



JOHNS HOPKINS
BLOOMBERG SCHOOL
of PUBLIC HEALTH

Lecture 8

Advanced inference in multiple linear regression

The material in this video is subject to the copyright of the owners of the material and is being provided for educational purposes under rules of fair use for registered students in this course only. No additional copies of the copyrighted work may be made or distributed.

Review of where we left off

1. We have established the multiple linear regression model:

$$\underbrace{Y_{n \times 1}} = \underbrace{X_{n \times (p+1)}}_{\substack{\mu \\ n \times 1}} \underbrace{\beta_{(p+1) \times 1}} + \underbrace{\epsilon_{n \times 1}}_{\substack{\mu \\ n \times 1}} \sim MVN(0_{n \times 1}, \sigma^2 I_{n \times n})$$

2. We know that:

$$\hat{\beta} \text{ satisfies } \underbrace{X'(Y - X\beta) = 0}_{\text{sums of squared residuals}} \text{ and minimizes } \sum_{i=1}^n (y_i - x_i' \beta)^2$$

$$V = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

3. We have defined:

- $\hat{Y} = X\hat{\beta} = HY$, where $H = X(X'X)^{-1}X'$
- $\hat{R} = Y - \hat{Y} = Y - X\hat{\beta} = (I - H)Y$

4. Then we showed that:

- $\hat{\beta} \sim MVN(\beta, \sigma^2(X'X)^{-1})$
- $\hat{Y} \sim MVN(X\beta, \sigma^2 H)$
- $\hat{R} \sim MVN(0, \sigma^2(I - H))$

Possible inference: single regression coefficient

Target	Estimate ~ Sampling Distn	95% CI for target	Test statistic for H0: Target = 0
β_j	$\hat{\beta}_j \sim N(\beta_j, [\sigma^2(X'X)^{-1}]_{jj})$	$\hat{\beta}_j \pm t \times \hat{se}(\hat{\beta}_j)$	$\frac{\hat{\beta}_j}{\hat{se}(\hat{\beta}_j)}$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} \sim \text{MVN} \left(\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \begin{bmatrix} \text{Var}(\beta_0) & \text{Cov}(\beta_0, \beta_1) & & \\ \text{Cov}(\beta_0, \beta_1) & \text{Var}(\beta_1) & & \\ & & \ddots & \\ & & & \text{Var}(\beta_p) \end{bmatrix} \right)$$

$\sigma^2 (X'X)^{-1}$

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



Example: inference for single regression coefficient

Nepali data $Y_i = \text{arm circumference}$

$$\text{age} = \begin{cases} \text{age} \\ (\text{age} - 6)^+ \end{cases}$$

$$Y_i = \beta_0 + \beta_1 \text{age}_i + \beta_2 (\text{age}_i - 6)^+ + \varepsilon_i$$

linear
spline
model with
knot at
6 months

β_1 = linear growth rate per month of
age among children under 6 months

$$H_0: \beta_1 = 0$$

95% CI for β_1

$$H_A: \beta_1 \neq 0$$

$$\text{test statistic} = .31 / .09 = 3.361$$

```
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept) 11.12089    0.50959  21.823 < 2e-16 ***
## age         0.31141     0.09264   3.361 0.000945 ***
## agesp6     -0.27958    0.09441  -2.961 0.003473 **
## ---
```

$$.31 \pm t \times .09$$

$$df = n - p - 1$$

$$n - p - 2$$

Possible inference: linear combination of coefficients

Target	Estimate ~ Sampling Distn	95% CI for target	Test statistic for H0: Target = 0
β_j	$\hat{\beta}_j \sim N(\beta_j, [\sigma^2(X^T X)^{-1}]_{jj})$	$\hat{\beta}_j \pm t \times \hat{se}(\hat{\beta}_j)$	$\frac{\hat{\beta}_j}{\hat{se}(\hat{\beta}_j)}$
<u>$A\beta$</u>	<u>$A\hat{\beta} \sim N(A\beta, \sigma^2 A(X^T X)^{-1} A^T)$</u>	<u>$A\hat{\beta} \pm t \times \hat{se}(A\hat{\beta})$</u>	$\left[\frac{A\hat{\beta}_j}{\hat{se}(A\hat{\beta}_j)} \right]$

