac{(SSE_0 - SSE_1)/h}{SSE_1/(n - p - 1)} \sim F_{h, n-p-1}(\lambda)$$

where independent

$$\lambda = \frac{\beta_2' \{ \mathbf{X}_2' \mathbf{X}_2 - \mathbf{X}_2' \mathbf{X}_1 (\mathbf{X}_1 \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2 \} \beta_2}{\sigma^2}$$

$$SSE_0 = (SSE_0 - SSE_1) + SSE_1$$

$$E\left(\frac{SSE_0}{n-p-1+h}\right) = \sigma^2 \text{ under } H_0$$

$$E\left(\frac{SSE_1}{n-p-1}\right) = \sigma^2 \text{ under } H_1$$

dependent

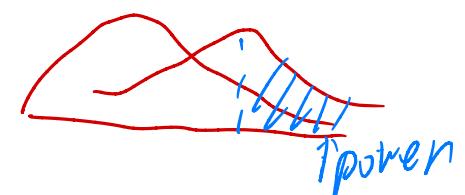
$$\frac{(SSE_0 - SSE_1)/h}{SSE_0/(n-p-1+h)}$$

wrong !!

$$\begin{aligned}\text{Power} &= \Pr[\text{Reject } H_0 \mid H_0 \text{ is fake}] \\ &= \Pr[F > F_{h, n-p-1, \alpha} \mid H_0 \text{ is fake}]\end{aligned}$$

$\hookrightarrow$   $H_0$  should be rejected

$\nearrow$  power ↑ more data,  $n$   
 $\curvearrowright$  what happens to  $H_1$   
 $\rightarrow$   $\uparrow$



**Theorem 11** Let  $y \sim N(\mu, V)$ ,  $r$  be the rank of a constant matrix  $A$ , and

$\lambda = \mu' A \mu$ . Then

$$\begin{matrix} |X|_n & |Y|_n & |A|_r \\ y' A y \sim \chi_r^2(\lambda) \end{matrix}$$

if and only if  $AV$  is idempotent (i.e.,  $AVAV = AV$ )

$F_{h, n-p-1, \alpha}$   
 $n$  double

what will happen

**Corollary 1** Let  $y \sim N(\mu, V)$  with  $V$  of rank  $n$ . Then  $y' V^{-1} y \sim \chi_n^2(\lambda)$  with  
 $\lambda = \mu' V^{-1} \mu$ .

Test starts from  $\mathbf{C}\hat{\beta}$

## Testing $H_0 : \mathbf{C}\beta = 0$

- Suppose  $\mathbf{y} \sim N_n(\mathbf{X}\beta, \sigma^2\mathbf{I})$  and  $\mathbf{C}$  is  $q \times (p+1)$  of rank  $q \leq p+1$ . Then

$$\underline{\mathbf{C}\hat{\beta}} \sim N_q[\mathbf{C}\beta, \sigma^2 \underline{\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'}]$$

- Define

$$\begin{aligned} \text{SSH} &= (\mathbf{C}\hat{\beta})' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} \mathbf{C}\hat{\beta} \\ \lambda &= (\mathbf{C}\beta)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} \mathbf{C}\beta / \sigma^2 \end{aligned}$$

← Using Thm 11 up page  
Corollary 1

- Then

- $\text{SSH}/\sigma^2 \sim \chi_q^2(\lambda)$
- $\text{SSE}/\sigma^2 \sim \chi_{n-p-1}^2$
- SSH and SSE are independent and

$$SSA = \frac{n \|\mathbf{C}\hat{\beta}\|^2}{\text{Var}(\mathbf{C}\hat{\beta})}$$

$$F = \frac{\text{SSH}/q}{\text{SSE}/(n-p-1)} \sim F_{q,n-p-1}(\lambda)$$

## Testing $H_0 : \mathbf{C}\beta = \mathbf{t}$

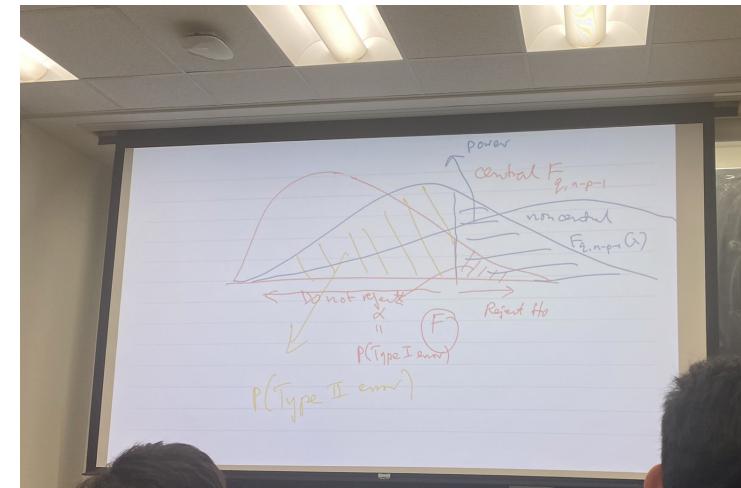
Suppose  $\mathbf{y} \sim N_n(\mathbf{X}\beta, \sigma^2\mathbf{I})$  and  $\mathbf{C}$  is  $q \times (p+1)$  of rank  $q \leq p+1$ . Then  $\mathbf{C}\hat{\beta} - \mathbf{t} \sim N_q[\mathbf{C}\beta - \mathbf{t}, \sigma^2\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']$ . Define

$$\begin{aligned} \text{SSH} &= (\mathbf{C}\hat{\beta} - \mathbf{t})' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\mathbf{C}\hat{\beta} - \mathbf{t}) \\ \lambda &= (\mathbf{C}\beta - \mathbf{t})' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\mathbf{C}\beta - \mathbf{t}) / \sigma^2 \end{aligned} \quad (1)$$

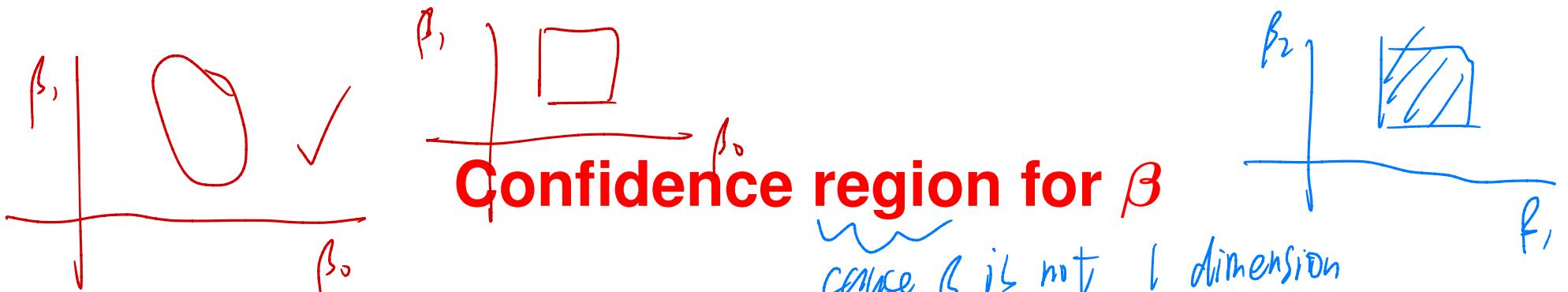
Then

$$\text{SSH} = \frac{\|\mathbf{C}\hat{\beta}\|^2}{\text{Var}(\mathbf{C}\hat{\beta})}$$

1.  $\text{SSH}/\sigma^2 \sim \chi_q^2(\lambda)$
2.  $\text{SSE}/\sigma^2 \sim \chi_{n-p-1}^2$
3. SSH and SSE are independent and



$$F = \frac{\text{SSH}/q}{\text{SSE}/(n-p-1)} = \frac{\text{SSH}/q}{s^2} \sim F_{q,n-p-1}(\lambda) \quad (2)$$



- Recall that  $\beta$  is  $(p + 1) \times 1$
- Put  $C = I$  and  $t = \beta$  in (1), so that  $q = p + 1$
- It follows from (2) that

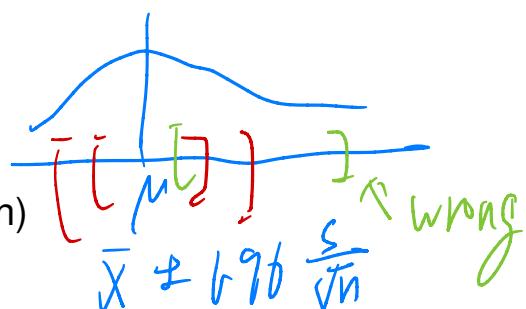
$$W = s^{-2}(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) / (p + 1) \sim F_{p+1, n-p-1}$$

- Therefore

$$P\{s^{-2}(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) / (p + 1) \leq F_{p+1, n-p-1; \alpha}\} = 1 - \alpha$$

and a  $100(1 - \alpha)\%$  confidence region for  $\beta$  is

$$(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) \leq (p + 1)s^2 F_{p+1, n-p-1; \alpha}$$



check  $|\bar{x} - \mu| \leq 1.96 \frac{s}{\sqrt{n}}$

October 11, 2023

# 95% confidence region for $(\beta_0, \beta_1)$ from fitting

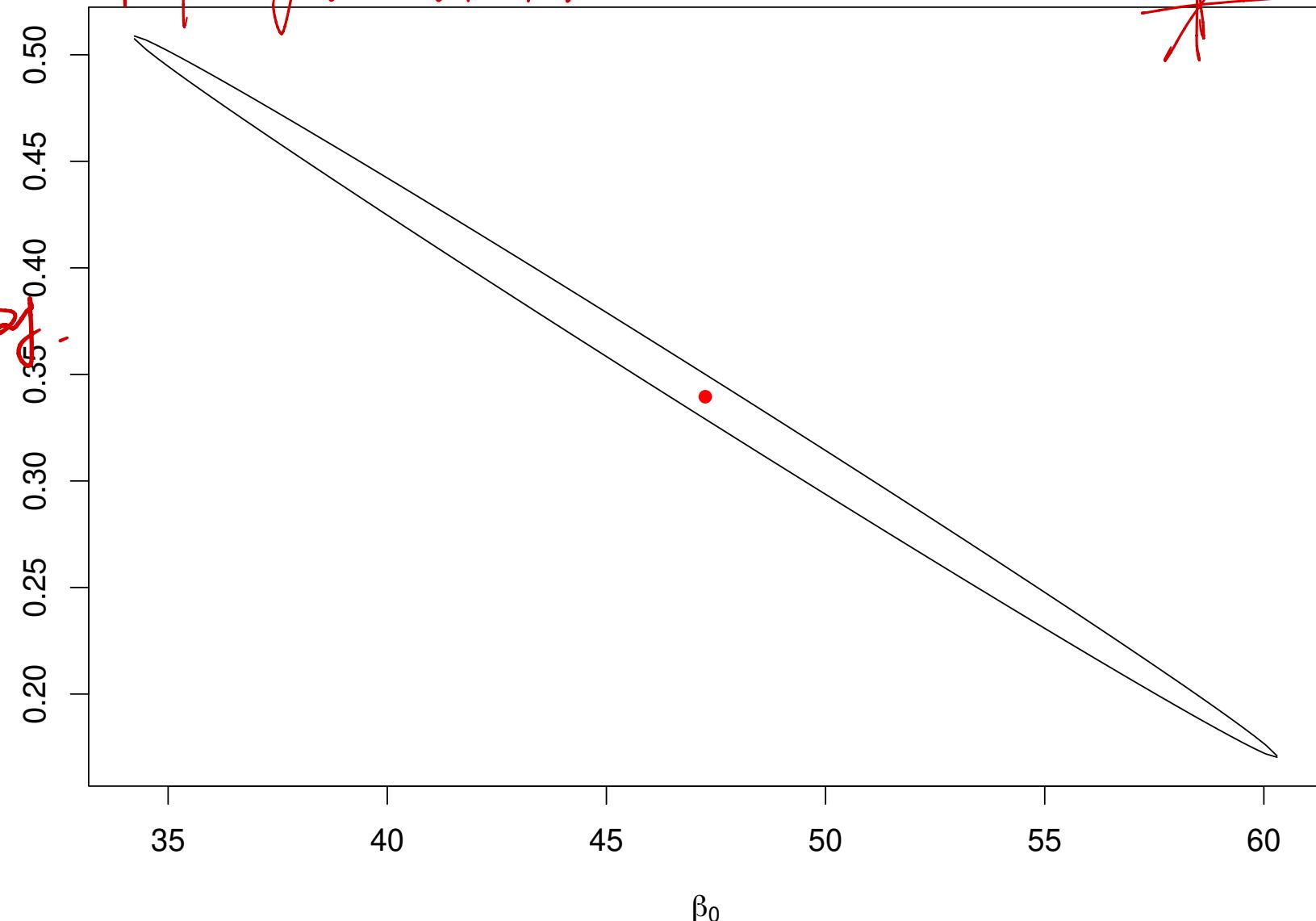
$$\text{cov}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

height vs reach in boxer data

$\beta_0, \beta_1$  negative correlated



还有一个结论  
是关于系数  
 $\beta_0$  和  $\beta_1$   
的范围不等式



## Confidence intervals for $\beta_j$ and $\mathbf{a}'\boldsymbol{\beta}$

- Recall from Theorem 3 that  $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$
- Let  $G = (\mathbf{X}'\mathbf{X})^{-1}$  and  $g_{ij}$  be the  $(i, j)$ th element of  $G$
- Then  $\hat{\beta}_j \sim N(\beta_j, \sigma^2 g_{jj})$  and  $(\hat{\beta}_j - \beta_j)s^{-1}/\sqrt{g_{jj}} \sim t_{n-p-1}$
- A  $100(1 - \alpha)\%$  confidence interval for  $\beta_j$  is
- Similarly a  $100(1 - \alpha)\%$  confidence interval for  $\mathbf{a}'\boldsymbol{\beta}$  is

7. 从圖中得知

$$\hat{\beta}_j \pm s\sqrt{g_{jj}} t_{n-p-1; \alpha/2}$$

$$\mathbf{a}'\hat{\boldsymbol{\beta}} \pm s\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}} t_{n-p-1; \alpha/2}$$

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

$$\begin{aligned} & \frac{\hat{\beta}_j - \beta_j}{\text{SD}(\hat{\beta}_j - \beta_j)} \\ & \text{SD}(\hat{\beta}_j - \beta_j) \\ & = \hat{s} \sqrt{g_{jj}} \end{aligned}$$

(3)

"A" means  
"estimate me"

## Confidence interval for $E(y | \mathbf{x} = \mathbf{x}_0)$

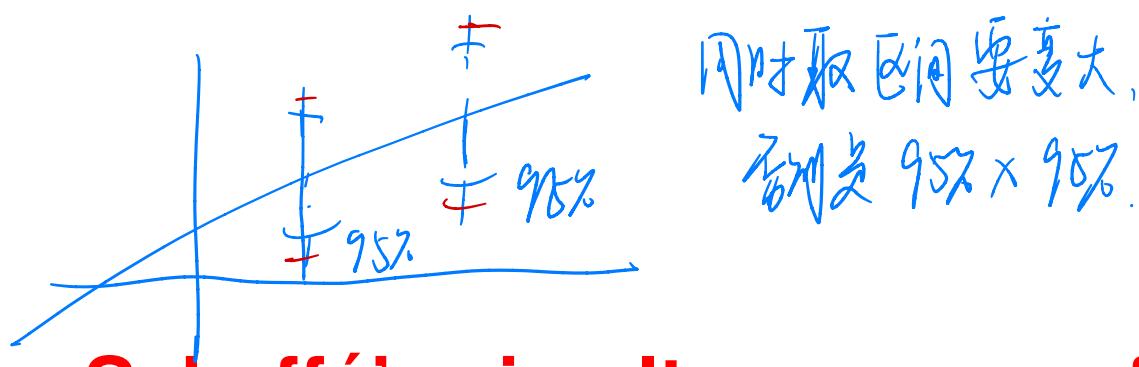
- Since  $E(y | \mathbf{x} = \mathbf{x}_0) = \mathbf{x}'_0 \boldsymbol{\beta}$ , an estimate of  $E(y | \mathbf{x} = \mathbf{x}_0)$  is

$$\widehat{E(y | \mathbf{x}_0)} = \mathbf{x}'_0 \hat{\boldsymbol{\beta}}$$

- Putting  $\mathbf{a} = \mathbf{x}_0$  in (3) yields the  $100(1 - \alpha)\%$  interval for  $E(y | \mathbf{x} = \mathbf{x}_0)$

$$\mathbf{x}'_0 \hat{\boldsymbol{\beta}} \pm s \sqrt{\mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0} t_{n-p-1; \alpha/2}$$

$$P \left( \text{---} \downarrow \text{---} \ni \mathbf{x}'_0 \hat{\boldsymbol{\beta}} \right) = 1 - \alpha$$



## Scheffé's simultaneous confidence band

for  $\widehat{E(y|x)} = \mathbf{x}'\hat{\boldsymbol{\beta}}$  for all  $\mathbf{x}$

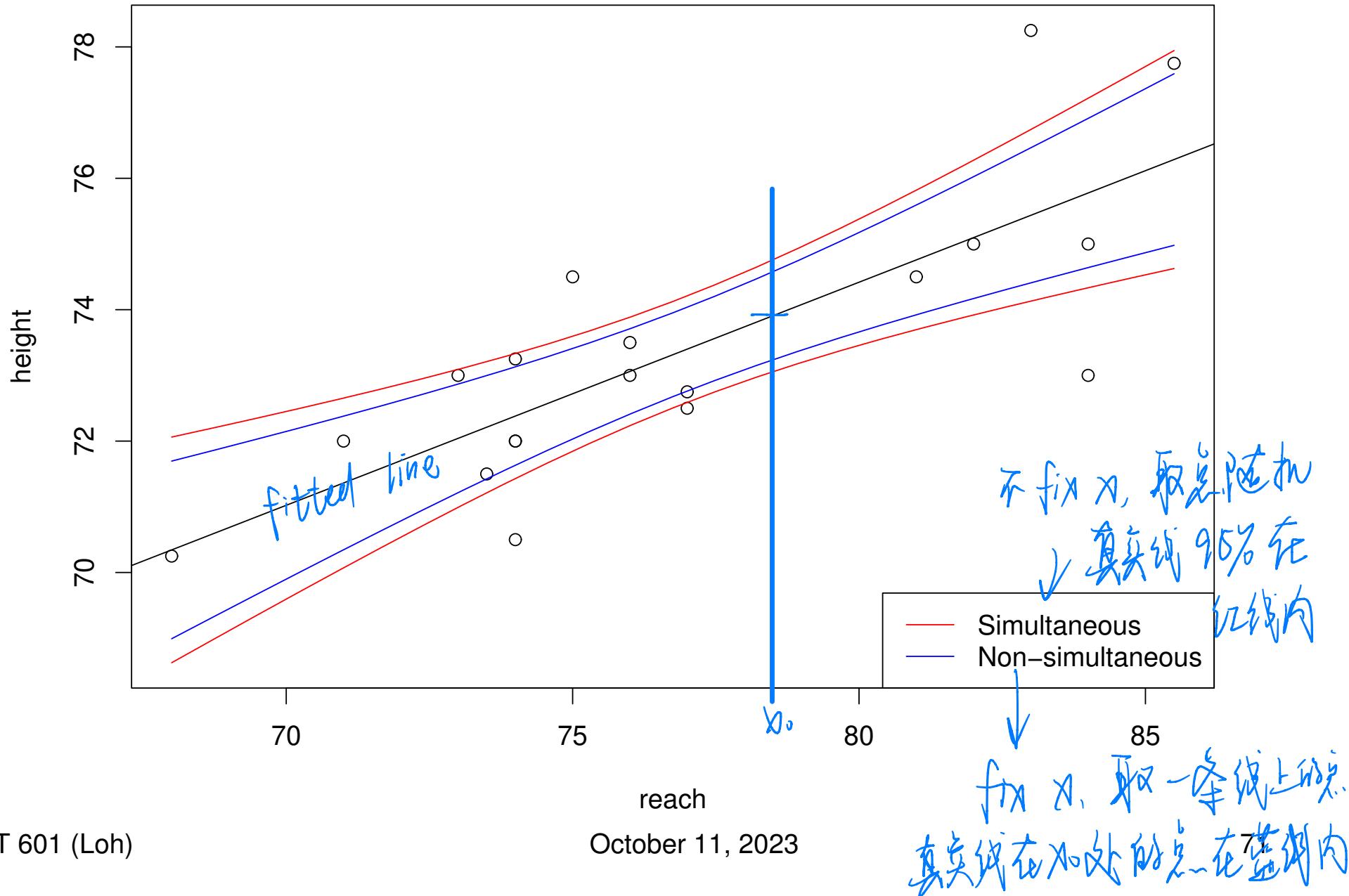
- Scheffé showed that

$$P \left( \mathbf{x}'\boldsymbol{\beta} \in \mathbf{x}'\hat{\boldsymbol{\beta}} \pm s \sqrt{(p+1)\mathbf{x}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}} F_{p+1,n-p-1;\alpha} \quad \text{for all } \mathbf{x} \right) = 1 - \alpha$$

- Hence a  $100(1 - \alpha)\%$  simultaneous confidence band for  $E(y|x) = \mathbf{x}'\boldsymbol{\beta}$  is

$$\mathbf{x}'\hat{\boldsymbol{\beta}} \pm s \sqrt{(p+1)\mathbf{x}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}} F_{p+1,n-p-1;\alpha}$$

# 95% simultaneous and non-simultaneous confidence bands for boxer data



## Homework 2 (due 11:59 pm Tue Oct 10 in Canvas)

1. Reproduce the plot of 95% simultaneous confidence band for  $E(y | \mathbf{x})$ .
2. Show that for testing  $H_0 : \beta_2 = 0$  vs  $H_1 : \beta_2 \neq 0$  on page 62, the value of

$$F = \frac{(\text{SSE}_0 - \text{SSE}_1)/h}{\text{SSE}_1/(n - p - 1)}$$

is the same as the value of the statistic

$$F = \frac{\text{SSH}/h}{\text{SSE}/(n - p - 1)}$$

on page 64

## Hints for HW 2 Problem #2

Let  $\hat{\beta}_r$  and  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  denote the LSEs of  $\beta$  under  $H_0$  (reduced model) and  $H_1$  (full model), respectively. Then

$$\text{SSE}_0 - \text{SSE}_1 = \text{SSR}_1 - \text{SSR}_0 = \hat{\beta}'\mathbf{X}'\mathbf{y} - \hat{\beta}_r'\mathbf{X}'\mathbf{y} = (\hat{\beta} - \hat{\beta}_r)'\mathbf{X}'\mathbf{y}$$

Prove that this is equal to  $\text{SSH} = (\mathbf{C}\hat{\beta})'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}\mathbf{C}\hat{\beta}$  by showing that

$$\hat{\beta}_r = \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}\mathbf{C}\hat{\beta}$$

which can be done using Lagrange multipliers to minimize

$$f = \underbrace{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)}_{+ \lambda'(\mathbf{C}\beta - \mathbf{0})}$$

with respect to  $\beta$  and  $\lambda$  as follows:

1. Differentiate  $f$  with respect to  $\lambda$  and  $\beta$  and set the results to 0 to obtain two equations
2. Eliminate  $\lambda$  from the equations and solve for  $\hat{\beta}_r$

# Prediction interval for a future observation

预测误差

- Let  $y_0$  be a future value of  $y$  at  $\mathbf{x} = \mathbf{x}_0$  and its predicted value be  $\hat{y}_0 = \mathbf{x}'_0 \hat{\beta}$
- By independence of  $y_0$  and  $\hat{y}_0$ ,

~~独立于~~  $y_0$  与  $\hat{y}_0$  独立

$$\begin{aligned}\text{var}(y_0 - \hat{y}_0) &= \text{var}(y_0) + \text{var}(\hat{y}_0) = \text{var}(\epsilon) + \text{var}(\mathbf{x}'_0 \hat{\beta}) \\ &= \sigma^2 + \sigma^2 \mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0 \\ &= \sigma^2 \{1 + \mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0\}\end{aligned}$$

- Hence, because  $y_0 - \hat{y}_0$  is normally distributed with mean 0,

$$\frac{y_0 - \hat{y}_0}{s \sqrt{1 + \mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0}} \sim t_{n-p-1}$$

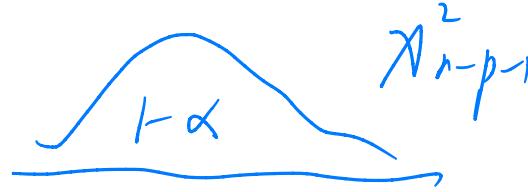
yielding  $P(|y_0 - \hat{y}_0| \leq s \sqrt{1 + \mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0} t_{n-p-1; \alpha/2}) = 1 - \alpha$