$A = 1 \times (p+1)$ vector = (c_0, c_1, \dots, c_p)

$$A_{1 \times (p+1)} \beta_{(p+1) \times 1} = c_0 \beta_0 + c_1 \beta_1 + \dots + c_p \beta_p$$

Example: linear growth rate per month of age among children greater than 6 months

target: $\beta_1 + \beta_2$ estimate: $\hat{\beta}_1 + \hat{\beta}_2$

$$\text{Var}(\hat{\beta}_1 + \hat{\beta}_2) = \text{Var}(\hat{\beta}_1) + \text{Var}(\hat{\beta}_2) + 2 \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) *$$

$$\text{Var}(aY + bX) = a^2 \text{Var}(Y) + b^2 \text{Var}(X) + 2ab \text{Cov}(X, Y)$$

Example: inference for linear combination of coefficients

Target	Estimate ~ Sampling Distn	95% CI for target	Test statistic for H0: Target = 0
β_j	$\hat{\beta}_j \sim N(\beta_j, [\sigma^2(X'X)^{-1}]_{jj})$	$\hat{\beta}_j \pm t \times \hat{se}(\hat{\beta}_j)$	$\frac{\hat{\beta}_j}{\hat{se}(\hat{\beta}_j)}$
$A\beta$	$A\hat{\beta} \sim N(A\beta, \sigma^2 A(X'X)^{-1}A')$	$A\hat{\beta} \pm t \times \hat{se}(A\hat{\beta})$	$\frac{A\hat{\beta}_j}{\hat{se}(A\hat{\beta}_j)}$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

$$A = \begin{pmatrix} 0 & 1 & 1 \end{pmatrix}_{1 \times 3}$$

$$A\hat{\beta} = 0\hat{\beta}_0 + 1\hat{\beta}_1 + 1\hat{\beta}_2$$

$$C_{ij} = \text{Cov}(\hat{\beta}_i, \hat{\beta}_j)$$

$$\sigma^2 (X'X)^{-1} = \begin{bmatrix} \text{Var}(\hat{\beta}_0) & C_{01} & C_{02} \\ C_{01} & \text{Var}(\hat{\beta}_1) & C_{12} \\ C_{02} & C_{12} & \text{Var}(\hat{\beta}_2) \end{bmatrix}_{3 \times 3}$$

$$\text{Var}(A\hat{\beta}) = A_{1 \times 3} \text{Var}(\hat{\beta})_{3 \times 3} A'_{3 \times 1} =$$

$$\begin{bmatrix} 0 \cdot \text{Var}(\hat{\beta}_0) + C_{01} + C_{02} & 0 \cdot C_{01} + \text{Var}(\hat{\beta}_1) + C_{12} & 0 \cdot C_{02} + C_{12} + \text{Var}(\hat{\beta}_2) \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

$$\text{Var}(\hat{\beta}_1) + C_{12} + C_{12} + \text{Var}(\hat{\beta}_2) = \text{Var}(\hat{\beta}_1) + \text{Var}(\hat{\beta}_2) + 2C_{12}$$

Example: inference for linear combination of coefficients

```
cc=complete.cases(select(d,age,arm))  
d.cc=filter(d,cc)  
d.cc = arrange(d.cc,age)  
reg1<-lm(data=d.cc, arm~age+agesp6)  $\Rightarrow$   
reg1.coef = reg1$coef  
reg1.vc = vcov(reg1)
```

Define the linear combination of betas

```
A = matrix(c(0,1,1),nrow=1,ncol=3)  $\rightarrow$ 
```

Estimate the A beta-hat

```
A %%% reg1.coef
```

```
##           [,1]
```

```
## [1,] 0.03182924
```

What is the statistical variance of the estimate

```
A %%% reg1.vc %%% t(A)
```

```
##           [,1]
```

```
## [1,] 1.985802e-05
```

What is the standard error of the estimate

```
sqrt(A %%% reg1.vc %%% t(A))
```

```
##           [,1]
```

```
## [1,] 0.004456234
```

$$\text{Var}(\hat{\beta}_1 + \hat{\beta}_2)$$

$$\text{se}(\hat{\beta}_1 + \hat{\beta}_2)$$



Example: inference for linear combination of coefficients

```
# Confirm these values!  
summary(glht(reg1, linfct = A))
```

```
##  
## Simultaneous Tests for General Linear Hypotheses  
##  
## Fit: lm(formula = arm ~ age + agesp6, data = d.cc)  
##  
## Linear Hypotheses:  
## Estimate Std. Error t value Pr(>|t|)  
## 1 == 0 0.031829 0.004456 7.143 2.12e-11 ***  
## ---  
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
## (Adjusted p values reported -- single-step method)
```