- Therefore a  $100(1 - \alpha)\%$  prediction interval for  $y_0$  is

$$\mathbf{x}'_0 \hat{\beta} \pm s \sqrt{1 + \mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0} t_{n-p-1; \alpha/2}$$

~~Q~~ ~~这个公式~~

Compare to page 69



## Confidence interval for $\sigma^2$

$$s^2 \leq \frac{\text{SSE}}{n-p-1}$$

- From Theorem 3,  $(n - p - 1)s^2/\sigma^2 = \text{SSE}/\sigma^2 \sim \chi^2_{n-p-1}$
- So

$$P\{\chi^2_{n-p-1;1-\alpha/2} \leq (n - p - 1)s^2/\sigma^2 \leq \chi^2_{n-p-1;\alpha/2}\} = 1 - \alpha$$

and a  $100(1 - \alpha)\%$  interval for  $\sigma^2$  is

$$\frac{(n - p - 1)s^2}{\chi^2_{n-p-1;\alpha/2}} \leq \sigma^2 \leq \frac{(n - p - 1)s^2}{\chi^2_{n-p-1;1-\alpha/2}}$$

$$\text{cov}(\mathbf{y}) = \begin{pmatrix} 6^2 & 26^2 \\ 26^2 & 6^2 \end{pmatrix}$$

## Generalized least squares

$\sigma^2 \mathbf{I}$

- Let  $\mathbf{y} = \mathbf{X}\beta + \epsilon$ ,  $E(\mathbf{y}) = \mathbf{X}\beta$  and  $\text{cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$ , where  $\mathbf{X}$  is a  $n \times k$  matrix of rank  $k$  and  $\mathbf{V}$  is a known positive definite matrix
- The generalized least squares estimate  $\hat{\beta}$  minimizes

$$(\mathbf{y} - \mathbf{X}\beta)' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta)$$

Definition

**Theorem 12** Let  $\mathbf{y} = \mathbf{X}\beta + \epsilon$ ,  $E(\mathbf{y}) = \mathbf{X}\beta$  and  $cov(\mathbf{y}) = \sigma^2\mathbf{V}$ , where  $\mathbf{X}$  is a  $n \times k$  matrix of rank  $k$  and  $\mathbf{V}$  is a known positive definite matrix. Then

1. Best linear unbiased estimate (BLUE) of  $\beta$  is

$$\hat{\beta} = (\mathbf{X}'\underline{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\underline{\mathbf{V}}^{-1}\mathbf{y}$$

and

$$cov(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}$$

2. An unbiased estimate of  $\sigma^2$  is

$$\begin{aligned}s^2 &= \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})}{n - k} \\ &= \frac{\mathbf{y}'\{\mathbf{V}^{-1} - \mathbf{V}^{-1}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\}\mathbf{y}}{n - k}\end{aligned}$$

3. If  $\mathbf{y}$  is  $N(\mathbf{X}\beta, \sigma^2\mathbf{V})$ ,  $\hat{\beta}$  is also the MLE

## Two proofs of BLUE

1. Differentiate  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$
2. (a) Since  $\mathbf{V}$  is positive definite,  $\mathbf{V} = \mathbf{P}\mathbf{P}'$  for a  $n \times n$  nonsingular matrix  $\mathbf{P}$   
(b) Multiply the model equation to get the model

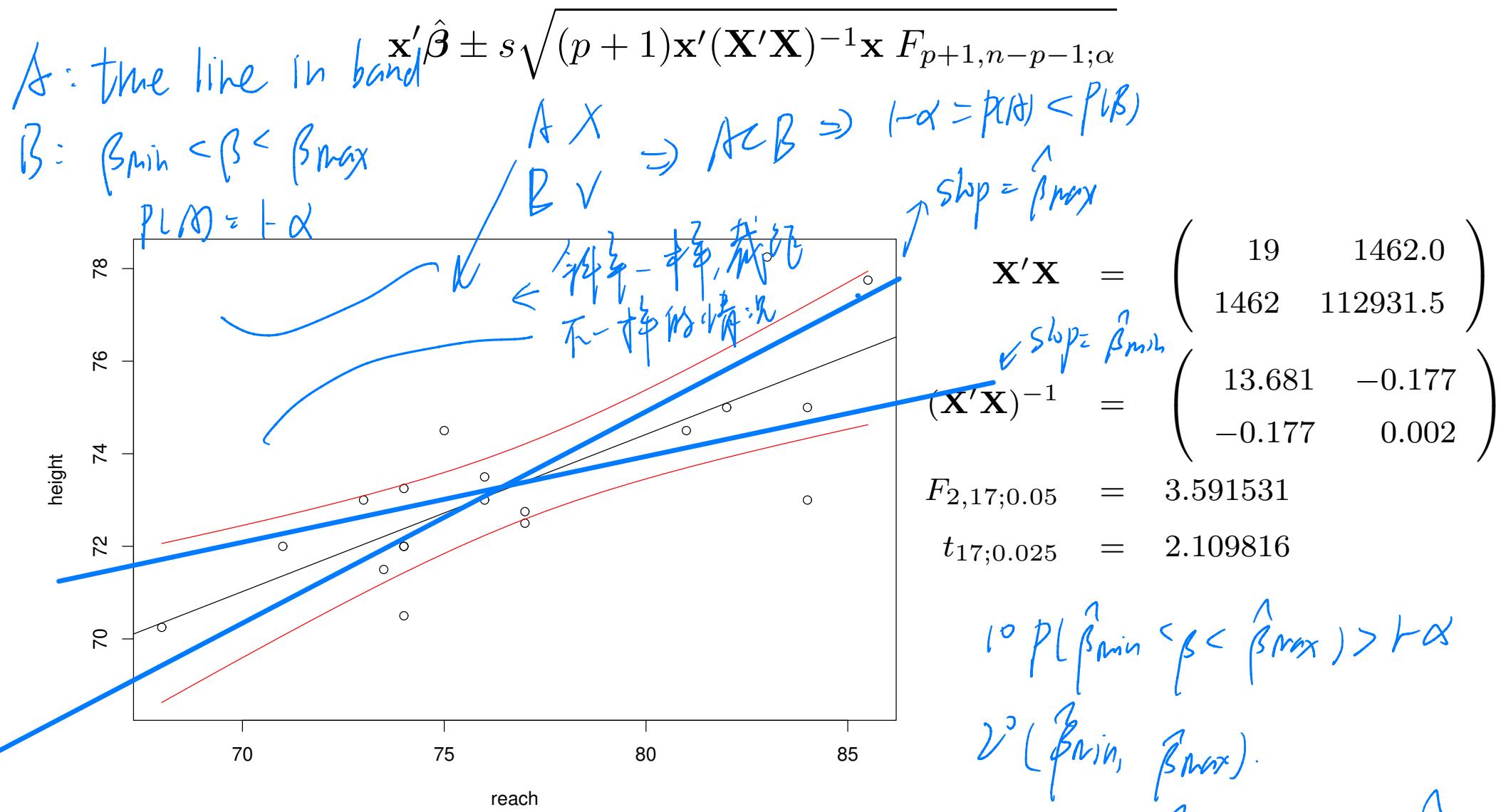
$$\mathbf{P}^{-1}\mathbf{y} = \mathbf{P}^{-1}\mathbf{X}\boldsymbol{\beta} + \mathbf{P}^{-1}\boldsymbol{\epsilon}$$

- (c) Then  $E(\mathbf{P}^{-1}\boldsymbol{\epsilon}) = \mathbf{0}$  and

$$\begin{aligned}\text{Cov}(\mathbf{P}^{-1}\boldsymbol{\epsilon}) &= \mathbf{P}^{-1}\text{Cov}(\boldsymbol{\epsilon})(\mathbf{P}^{-1})' = \mathbf{P}^{-1}\sigma^2\mathbf{V}(\mathbf{P}^{-1})' \\ &= \sigma^2\mathbf{P}^{-1}\mathbf{P}\mathbf{P}'(\mathbf{P}')^{-1} = \sigma^2\mathbf{I}\end{aligned}$$

- (d) From Theorem 2, the BLUE is

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= [(\mathbf{P}^{-1}\mathbf{X})'(\mathbf{P}^{-1}\mathbf{X})]^{-1}(\mathbf{P}^{-1}\mathbf{X})'\mathbf{P}^{-1}\mathbf{y} \\ &= [\mathbf{X}'(\mathbf{P}^{-1})'\mathbf{P}^{-1}\mathbf{X}]^{-1}\mathbf{X}'(\mathbf{P}^{-1})'\mathbf{P}^{-1}\mathbf{y} \\ &= [\mathbf{X}'(\mathbf{P}')^{-1}\mathbf{P}^{-1}\mathbf{X}]^{-1}\mathbf{X}'(\mathbf{P}')^{-1}\mathbf{P}^{-1}\mathbf{y} \\ &= [\mathbf{X}'(\mathbf{P}\mathbf{P}')^{-1}\mathbf{X}]^{-1}\mathbf{X}'(\mathbf{P}\mathbf{P}')^{-1}\mathbf{y} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}\end{aligned}$$



	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	47.25588	4.86946	9.705	2.40e-08 ***	
reach	0.33953	0.06316	5.376	5.03e-05 ***	

Residual standard error: 1.316 on 17 degrees of freedom

Multiple R-squared: 0.6296, Adjusted R-squared: 0.6078

$1^{\circ} P(\hat{\beta}_{\min} < \beta < \hat{\beta}_{\max}) > 1 - \alpha$   
 $2^{\circ} (\hat{\beta}_{\min}, \hat{\beta}_{\max})$ .

$\hat{\beta} \pm t_{n-2, \frac{\alpha}{2}} \text{SD}(\hat{\beta})$

$\downarrow$   
 is BLUE

$\beta$  smaller

$I_L$  98% t-interval

$I_I = \begin{cases} \emptyset & \text{with prob } 4\% \\ (-\infty, +\infty) & 96\% \end{cases} \Rightarrow$  有分數越大，不空集的期望越大

## Model error misspecification

- Suppose model is  $\mathbf{y} = \mathbf{X}\beta + \epsilon$  with  $\text{cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$  for some  $\mathbf{V} \neq \mathbf{I}$  and we assume that  $\text{cov}(\mathbf{y}) = \sigma^2 \mathbf{I}$
- Denoting the ordinary least squares (OLS) estimate by  $\hat{\beta}^* = \underline{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}}$ ,

$$E(\hat{\beta}^*) = \beta$$

$$\text{cov}(\hat{\beta}^*) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{V} \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}$$

$$\omega \mathbf{V}(\hat{\beta}^*) = A \text{cov}(\mathbf{y}) A'$$

- Thus OLS estimates remains unbiased but their variances tend to be larger

# Model structure misspecification

- Suppose the predictors can be partitioned into two sets

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = (\mathbf{X}_1, \mathbf{X}_2) \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix} + \boldsymbol{\epsilon} \\ &= \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon} \end{aligned}$$

- If we leave out  $\mathbf{X}_2\boldsymbol{\beta}_2$  when it should be included (i.e., when  $\boldsymbol{\beta}_2 \neq \mathbf{0}$ ), it is called underfitting
- If we include  $\mathbf{X}_2\boldsymbol{\beta}_2$  when it should be excluded (i.e., when  $\boldsymbol{\beta}_2 = \mathbf{0}$ ), it is called overfitting

**Theorem 13 (Underfitting).** If we fit the reduced model  $\mathbf{y} = \mathbf{X}_1\beta_1^* + \epsilon^*$ , where  $\mathbf{X}_1$  is  $n \times (p+1)$ , when the true model is  $\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \epsilon$  with  $\text{cov}(\mathbf{y}) = \sigma^2 \mathbf{I}$ , then the mean and covariance of OLS estimate  $\hat{\beta}_1^* = (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{y}$  are

$$\begin{aligned} E(\hat{\beta}_1^*) &= \beta_1 + \mathbf{A}\beta_2 \\ \text{cov}(\hat{\beta}_1^*) &= \sigma^2 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \end{aligned}$$

where  $\mathbf{A} = (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{X}_2$  and the variance estimate

$$s_1^2 = \frac{(\mathbf{y} - \mathbf{X}_1 \hat{\beta}_1^*)' (\mathbf{y} - \mathbf{X}_1 \hat{\beta}_1^*)}{n - p - 1}$$

has expected value

$$E(s_1^2) = \sigma^2 + \frac{\beta_2' \mathbf{X}'_2 \{ \mathbf{I} - \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \} \mathbf{X}_2 \beta_2}{n - p - 1} \geq \sigma^2$$

**Corollary 2** If  $\mathbf{X}'_1 \mathbf{X}_2 = 0$ , i.e., if the columns of  $\mathbf{X}_1$  are orthogonal to the columns of  $\mathbf{X}_2$ , then  $\hat{\beta}_1^*$  is unbiased

overfit

等同紙圖是否有 outlier

## Box-Cox transformation

Box-Cox method assumes that  $\{y_1^\lambda, y_2^\lambda, \dots, y_n^\lambda\}$  for some  $\lambda \neq 0$ , or  $\{\log y_1, \log y_2, \dots, \log y_n\}$  (corresponding to  $\lambda = 0$ )

1. are normally distributed
2. are mutually independent
3. have constant variance
4. satisfy a chosen model

$$y^\lambda \sim N(\beta, \sigma^2 I)$$

解决另外三纸图的问题

# Maximum likelihood estimation of Box-Cox

$$\lim_{\lambda \rightarrow 0} \frac{y^{\lambda}-1}{\lambda} = \lim_{\lambda \rightarrow 0} \frac{\log y - y^\lambda}{\lambda}$$

- Define

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

- Assume that there is  $\lambda$  such that  $\mathbf{y}^{(\lambda)} = \mathbf{X}\boldsymbol{\beta} + N(0, \sigma^2)$
- Likelihood function is

$$\frac{1}{(2\pi)^{n/2}\sigma^n} \exp \left\{ -\frac{(\mathbf{y}^{(\lambda)} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y}^{(\lambda)} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2} \right\} J(\lambda, \mathbf{y})$$

with Jacobian

fix  $\lambda$  got  $\hat{\boldsymbol{\beta}}$ ,  $\hat{\sigma}^2$

$$J(\lambda, \mathbf{y}) = \prod_{i=1}^n \frac{dy_i^{(\lambda)}}{dy_i} = \prod_{i=1}^n y_i^{\lambda-1}$$

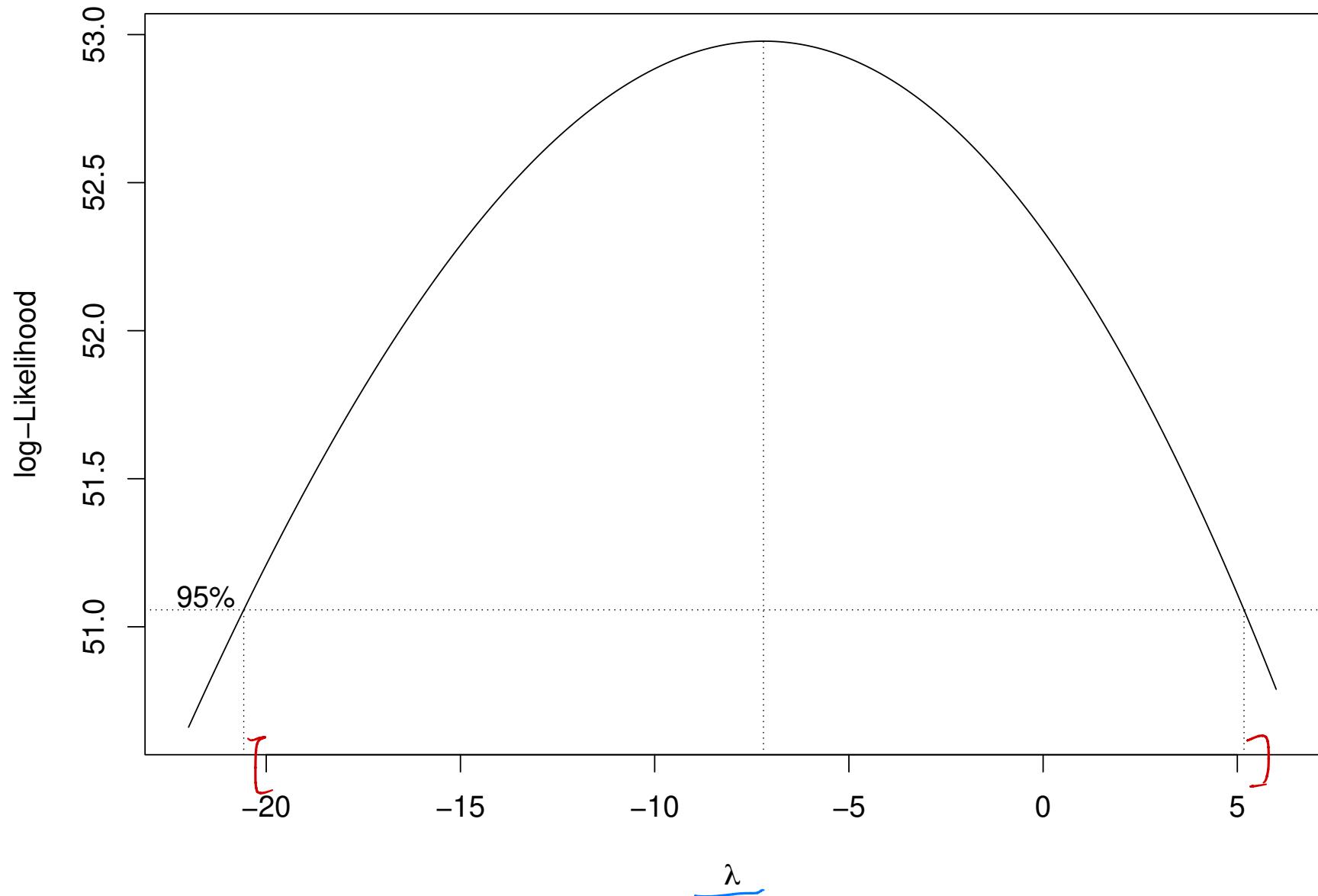
$\Downarrow$  Jacob 等於  $y^{(\lambda)}$  的方程,

$$(y_1^{\lambda-1}, \dots, y_n^{\lambda-1})$$

- Estimated  $\lambda$  is the MLE

- R function is called `boxcox` in the MASS package

# boxcox plot for boxer data using all predictors



# Solution without using boxcox

- Given  $\lambda$ , let  $(\hat{\beta}_\lambda, \hat{\sigma}_\lambda)$  maximize  $\ell(\lambda, \beta, \sigma)$
- Then  $(\hat{\beta}_\lambda, \hat{\sigma}_\lambda)$  are the LSEs for the  $\mathbf{y}^{(\lambda)}$  data, i.e.,

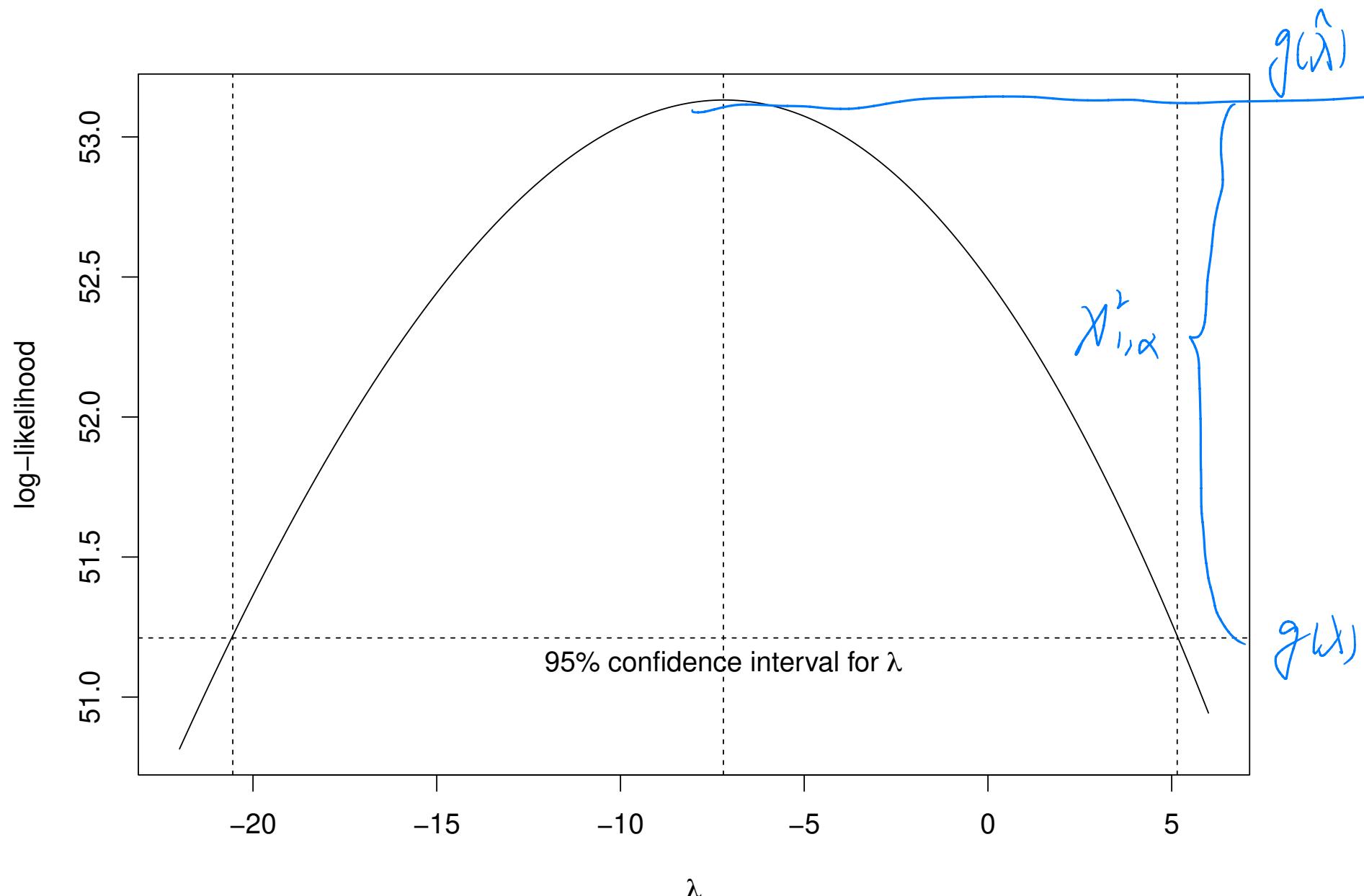
$$\hat{\beta}_\lambda = \underbrace{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}^{(\lambda)}}_{\text{blue underline}}, \quad \hat{\sigma}_\lambda^2 = \frac{(\mathbf{y}^{(\lambda)} - \mathbf{X}\hat{\beta}_\lambda)'(\mathbf{y}^{(\lambda)} - \mathbf{X}\hat{\beta}_\lambda)}{n - p - 1} \quad \text{blue underline} \quad n$$

and

$$\begin{aligned}\ell(\lambda, \hat{\beta}_\lambda, \hat{\sigma}_\lambda) &= \frac{\exp(-(n-p-1)/2)}{(2\pi)^{n/2}\hat{\sigma}_\lambda^n} J(\lambda, \mathbf{y}) \\ &= \hat{\sigma}_\lambda^{-n} \prod_{i=1}^n y_i^{\lambda-1} \frac{\exp(-(n-p-1)/2)}{(2\pi)^{n/2}}\end{aligned}$$

- Let  $g(\lambda) = \log \ell(\lambda, \hat{\beta}_\lambda, \hat{\sigma}_\lambda) = -n \log(\hat{\sigma}_\lambda) + (\lambda - 1) \sum_{i=1}^n \log(y_i) - \{(n - p - 1) + n \log(2\pi)\}/2$
- MLE  $\hat{\lambda}$  maximizes  $g(\lambda)$
- 100(1 -  $\alpha$ )% confidence limits for  $\lambda$  are roots of  $g(\lambda) = g(\hat{\lambda}) - \chi^2_{1;\alpha}/2$

# Manual plot for boxer data using all predictors



# Homework 3

## (due 11:59 PM Tue Oct 17)

1. Reproduce the confidence region plot on page 67
2. Show that the Box-Cox  $\lambda$  can be computed as follows:
  - (a) Let  $\dot{y} = \sqrt[n]{y_1 y_2 \dots y_n}$  denote the geometric mean of the  $y$  observations
  - (b) For each  $\lambda$ , transform the  $y_i$  to
$$z_i(\lambda) = \begin{cases} \frac{y_i^\lambda - 1}{\lambda \dot{y}^{\lambda-1}}, & \lambda \neq 0 \\ \dot{y} \log y_i, & \lambda = 0 \end{cases}$$
  - (c) Fit the desired linear regression model to  $z_1(\lambda), z_2(\lambda), \dots, z_n(\lambda)$  and compute its residual sum of squares  $S(\lambda)$
  - (d) The Box-Cox  $\lambda$  minimizes  $S(\lambda)$

$$Y \sim \text{binomial}(1, p), p = E(Y), P = X\beta \rightarrow \text{适合} \quad Y = \beta_0 + \beta_1 X + \varepsilon \quad P = \beta_0 + \beta_1 X$$