$$\frac{0.0318}{0.0044} = 7.143 \Rightarrow p\text{value}$$

```
# 95% CI for beta1 + beta2  
A %>% reg1.coef - qt(0.975,df=summary(reg1)$df[2]) * sqrt(A %>% reg1.vc %>% t(A))
```

```
## [1,]  
## [1,] 0.02303672
```

```
A %>% reg1.coef + qt(0.975,df=summary(reg1)$df[2]) * sqrt(A %>% reg1.vc %>% t(A))
```

```
## [1,]  
## [1,] 0.04062177
```

95% CI for $\beta_1 + \beta_2$

.023 to .041

Example: inference for linear combination of coefficients

```
# Confirm these values!  
summary(glht(reg1, linfct = A))  
  
##  
## Simultaneous Tests for General Linear Hypotheses  
##  
## Fit: lm(formula = arm ~ age + agesp6, data = d.cc)  
##  
## Linear Hypotheses:  
## Estimate Std. Error t value Pr(>|t|)  
## 1 == 0 0.031829 0.004456 7.143 2.12e-11 ***  
## ---  
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
## (Adjusted p values reported -- single-step method)  
  
# Hypothesis test of H0: beta1 + beta2 = 0  
test.stat = (A %*% reg1.coef) / sqrt(A %*% reg1.vc %*% t(A))  
test.stat  
  
## [1]  
## [1,] 7.142632  
  
2 * pt(abs(test.stat),df=summary(reg1)$df[2],lower.tail=FALSE)  
  
## [1]  
## [1,] 2.124636e-11
```

Possible inference: Non-linear function of a coefficient

Target	Estimate ~ Sampling Distn	95% CI for target	Test statistic for H0: Target = 0
β_j	$\hat{\beta}_j \sim N(\beta_j, [\sigma^2(X'X)^{-1}]_{jj})$	$\hat{\beta}_j \pm t \times \hat{se}(\hat{\beta}_j)$	$\frac{\hat{\beta}_j}{\hat{se}(\hat{\beta}_j)}$
$A\beta$	$A\hat{\beta} \sim N(A\beta, \sigma^2 A(X'X)^{-1}A')$	$A\hat{\beta} \pm t \times \hat{se}(A\hat{\beta})$	$\frac{A\hat{\beta}_j}{\hat{se}(A\hat{\beta}_j)}$
$g(\beta_j)$	$g(\hat{\beta}_j) \sim N(g(\beta_j), \underbrace{[g'(\beta_j)]^2}_{\text{continuous at its 1st derivative}} [\sigma^2(X'X)^{-1}]_{jj})$	$g(\hat{\beta}_j) \pm t \times \hat{se}(g(\hat{\beta}_j))$	$\frac{g(\hat{\beta}_j)}{\hat{se}(g(\hat{\beta}_j))}$

g is a function that is continuous at its 1st derivative

Example: We say that Y follows a log-normal distribution
if $\log(Y) \sim N(\mu, \sigma^2) \Rightarrow E(Y) = \exp(\mu + \sigma^2/2)$
median(Y) = $\exp(\mu)$

PS2 NMES = $Y_i = \text{expenditures}_i + 1$

$\log(Y_i) = \beta_0 + \beta_1(\text{age} - 65) + \beta_2(\text{age} - 75)^+ + \beta_3(\text{age} - 85)^+ + \varepsilon_i$

Estimate the median expenditure for 65 year olds

Example: inference for non-linear function of a coefficient

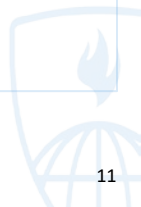
median expenditure for 65 year olds

$$E[\log(Y_i)] \text{ at } 65 = \beta_0 = \exp(\beta_0) - 1$$

$$g(\beta_0) = \exp(\beta_0) - 1 \quad \text{estimate: plug in } \underbrace{\exp(\hat{\beta}_0) - 1}$$

$$\frac{\partial}{\partial \beta_0} g(\beta_0) = \frac{\partial}{\partial \beta_0} \exp(\beta_0) - 1 = \underbrace{\exp(\beta_0)}$$

$$\text{Var}(g(\hat{\beta}_0)) = \exp(\hat{\beta}_0)^2 \left[\hat{\sigma}^2 (X'X)^{-1} \right]_{11}$$



Univariate delta method

Assuming the function g is continuous at its first derivative. The delta method is derived from the first order approximation to Taylor series using Taylor's theorem.