**Generalized linear models**

- Linear regression with normal errors assumes that  $\mathbf{Y}$  is normally distributed with constant variance and mean  $\mu = E(\mathbf{Y}) = \mathbf{X}\beta$
- Advantage:  $\beta$  gives effects of each predictor, controlling for the others
- A generalized linear model extends this idea to other distributions by introducing a link function
- Three useful link functions:

**Normal:**  $\eta(\mu) = \mu$  (identity function)

**Bernoulli with  $\mu = p$ :**  $\eta(p) = \log(p/(1-p))$  (logit function)

**Poisson:**  $\eta(\mu) = \log(\mu)$  (log function)

$$\log(\mu) = \beta_0 + \beta_1 X \quad (-\infty, +\infty)$$

$$\mu = e^{\beta_0 + \beta_1 X}$$

$\eta(\mu) = \mathbf{X}\beta \quad \leftarrow \text{keep}$

$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X + \varepsilon \quad t(-\infty, \infty)$

$p = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$

# Logistic regression example: Blowdown data

- A severe thunderstorm occurred in July 4, 1999, over 477,000 acres of the Boundary Waters Canoe Area Wilderness in northeastern Minnesota (data in `blowdown.csv` on Canvas; also in R package `alr4`)
- Observations from 3666 trees were collected on the following variables on each tree:
  1. whether it was blown down ( $Y = 1$ ) or not ( $Y = 0$ )
  2. its trunk diameter  $D$  in centimeters
  3. its species  $S$  (A-aspen, BA-black ash, BF-balsam fir, BS-black spruce, C-cedar, JP-jack pine, PB-paper birch, RM-red maple, RP-red pine)
  4. local intensity  $L$  of the storm (measured by fraction of damaged trees in its vicinity)

Normal  
t  
χ<sup>2</sup>  
F  
z

## Logistic regression theory

- $Y$  be Bernoulli with  $p = P(Y = 1)$ , let  $\mathbf{X} = (X_1, \dots, X_p)$  be the predictor variables, and define the logit link function  $\text{logit}(p) = \log(p/(1 - p))$
- The *multiple linear logistic regression* model assumes that

$$\text{logit}(p) = \beta_0 + \sum_{j=1}^p \beta_j X_j$$

or equivalently,  $p = \exp(\beta_0 + \sum_{j=1}^p \beta_j X_j) / \{1 + \exp(\beta_0 + \sum_{j=1}^p \beta_j X_j)\}$

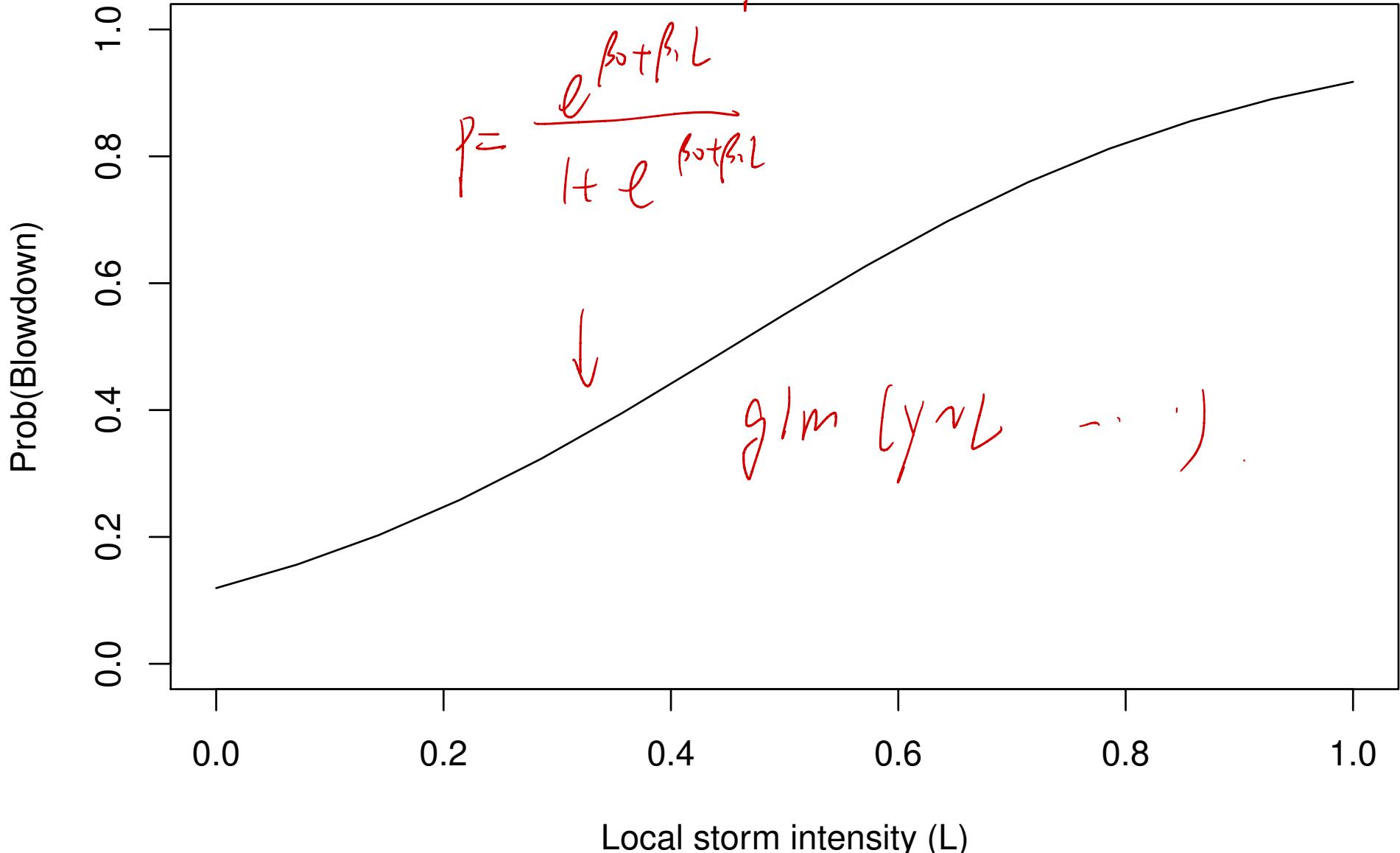
- Let  $n$  be the sample size, let  $(x_{i1}, \dots, x_{ip}, y_i)$  be the values of  $(X_1, \dots, X_p, Y)$  for the  $i$ th observation and let  $p_i = P(Y = y_i | \mathbf{X} = \mathbf{x}_i)$
- $\beta_0, \beta_1, \dots, \beta_k$  are estimated by maximizing the likelihood function

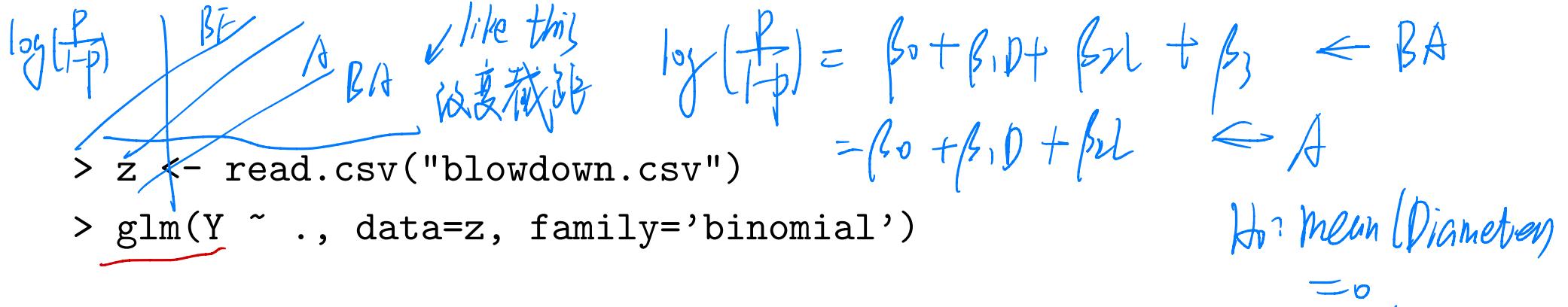
$$\prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} = \frac{\exp\{\sum_{i=1}^n y_i (\beta_0 + \sum_{j=1}^p \beta_j x_{ij})\}}{\prod_i \{1 + \exp(\beta_0 + \sum_j \beta_j x_{ij})\}}$$

# Logistic regression model using only $L$

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 L$$

$$P = \frac{e^{\beta_0 + \beta_1 L}}{1 + e^{\beta_0 + \beta_1 L}}$$

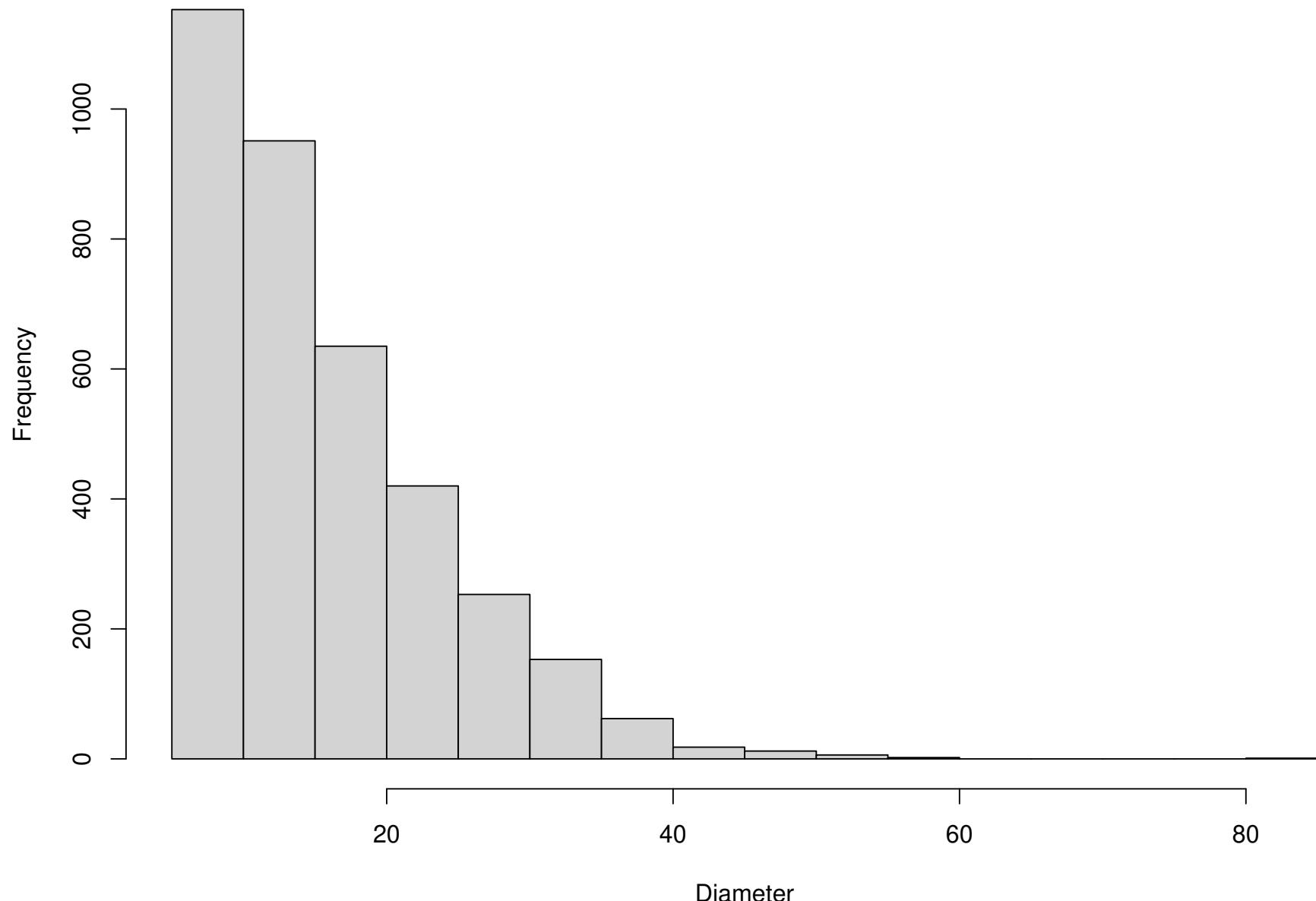




```
> z <- read.csv("blowdown.csv")
> glm(Y ~ ., data=z, family='binomial')
```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-2.887322	0.213111	-13.548	< 2e-16 ***
Diameter	0.075279	0.006872	10.955	< 2e-16 ***
Intensity	4.632263	0.209640	22.096	< 2e-16 ***
SpeciesBA	-2.337078	0.485614	-4.813	1.49e-06 ***
SpeciesBF	-0.360890	0.173417	-2.081	0.0374 *
SpeciesBS	0.057240	0.151798	0.377	0.7061
SpeciesC	-2.130241	0.215248	-9.897	< 2e-16 ***
SpeciesJP	1.076664	0.178150	6.044	1.51e-09 ***
SpeciesPB	-1.834158	0.185793	-9.872	< 2e-16 ***
SpeciesRM	-1.912548	0.300052	-6.374	1.84e-10 ***
SpeciesRP	-0.145442	0.424336	-0.343	0.7318

## Histogram of Diameter



# Using log(Diameter) in place of Diameter

```
> z$logDiameter <- log(z$Diameter)
```

```
> glm(Y ~ Diameter, data=z, family='binomial')
```

$$\log\left(\frac{P_1}{1-P_1}\right) = \beta_0 + \beta_1 X$$

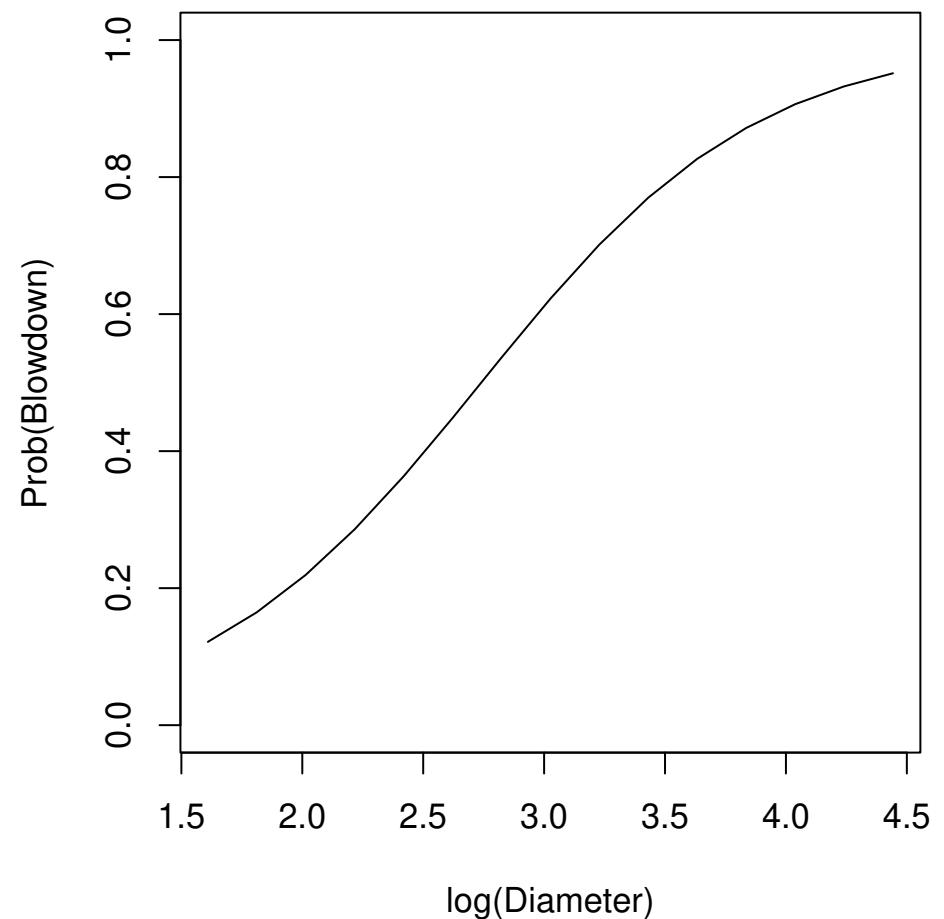
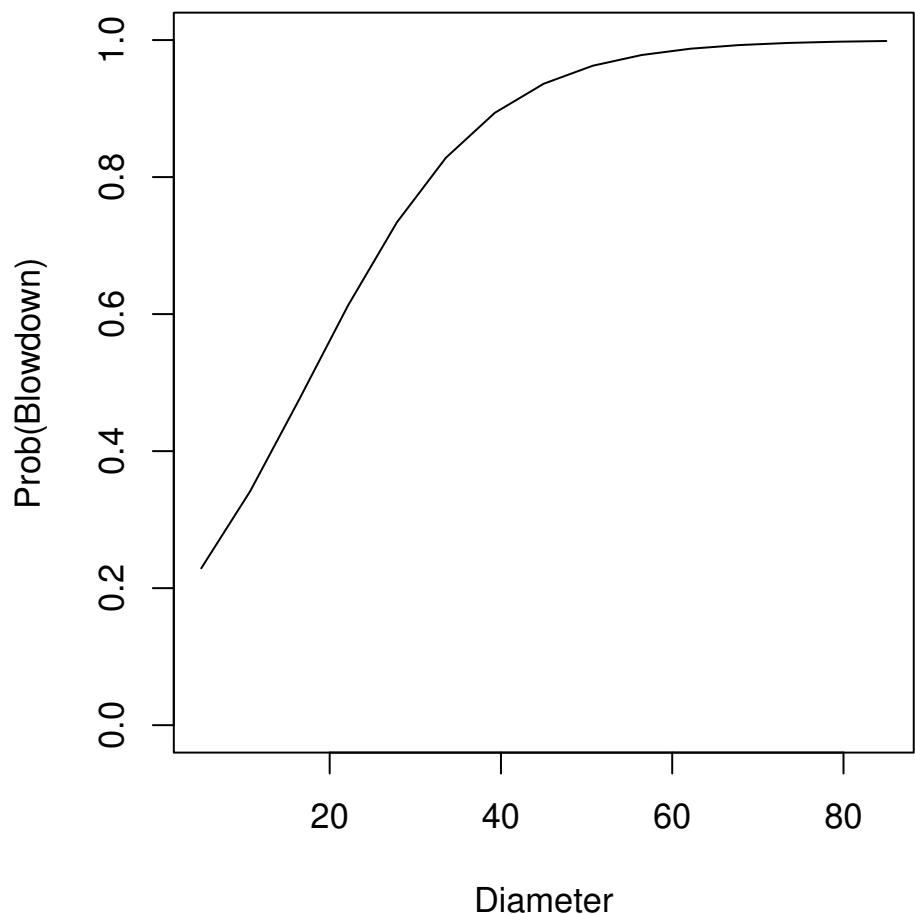
$$\log\left(\frac{P_2}{1-P_2}\right) = \beta_0 + \beta_1 (X+1)$$

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-5.9971951	0.3748409	-15.999	< 2e-16	***
Intensity	4.6288861	0.2128451	21.748	< 2e-16	***
SpeciesBA	-2.2427869	0.4935773	-4.544	5.52e-06	***
SpeciesBF	0.0002284	0.1789331	0.001	0.999	
SpeciesBS	0.1672262	0.1517509	1.102	0.270	
SpeciesC	-2.0765125	0.2162343	-9.603	< 2e-16	***
SpeciesJP	1.0399651	0.1787626	5.818	5.97e-09	***
SpeciesPB	-1.7235679	0.1864623	-9.244	< 2e-16	***
SpeciesRM	-1.7956738	0.3019339	-5.947	2.73e-09	***
SpeciesRP	0.0031381	0.4131721	0.008	0.994	
logDiameter	1.5813423	0.1114597	14.188	< 2e-16	***

$$\log\left(\frac{P_2}{1-P_2} / \frac{P_1}{1-P_1}\right) = \beta_0 + \beta_1$$

$\hat{f}_1$

# Logistic regression on Diameter and $\log(\text{Diameter})$



# **Homework 4**

## **(due 11:59 PM Tue Oct 24)**

1. Show that the likelihood function for the Box-Cox transformation on page 86 is invariant to scale multiplication (e.g., lbs  $\rightarrow$  kg).
2. Follow the steps on page 86 to reproduce the plot on page 87. Submit your R code.

# UMASS aids study

看 Canvas 的 data set

- Data from randomized trials on treatment for drug abuse at University of Massachusetts (Hosmer et al., 2013); data in `uis.csv` on Canvas
- Main goal to compare treatment programs of two different durations in reduction of drug abuse and in prevention of high-risk HIV behavior
- Variables are
  1. Subject ID
  2. Age at enrollment
  3. BECK depression score at admission
  4. IV drug use history at admission (1=never, 2=previous, 3=recent)
  5. Number of prior drug treatments
  6. Subject's race (0=white, 1=other)
  7. Treatment assignment (0=short, 1=long)
  8. Treatment site (0=A, 1=B)
  9. Patient's status at end of treatment program (1=remained drug free, 0=otherwise)

## R code

```

> z <- read.csv("uis.csv") 改成 factor 变量.
> z$SITE <- as.factor(z$SITE)
> z$RACE <- as.factor(z$RACE)
> z$TREAT <- as.factor(z$TREAT)  $\log\left(\frac{P}{1-P}\right) = \beta_0 + \dots + \beta_1 TREAT$ 
> z$IVHX <- as.factor(z$IVHX)
> glm(DFREE ~ . - ID, data=z, family='binomial')
↑ 原因

```

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-2.4111283	0.5983427	-4.030	5.59e-05	***
AGE	0.0504143	0.0174057	2.896	0.00377	**
BECK	0.0002759	0.0107982	0.026	0.97961	
IVHX2	-0.6036962	0.2875974	-2.099	0.03581	*
IVHX3	-0.7336591	0.2549893	-2.877	0.00401	**
NDRUGTX	-0.0615329	0.0256441	-2.399	0.01642	*
RACE1	0.2260262	0.2233685	1.012	0.31159	
TREAT1	+ 0.4424802	0.1992922	2.220	0.02640	*
SITE1	0.1489209	0.2176062	0.684	0.49375	

# Logistic regression with aggregated data: lung cancer example

- File `lungcancer.csv` gives number of lung cancer cases by age group and their population sizes in four Danish cities (data also from `eba1977` in R package `ISwR`)

city	age	pop	cases
Fredericia:6	40-54:4	Min. : 509.0	Min. : 2.000
Horsens :6	55-59:4	1st Qu.: 628.0	1st Qu.: 7.000
Kolding :6	60-64:4	Median : 791.0	Median :10.000
Vejle :6	65-69:4	Mean :1100.3	Mean : 9.333
	70-74:4	3rd Qu.: 954.8	3rd Qu.:11.000
	75+ :4	Max. :3142.0	Max. :15.000

- cases/pop is rate of lung cancer per capita by age group in each city

## lungcancer.csv

"city","age","pop","cases"  
"Fredericia","40-54",3059,11  
"Horsens","40-54",2879,13  
"Kolding","40-54",3142,4  
"Vejle","40-54",2520,5  
"Fredericia","55-59",800,11  
"Horsens","55-59",1083,6  
"Kolding","55-59",1050,8  
"Vejle","55-59",878,7  
"Fredericia","60-64",710,11  
"Horsens","60-64",923,15  
"Kolding","60-64",895,7  
"Vejle","60-64",839,10  
"Fredericia","65-69",581,10  
"Horsens","65-69",834,10  
"Kolding","65-69",702,11  
:  
:

# Logistic regression R code for lung cancer data

binomial 分布的次數 n

```
> z <- read.csv("lungcancer.csv")
> z$city <- as.factor(z$city); z$age <- as.factor(z$age)
> glm(cases/pop ~ city + age, weight=pop, family='binomial', data=z)
```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-5.6262	0.2008	-28.021	< 2e-16 ***
cityHorsens	-0.3345	0.1827	-1.830	0.0672 .
cityKolding	-0.3764	0.1890	-1.991	0.0465 *
cityVejle	-0.2760	0.1891	-1.459	0.1444
age55-59	1.1070	0.2490	4.445	8.77e-06 ***
age60-64	1.5291	0.2325	6.577	4.81e-11 ***
age65-69	1.7819	0.2305	7.732	1.06e-14 ***
age70-74	1.8727	0.2365	7.918	2.42e-15 ***
age75+	1.4289	0.2512	5.688	1.29e-08 ***

$$\log(\text{rate}/(1 - \text{rate})) = -5.6262 - 0.3345 \times I(\text{Horsens}) + \dots + 1.4289 \times I(\text{age75+})$$

$$\log(\mu) = \mathbf{X}\beta, \quad \log\left(\frac{P}{1-P}\right) = \mathbf{X}\gamma \quad \text{logistic}$$