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

In statistical applications, we are interested in finding the distribution of $g(\hat{\theta})$ where $\hat{\theta}$ follows a normal distribution.

Applying the first order Taylor expansion to $g(\hat{\theta})$ about the mean θ , we get:

$$g(\hat{\theta}) \approx g(\theta) + g'(\theta)(\hat{\theta} - \theta)$$

Then, $E(g(\hat{\theta})) = g(\theta) + g'(\theta)(E(\hat{\theta}) - \theta) = g(\theta) + g'(\theta - \theta) = g(\theta)$ and $Var(g(\hat{\theta})) = g'(\theta)^2 Var(\hat{\theta})$.

$$E[g(\theta) + g'(\theta)(\hat{\theta} - \theta)] \quad \downarrow \quad g'(\theta)(\theta - \theta)$$

0



Possible inference: non-linear function of coefficients

Target	Estimate ~ Sampling Distn	95% CI for target	Test statistic for H0: Target = 0
β_j	$\hat{\beta}_j \sim N(\beta_j, [\sigma^2(X'X)^{-1}]_{jj})$	$\hat{\beta}_j \pm t \times \hat{se}(\hat{\beta}_j)$	$\frac{\hat{\beta}_j}{\hat{se}(\hat{\beta}_j)}$
$A\beta$	$A\hat{\beta} \sim N(A\beta, \sigma^2 A(X'X)^{-1}A')$	$A\hat{\beta} \pm t \times \hat{se}(A\hat{\beta})$	$\frac{A\hat{\beta}_j}{\hat{se}(A\hat{\beta}_j)}$
univariate delta $g(\beta_j)$	$g(\hat{\beta}_j) \sim N(g(\beta_j), [g'(\beta_j)]^2 [\sigma^2(X'X)^{-1}]_{jj})$	$g(\hat{\beta}_j) \pm t \times \hat{se}(g(\hat{\beta}_j))$	$\frac{g(\hat{\beta}_j)}{\hat{se}(g(\hat{\beta}_j))}$
multivariate delta $g(\beta)$	$g(\hat{\beta}) \sim N(g(\beta), g'(\beta)'[\sigma^2(X'X)^{-1}]g'(\beta))$	$g(\hat{\beta}) \pm t \times \hat{se}(g(\hat{\beta}))$	$\frac{g(\hat{\beta})}{\hat{se}(g(\hat{\beta}))}$

Example: Y_i = arm circumference

$$Y_i = \beta_0 + \beta_1 \text{age}_i + \beta_2 (\text{age}_i - 6)^+ + \varepsilon_i$$

monthly growth
rate

< 6 months β_1
 ≥ 6 months $\beta_1 + \beta_2$

$$\text{relative growth rate} = \frac{\beta_1 + \beta_2}{\beta_1}$$

$$= 1 + \frac{\beta_2}{\beta_1}$$

Example: inference for non-linear function of coefficients

$$g(\beta) = 1 + \beta_2/\beta_1$$

estimate \Rightarrow plug in estimate
 $1 + \hat{\beta}_2/\hat{\beta}_1$

$$g'(\beta) = \begin{bmatrix} \frac{d}{d\beta_0} g \\ \frac{d}{d\beta_1} g \\ \frac{d}{d\beta_2} g \end{bmatrix} = \begin{bmatrix} 0 \\ -\beta_2/\beta_1^2 \\ 1/\beta_1 \end{bmatrix}_{3 \times 1}$$

$\text{Var}(g(\hat{\beta})) ??$

$$\text{Var}(g(\hat{\beta})) = \underbrace{[g'(\tilde{\beta})]'}_{1 \times 3} \underbrace{V_{\hat{\beta}}}_{3 \times 3} \underbrace{[g'(\hat{\beta})]}_{3 \times 1}$$