Poisson  $\rightarrow \mu = \mathbf{X}\beta$  normal

## Poisson regression

$$\mu = E[y|x]$$

- Recall that a Poisson variable  $Y$  with mean  $\mu$  has density

$$P(Y = y) = \exp(-\mu)\mu^y/y!$$

$\rightarrow$  这里要求  $\mu > 0$  不要求  $0 < \mu < 1$

- Poisson regression assumes that  $\log \mu = \mathbf{X}\beta$ , or equivalently,  $\mu = \exp(\mathbf{X}\beta)$
- Likelihood function is

$$\ell(\beta) = \prod_{i=1}^n \frac{\exp\{-\exp(\mathbf{X}_i\beta)\} \exp(y_i \mathbf{X}_i\beta)}{y_i!}$$

- R function for fitting Poisson regression is `glm` with `family='poisson'`

## Crab data

The data file `crab.csv` contains observations on the following variables for 173 female crabs:

刀背蟹

**color.** color of crab (2=light medium, 3=medium, 4=dark medium, 5=dark)

**spine.** spine condition (1=both good, 2=one worn or broken, 3=both worn or broken)

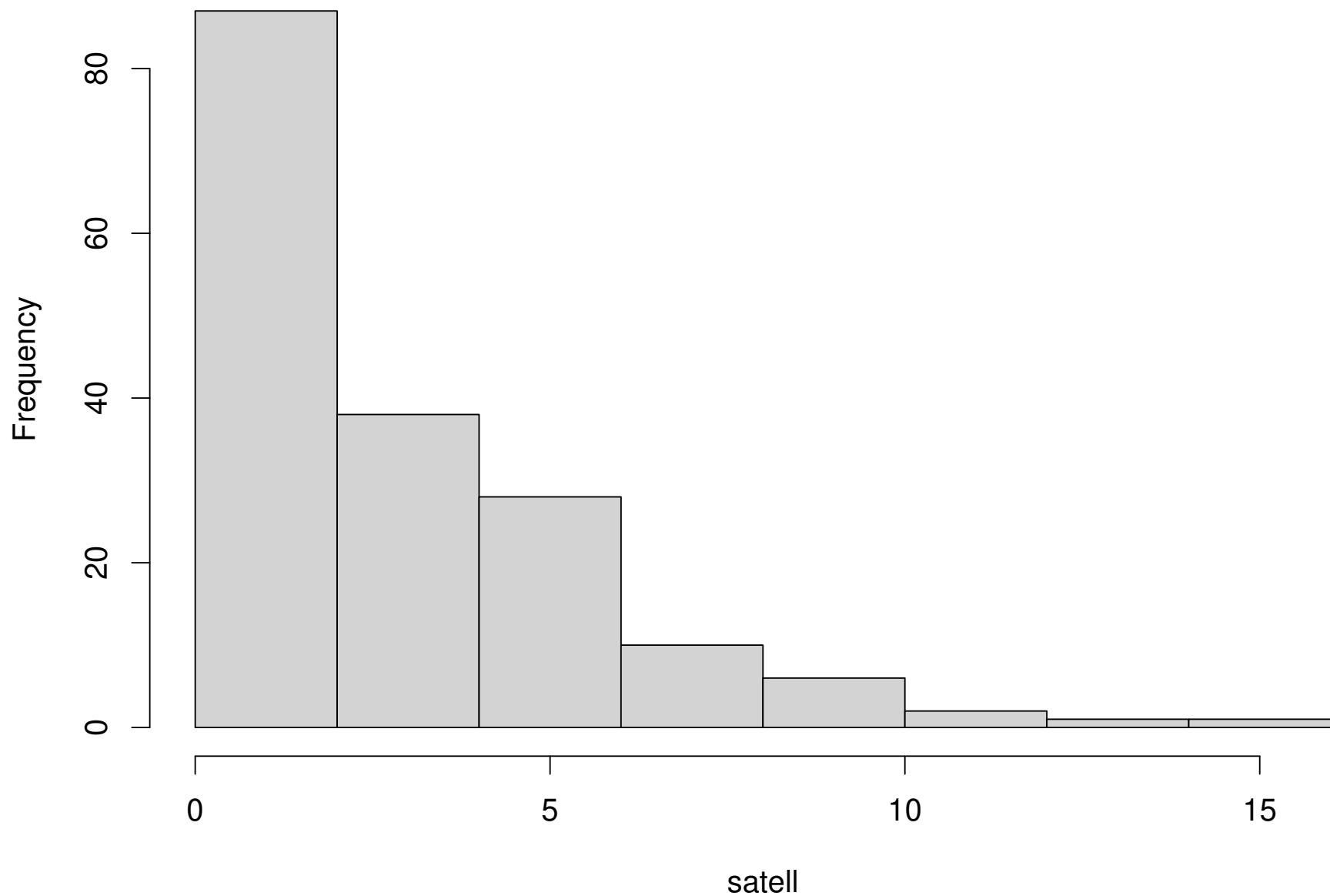
**width.** crab carapace width in cm

**satell.** number of satellites  $\text{Y} \leftarrow \text{Poisson}$

**weight.** crab weight in grams

**y.** whether crab has satellites (1=yes, 0=no)  $\leftarrow \text{logistic}$

## Histogram of satell



# Regression of satell on covariates ignoring y

```
> z <- read.csv("crab.csv")
> z$color <- as.factor(z$color)
> z$spine <- as.factor(z$spine)
> glm(satell ~ . - y, data=z, family='poisson')
```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.3618003	0.9665506	-0.374	0.70817
color3	-0.2648512	0.1681107	-1.575	0.11515
color4	-0.5137051	0.1953624	-2.629	0.00855 **
color5	-0.5308601	0.2269157	-2.339	0.01931 *
spine2	-0.1503718	0.2135754	-0.704	0.48139
spine3	0.0872826	0.1199287	0.728	0.46674
width	0.0167487	0.0489197	0.342	0.73207
weight	0.0004965	0.0001663	2.986	0.00283 **
---				

Null deviance: 632.79 on 172 degrees of freedom

Residual deviance: 549.59 on 165 degrees of freedom

highly correlated  
by y is → significant  
第一个 not weight to width  
还要看，单时  
P值小三个数量级  
106

$$\frac{M}{n} = \text{rate} > |$$

↑ worry

$Y \sim \text{Bin}(n, p)$ . If  $np \rightarrow c \stackrel{n \uparrow}{\downarrow} p \downarrow$  as  $n \rightarrow \infty$ , then Poisson  
 If  $p$  is fixed but  $n \rightarrow \infty$ , then Normal  
 这种情况下  $n$  不是泊松分布,  $Y$  也不不是正态分布

## Poisson regression with offset

- Let  $Y$  be the number of cases in a population of size  $n$  and let  $\mu = E(Y)$
- $Y$  is binomial but it is approximately Poisson if  $n$  is large
- The model  $\log(\mu/n) = \mathbf{X}\beta$  can be written as a Poisson regression model

$$\log(\mu) = \mathbf{X}\beta + \log(n) \quad \leftarrow \begin{cases} \beta_0, \dots, \beta_k \\ \beta_{k+1} = 1 \end{cases}$$

- $\log(n)$  is called an offset variable; it acts as a predictor variable with regression coefficient fixed at 1

## Insurance data (from MASS package)

**District.** district of residence of policyholder (1–4): 4 corresponds to major cities

**Group.** ordered factor: group of car with levels <1 litre, 1–1.5 litre, 1.5–2 litre, >2 litre

**Age.** ordered factor: age of insured in 4 groups labelled <25, 25–29, 30–35, >35

**Holders.** numbers of policy holders

**Claims.** numbers of claims

# Insurance data

District		Group	Age	Holders	Claims
1	1	<11	<25	197	38
2	1	<11	25-29	264	35
3	1	<11	30-35	246	20
4	1	<11	>35	1680	156
5	1	1-1.51	<25	284	63
6	1	1-1.51	25-29	536	84
7	1	1-1.51	30-35	696	89
8	1	1-1.51	>35	3582	400
9	1	1.5-21	<25	133	19
10	1	1.5-21	25-29	286	52
11	1	1.5-21	30-35	355	74
12	1	1.5-21	>35	1640	233
13	1	>21	<25	24	4
14	1	>21	25-29	71	18
15	1	>21	30-35	99	19
16	1	>21	>35	452	77
:					

$$\log\left(\frac{M}{n}\right) = X\beta$$

$$\frac{M}{n} = e^{X\beta} \sim \mathcal{N}$$

$$\log\left(\frac{\text{E(Claim)}}{\text{Holders}}\right)$$

$$= \log(\mu) + \log(\text{Holders})$$

# R code for insurance data

```
> library('MASS')  
> Insurance$Group <- factor(Insurance$Group, ordered=FALSE)  
> Insurance$Age <- factor(Insurance$Age, ordered=FALSE)  
> glm(Claims ~ District+Group+Age+offset(log(Holders)),  
      data=Insurance, family=poisson)
```

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-1.82174	0.07679	-23.724	< 2e-16	***
District2	0.02587	0.04302	0.601	0.547597	not significant compared with District 1
District3	0.03852	0.05051	0.763	0.445657	
District4	0.23421	0.06167	3.798	0.000146	***
Group1-1.51	0.16134	0.05053	3.193	0.001409	**
Group1.5-21	0.39281	0.05500	7.142	9.18e-13	***
Group>21	0.56341	0.07232	7.791	6.65e-15	***
Age25-29	-0.19101	0.08286	-2.305	0.021149	*
Age30-35	-0.34495	0.08137	-4.239	2.24e-05	***
Age>35	-0.53667	0.06996	-7.672	1.70e-14	***

# Fitting lung cancer data with Poisson regression (Poisson approximation to binomial)

```

> z <- read.csv("lungcancer.csv")
> z$city <- as.factor(z$city); z$age <- as.factor(z$age)
> glm(cases ~ city + age, offset=log(pop), family='poisson', data=z)

```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-5.6321	0.2003	-28.125	< 2e-16 ***
cityHorsens	-0.3301	0.1815	-1.818	0.0690 .
cityKolding	-0.3715	0.1878	-1.978	0.0479 *
cityVejle	-0.2723	0.1879	-1.450	0.1472
age55-59	1.1010	0.2483	4.434	9.23e-06 ***
age60-64	1.5186	0.2316	6.556	5.53e-11 ***
age65-69	1.7677	0.2294	7.704	1.31e-14 ***
age70-74	1.8569	0.2353	7.891	3.00e-15 ***
age75+	1.4197	0.2503	5.672	1.41e-08 ***

*✓ Rate in logistic regression*

$$\log(\text{cases}/\text{pop}) = -5.6321 - 0.3301 \times I(\text{Horsens}) + \dots + 1.4197 \times I(\text{age75+})$$

# Survival and hazard functions

- Let  $\underline{U}$  be survival time of a subject with density  $f(u)$  and cdf

$$F(u) = P(\underline{U} \leq u) = \int_{-\infty}^u f(x) dx$$

- Survival probability function** is  $S(u) = P(\underline{U} > u) = 1 - F(u)$
- Hazard rate** (instantaneous rate of death) at time  $u$  is

$$\lambda(u) = f(u)/S(u) \tag{4}$$

- Cumulative hazard** function at time  $u$  is

$$\Lambda(u) = \int_{-\infty}^u \lambda(x) dx = \int_{-\infty}^u f(x)/S(x) dx = -\log(S(u))$$

and

$$S(u) = \exp(-\Lambda(u)) \tag{5}$$

$$u \longleftrightarrow X$$

## Cox proportional hazards (PH) regression

- Assume that given covariate vector  $\mathbf{x}$ ,  $\lambda(u, \mathbf{x}) = \lambda_0(u) \exp(\mathbf{x}'\boldsymbol{\beta})$ , where  $\lambda_0(u)$  is an unknown **baseline hazard** function  
*Note:  $\lambda_0$ 里面*
- log hazard ratio** (also called log relative risk) does not depend on  $u$ :  
$$\log\{\lambda(u, \mathbf{x})/\lambda_0(u)\} = \mathbf{x}'\boldsymbol{\beta}$$
- Hazards of two subjects with covariates  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are **proportional**:  
$$\lambda(u, \mathbf{x}_2)/\lambda(u, \mathbf{x}_1) = \exp\{(\mathbf{x}_2 - \mathbf{x}_1)' \underbrace{\boldsymbol{\beta}}_{\text{只与 } \mathbf{x} \text{ 有关}}\}$$
- Let  $\mathbf{x}_1 = (a_1, a_2, \dots, a_k)$  and  $\mathbf{x}_2 = (b_1, b_2, \dots, b_k)$ , such that  $b_j = a_j + 1$  and  $b_i = a_i$  for  $i \neq j$ . Then  
$$\beta_j = \log \left\{ \frac{\lambda(u, \mathbf{x}_2)}{\lambda(u, \mathbf{x}_1)} \right\}$$

i.e.,  $\beta_j = \log$  relative risk due to increasing  $j$ th element of  $\mathbf{x}$  by 1

- Let  $S(u, \mathbf{x})$ ,  $\lambda(u, \mathbf{x})$  denote survival probability and hazard functions at  $\mathbf{x}$
- PH model assumes  $\lambda(u, \mathbf{x}) = \lambda_0(u) \exp(\mathbf{x}'\boldsymbol{\beta})$
- Let  $\Lambda_0(u) = \int_{-\infty}^u \lambda_0(z) dz$  be the **baseline cumulative hazard** and let  $\Lambda(u, \mathbf{x}) = \int_{-\infty}^u \lambda(z, \mathbf{x}) dz = \Lambda_0(u) \exp(\mathbf{x}'\boldsymbol{\beta})$  denote the **cumulative hazard**
- Then by (4) and (5),

$$\begin{aligned}
 S(u, \mathbf{x}) &= \exp\{-\Lambda(u, \mathbf{x})\} \\
 &= \exp\{-\Lambda_0(u) \exp(\mathbf{x}'\boldsymbol{\beta})\} \\
 f(u, \mathbf{x}) &= \lambda(u, \mathbf{x})S(u, \mathbf{x}) \\
 &= \lambda_0(u) \exp\{\mathbf{x}'\boldsymbol{\beta} - \Lambda_0(u) \exp(\mathbf{x}'\boldsymbol{\beta})\}
 \end{aligned}$$

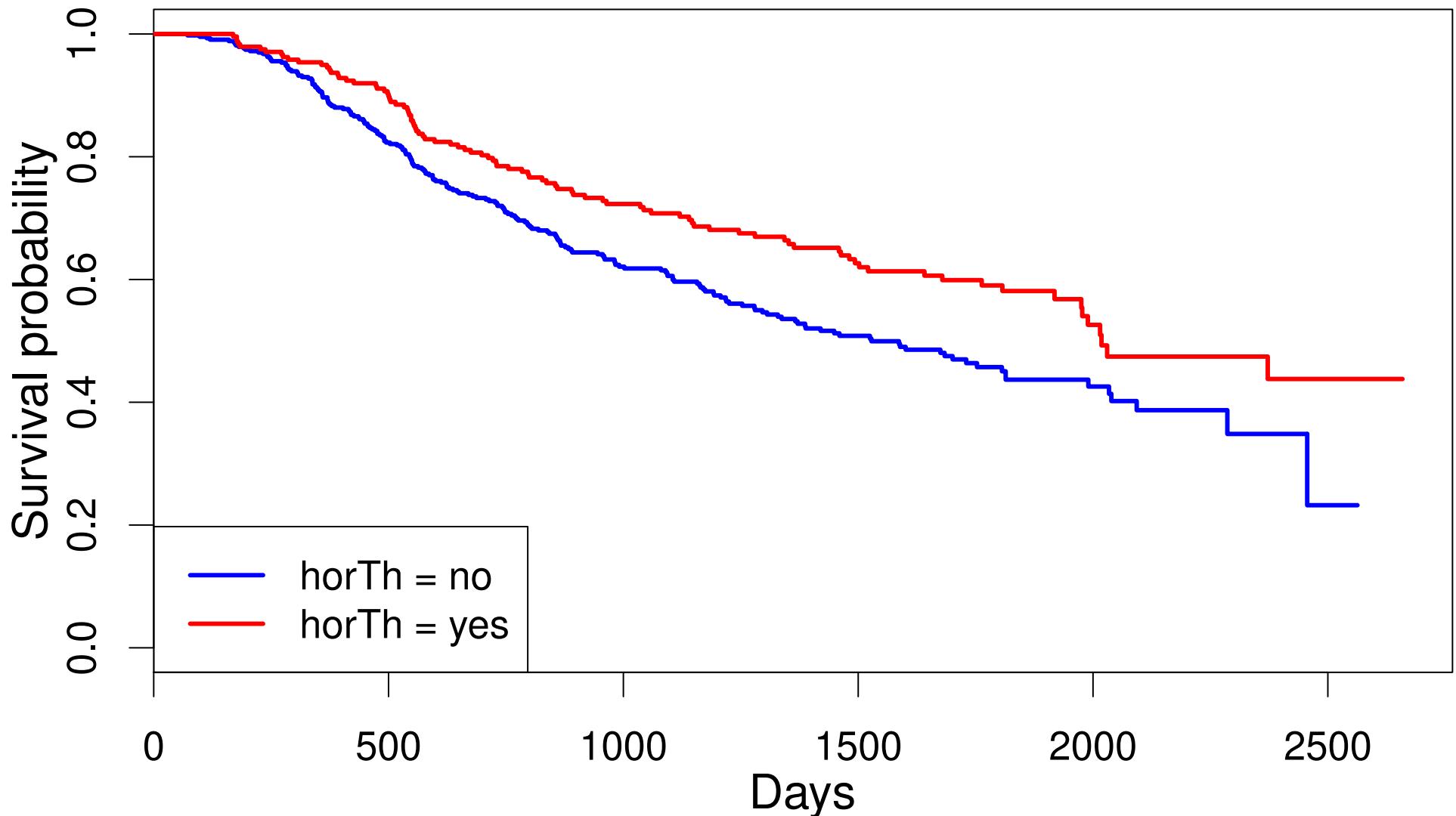
- Let  $u_i$  and  $\mathbf{x}_i$  denote the survival time and covariate vector of subject  $i$
- Let  $s_i$  be an independent observation from some censoring distribution and let  $\delta_i = I(u_i < s_i)$  be the event (e.g., death) indicator
- Observed data for subject  $i$  is  $(y_i, \delta_i, \mathbf{x}_i)$ , where  $y_i = \min(u_i, s_i)$
- $\hat{\beta}$  maximizes the log-likelihood

$$\begin{aligned}
& \log \left\{ \prod_{i=1}^n f(y_i, \mathbf{x}_i)^{\delta_i} S(y_i, \mathbf{x}_i)^{1-\delta_i} \right\} \\
&= \sum_{i=1}^n \delta_i \log \{f(y_i, \mathbf{x}_i)/S(y_i, \mathbf{x}_i)\} + \sum_{i=1}^n \log S(y_i, \mathbf{x}_i) \\
&= \sum_{i=1}^n \delta_i \{\log \lambda_0(y_i) + \mathbf{x}'_i \beta\} - \sum_{i=1}^n \Lambda_0(y_i) \exp(\mathbf{x}'_i \beta) \\
&= \sum_{i=1}^n \delta_i \{\log \Lambda_0(y_i) + \mathbf{x}'_i \beta\} - \sum_{i=1}^n \Lambda_0(y_i) \exp(\mathbf{x}'_i \beta) \\
&\quad + \sum_{i=1}^n \delta_i \log \{\lambda_0(y_i)/\Lambda_0(y_i)\}
\end{aligned}$$

# Breast cancer example

- Randomized clinical trial of 672 subjects with primary node positive breast cancer (Schumacher et al., 1994); 14 subjects with censored times less than smallest uncensored time excluded; data from TH.data R package
- Response is recurrence-free survival time (8–2659 days, 299 uncensored, 387 censored)  
*still cancer free ↑ in between no cancer.*
- Eight predictor variables: (trial停止的時候)
  1. **horTh** (hormone therapy, yes/no)
  2. **age** (21–80 years)
  3. **tsize** (tumor size, 3–120 mm)
  4. **pnodes** (number of positive lymph nodes, 1–51)
  5. **progrec** (progesterone receptor status, 0–2380 fmol)
  6. **estrec** (estrogen receptor status, 0–1144 fmol)
  7. **menostat** (menopausal status, pre/post)
  8. **tgrade** (tumor grade, 1, 2, 3)

# Kaplan-Meier survival curves



# R code for Kaplan-Meier plot

```
library(survival)
library(MASS)
z <- read.csv("cancerdata.csv")
leg.txt <- c("horTh = no","horTh = yes")
leg.col <- c("blue","red")
leg.lty <- rep(1,2)
y <- z$time
stat <- z$death
treat <- z$horTh
fit.bytreat <- survfit(Surv(y,stat) ~ treat, conf.type="none")
plot(fit.bytreat,conf.int=FALSE,col=leg.col,mark.time=FALSE,lwd=3)
mtext("Survival probability",side=2,line=2,cex=1.7)
mtext("Days",side=1,line=2,cex=1.7)
legend("bottomleft",legend=leg.txt,lty=leg.lty,col=leg.col,lwd=3)
```

$0 \rightarrow 1$ ,  $X \rightarrow X + 1$

## R code for Cox PH model

```
> summary(coxph(Surv(time, death) ~ ., data=z))
```

horThyes = 1.

less likely die

	coef	exp(coef)	se(coef)	z	Pr(> z )	
horThyes	-0.3372029	0.7137640	0.1289618	-2.615	0.008929	**
age	-0.0093924	0.9906516	0.0092733	-1.013	0.311136	
menostatPre	-0.2672772	0.7654609	0.1833366	-1.458	0.144882	
tsize	0.0077164	1.0077463	0.0039497	1.954	0.050739	.
tgrade	0.2802894	1.3235128	0.1060553	2.643	0.008221	**
pnodes	0.0498939	1.0511596	0.0074094	6.734	1.65e-11	***
progres	-0.0022378	0.9977647	0.0005758	-3.887	0.000102	***
estrec	0.0001674	1.0001674	0.0004477	0.374	0.708431	

similar pattern  
correlated

still significant



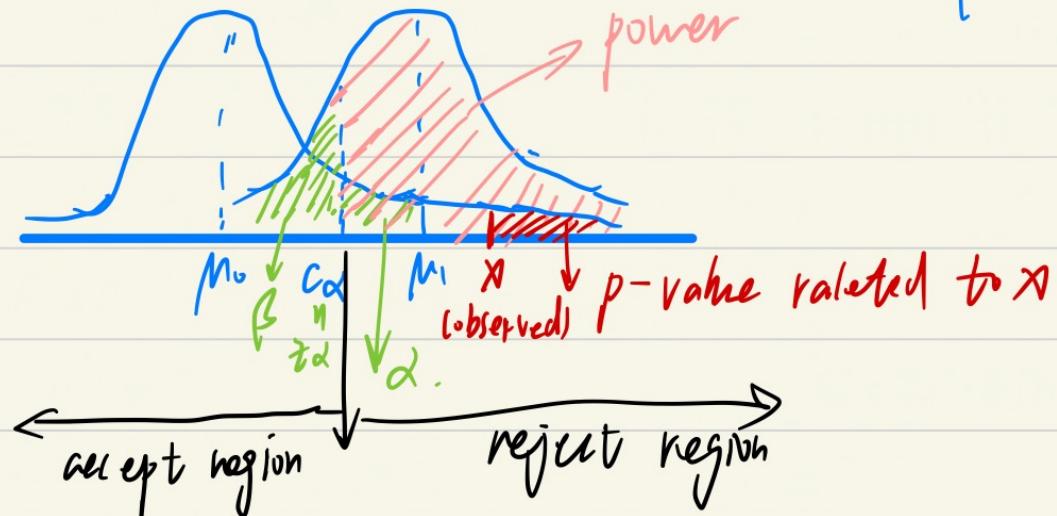
\*\*

partial effect

Consider  $X \sim (\mu, 1^2)$ ,  $H_0: \mu = 0$ ,  $H_1: \mu = 1$ ,  $\alpha = P[\text{Type I error}]$

$H_0$        $H_1$

$\beta = P[\text{Type II error}]$



若从  $6^2$  变的情况