Example: non-linear function of coefficients

```
reg.coeff = reg$coeff
reg.vc = vcov(reg)

# Compute the estimate of g(beta)
g.est = 1 + reg.coeff[3]/reg.coeff[2]
# Define the vector of the derivative of g(beta) wrt beta
g.prime = matrix(c(0,-reg.coeff[3]/reg.coeff[2]^2,1/reg.coeff[2]),nrow=3,ncol=1)
g.prime

##           [,1]
## [1,] 0.000000
## [2,] 2.883012
## [3,] 3.211236

# Compute the variance of g(beta.hat)
g.var = t(g.prime) %*% reg.vc %*% g.prime
g.est

##      agesp6
## 0.1022112

g.est - qt(0.975,df=summary(reg)$df[2]) * sqrt(g.var)

##           [,1]
## [1,] 0.02689796

g.est + qt(0.975,df=summary(reg)$df[2]) * sqrt(g.var)

##           [,1]
## [1,] 0.1775244
```

$1 + \hat{\beta}_2 / \hat{\beta}_1 = 0.10$
95% CI for $1 + \beta_2 / \beta_1$
0.027 to 0.18

Comparing nested MLR models

model of interest

$$Y_i = \beta_0 + \underbrace{\beta_1 X_{1i} + \dots + \beta_p X_{pi}}_{p+1} + \underbrace{\beta_{p+1} X_{p+1i} + \dots + \beta_{p+s} X_{p+si}}_s + \varepsilon_i$$

extended
model

Interested in comparing the extended model to
simpler model nested within the extended model
null model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_p X_{pi}$$

$$H_0: \beta_{p+j} = 0, \text{ for } j=1, \dots, s$$

$$H_A: \text{at least one } \beta_{p+j} \neq 0$$

F-test for nested models; ANOVA method

Define:

- ▶ $R_N = R_0 = (I - H_N)Y \rightarrow$ residual from fitting the null model
- ▶ $R_E = R_E = (I - H_E)Y \rightarrow$ residuals from fitting the extended model
- ▶ $\Delta = R_N - R_E$

You can show the following results (which we will not do in class):

- $H_E - H_N$ is idempotent with rank (s)
- $H_E - H_N$ is orthogonal to $(I - H_E)Y$
- $\frac{\Delta' \Delta / s}{R_E' R_E / (n - p - s - 1)} \sim \mathcal{F}_{df1=s, df2=n-p-s-1}$

$$\frac{[SS_{\text{residual null}} - SS_{\text{residual ext}}] / s}{SS_{\text{residual extended}} / (n - p - s - 1)}$$

$$SS_{\text{residual extended}} / (n - p - s - 1)$$

Examples: nested MLR model comparisons

Consider the medical expenditure data you are analyzing for Problem Set 2. Define $\underline{Y} = \log(\text{medical expenditures} + 1)$ and let $\underline{X_1} = \text{age} - 65$ and $\underline{X_2} = \text{male}$ (indicator 1 = male, 0 = female). Define three models:

Model	Xs	residual df	SS(residual)
A	X_1, X_2	5691	31332.38
B	$X_1, (X_1 - 10)^+, (X_1 - 20)^+, X_2$	5689	31314.59
C	$[X_1, (X_1 - 10)^+, (X_1 - 20)^+] \times X_2$	5686	31299.23

model A: $E(Y | X_1, X_2) = \beta_0 + \beta_1(\text{age} - 65) + \beta_2 \text{male}$

model B $E(Y | X_1, X_2) = \beta_0 + \beta_1(\text{age} - 65) + \beta_2 \text{male}$
 $+ \beta_3 (\text{age} - 75)^+ + \beta_4 (\text{age} - 85)^+$

model C $E(Y | X_1, X_2) = \text{model B} +$
 $\beta_5 (\text{age} - 65) \text{male} + \beta_6 (\text{age} - 75)^+ \text{male}$
 $+ \beta_7 (\text{age} - 85)^+ \text{male}$

Example 1: nested MLR model comparisons

After adjusting for gender, is the average log expenditure a linear function of age?

null model A
ext model B

H0: $\beta_3 = 0$ and $\beta_4 = 0$

HA: at least β_3 or $\beta_4 \neq 0$

Model	Xs	residual df	SS(residual)	MS	F
A	X_1, X_2	5691	31332.38		
B	$X_1, (X_1 - 10)^+, (X_1 - 20)^+, X_2$	5689	31314.59	5.50	
Change		2	17.79	8.90	$\frac{8.90}{5.50} = 1.62$

Compute the P-value as: $Pr(\mathcal{F}_{2,5689} > 1.62) = 0.199$.

$$F = \frac{[31332.38 - 31314.59]/2}{31314.59/5689} = 1.62 \quad F(2, 5689)$$

Example 1: nested MLR model comparisons

```
load("C:\\Users\\Elizabeth\\Dropbox\\Biostat6532020\\Problem Set 2\\nmes.rdata")
d = nmes %>% select(names(.)[c(1,2,3,15)]) %>% filter(.,lastage>=65)
d = mutate(d,
  logy = log(totalexp+1),
  agec=lastage-65,
  agesp1 = ifelse(lastage-75>0, lastage-75,0),
  agesp2 = ifelse(lastage-85>0, lastage-85,0)
)
reg0 = lm(logy~agec+male,data=d) A
reg1 = lm(logy~agec+agesp1+agesp2+male,data=d) B
reg2 = lm(logy~(agec+agesp1+agesp2)*male,data=d) C
```

Questoin 1: using anova function

```
anova(reg0,reg1)
```

null extended

```
## Analysis of Variance Table
```

```
##
```

```
## Model 1: logy ~ agec + male
```

- null

```
## Model 2: logy ~ agec + agesp1 + agesp2 + male
```

- ext

```
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
```

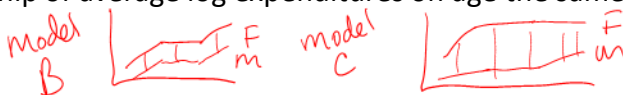
```
## 1     5691 31332
```

```
## 2     5689 31315  2      17.79 1.6159 0.1988
```

underlined

Example 2: nested MLR model comparisons

Is the non-linear relationship of average log expenditures on age the same for males and females? i.e. are the curves parallel?



- Equivalently: Is the difference between the average log expenditure for males and females the same at all ages?

H0: $\beta_5 = 0, \beta_6 = 0, \beta_7 = 0$

HA: at least one is non-zero

Model	Xs	residual df	SS(residual)
A	X_1, X_2	5691	31332.38
B	$X_1, (X_1 - 10)^+, (X_1 - 20)^+, X_2$	5689	31314.59
C	$[X_1, (X_1 - 10)^+, (X_1 - 20)^+] \times X_2$	5686	31299.23

$$F = \frac{[31314.59 - 31299.23]/3}{31299.23/5686}$$

Example 2: nested MLR model comparisons

Model	Xs	residual df	SS(residual)
A	X_1, X_2	5691	31332.38
B	$X_1, (X_1 - 10)^+, (X_1 - 20)^+, X_2$	5689	31314.59
C	$[X_1, (X_1 - 10)^+, (X_1 - 20)^+] \times X_2$	5686	31299.23

Question 2:

```
anova(reg1, reg2)
```

Analysis of Variance Table

##

Model 1: $\text{logy} \sim \text{agec} + \text{agesp1} + \text{agesp2} + \text{male}$

Model 2: $\text{logy} \sim (\text{agec} + \text{agesp1} + \text{agesp2}) * \text{male}$

Res.Df RSS Df Sum of Sq F Pr(>F)

1 5689 31315

2 5686 31299 3 15.36 0.9301 0.4252

Likelihood ratio tests for nested MLR models

Let \loglike_{ext} and \loglike_{null} be the values of the log likelihoods evaluated at the parameter estimates from the extended and null models, respectively.

Then to test H_0 :

Compute $2 \times \loglike_{ext} - 2 \times \loglike_{null} \sim \chi$, df = s

likelihood ratio
tests

\Rightarrow perform
better
type I, power

Compared to
model F test
when n is small

Examples: nested MLR model comparisons using LRT

Question 1: by hand

```
lr.test.stat = as.numeric(2 * logLik(reg1) - 2 * logLik(reg0))  
pchisq(lr.test.stat,df=2,lower.tail=FALSE)
```

```
## [1] 0.1985122
```

Question 1: Using lrtest function

```
#install.packages(lmtest)
```

```
library(lmtest)
```

```
lrtest(reg0,reg1)
```

```
## Likelihood ratio test
```

```
##
```

```
## Model 1: logy ~ agec + male
```

```
## Model 2: logy ~ agec + agesp1 + agesp2 + male
```

```
##   #Df LogLik Df   Chisq Pr(>Chisq)
```

```
## 1    4 -12934
```

```
## 2    6 -12933  2 3.2338    0.1985
```


Next time....

- ▶ Model checking for MLR models
- ▶ Key extensions for MLR models

