

## Question 2

2010, MS-2

$n = 20$  patients.       $n_1 = 10$  std  
                                  $n_2 = 10$  trt

a) Kaplan

r

c

$$a) E(\hat{\delta}_1 - \hat{\delta}_c) = E(\hat{\delta}_1) - E(\hat{\delta}_c)$$

$$= E(\bar{y}_f - \bar{y}_b) - E(\bar{x}_f - \bar{x}_b)$$

$$= (E\bar{y}_f - E\bar{y}_b) - (E\bar{x}_f - E\bar{x}_b)$$

$$= (\bar{y}_f - \bar{y}_b) - (\bar{x}_f - \bar{x}_b) = \delta_1 - \delta_c$$

$\Rightarrow$  unbiased

b)  $s^2$  is the same at baseline and F-U- for both samples

$$\text{Var}(\hat{\delta}_1 - \hat{\delta}_c) = \text{Var}\hat{\delta}_1 + \text{Var}\hat{\delta}_c - 2\text{cov}(\hat{\delta}_1, \hat{\delta}_c)$$

$$\text{cov}(\bar{y}_f, \bar{y}_b)$$

=

$$\rho = \frac{\text{cov}(x, y)}{\sqrt{\text{var } x \cdot \text{var } y}}$$

$$\Rightarrow \text{Var}(\hat{\delta}_1) = \text{Var}(\bar{y}_f) + \text{Var}(\bar{y}_b) - 2\text{cov}(\bar{y}_f, \bar{y}_b)$$

$$= \frac{s^2}{n} + \frac{s^2}{n} - 2 \frac{\rho_y \sqrt{s^2(s^2)}}{n}$$

=

Question 4 $n = 25$  homesdata on  $PM_{10}$  taken  $\Rightarrow$  use log $X_{Ii}$ ,  $X_{Oi}$  for  $i$  homes, the ln of  $PM_{10}$  [ ] indoors + outdoorsa) Estimate for the association btwn log outdoor and indoor  $PM_{10}$ 

$$\text{correlation} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$$

corr =  $\beta$  estimate of standardized,  
no int regression

$$= \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}} = \frac{\sum X_i Y_i - \sum X_i \bar{Y} - \sum Y_i \bar{X} + \bar{X} \bar{Y}}{\sqrt{(\sum X_i^2 - n\bar{X}^2)(\sum Y_i^2 - n\bar{Y}^2)}}$$

$$= 425 - \left( 95 \left( \frac{110}{25} \right) - 100 \left( \frac{95}{25} \right) + \left( \frac{110}{25} \right) \left( \frac{95}{25} \right) \right) \quad \text{complex}$$

$$\sqrt{\left( 390 - 25 \left( \frac{95}{25} \right)^2 \right) \left( 520 - 25 \left( \frac{110}{25} \right)^2 \right)}$$

$$= \frac{n(\sum XY) - \sum X \sum Y}{\sqrt{(n\sum X^2 - (\sum X)^2)(n\sum Y^2 - (\sum Y)^2)}} = 0.2676$$

Is using correlation  
right?

$$\text{test} \quad \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{n-2}$$

b)  $H_0: \mu_{out} = \mu_{in}$

$\bar{x}_{out} = \frac{110}{25} = 4.4$   $S_{out}^2 = \frac{1}{25} \left( \sum x_i^2 - \frac{(\sum x_i)^2}{n} \right) = 36$

$\bar{x}_{in} = \frac{95}{25} = 3.8$   $S_{in}^2 = \frac{1}{25} \left( \sum x_i^2 - \frac{(\sum x_i)^2}{n} \right) = 19$

$$\begin{aligned} \sum (x_i - \bar{x})^2 &= \sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ &= \sum x_i^2 - 2n\bar{x}^2 + n\bar{x}^2 \\ &= \sum x_i^2 - n\bar{x}^2 \end{aligned}$$

t-test for 2 sample means

$$\frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} =$$

\* I'm assuming the samples are independent  
~~\* assuming homogeneity of variances and~~

\* check what assumptions are necessary →

seems like 2 cov should be required too.

c) median b/c skewed.

- d) i. We would be missing potentially 3 low values of pollutant, and this could bias our results toward higher values
- ii. use a 'lower bound' or code missing as 0 or some other low number.

e) 400 children  $\cdot$   $PM_{10}$  [ ] in homes (quantiles)  
Y/N recent wheezing

$\Rightarrow$  The one-sided p-value indicates that the proportion of children wheezing increases w/ increasing tertile  
( $p = 0.0398$ , not highly significant).

(Two-sided test FTR).

---

f) dichotomized children.

estimate risk relative to other + 95% CI  
based on this is the risk sign?

$$\text{risk in } > 50 \text{ mg/m}^3 = \frac{26}{112}$$

$$\leq 50 \text{ mg/m}^3 = \frac{29}{288}$$

$$RR = 2.3054$$

$$CI\ 95\% = (1.4233, 3.7342)$$

$\Rightarrow$  limits do not include 1.

risk is significantly greater.

12010, MS-2

$n = 400$      $n_1 = 260$  men  $\begin{cases} 100 \text{ drug} \\ 100 \text{ none} \end{cases}$

ht, wt, age

<u>drug</u>	<u>gender</u>	<u><math>E(y_i)</math></u>
= 0	= 0	$\beta_0 + \beta_1(\text{age}) +$
= 0	= 1	$\beta_0 + \beta_1(\text{age}) + \beta_3$
= 1	= 0	$\beta_0 + \beta_1(\text{age}) + \beta_2$
= 1	= 1	$\beta_0 + \beta_1(\text{age}) + \beta_2 + \beta_3 + \beta_4$

$$b) \vec{y} = \vec{x} \vec{\beta} + \vec{e}$$

i.  $y = 400 \times 1$        $x = 400 \times 5$        $\beta = 5 \times 1$        $e = 400 \times 1$

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

iii  $\hat{\beta} = (b, \sigma^2(x'x)^{-1})$

c)	<u>source</u>	<u>df</u>	<u>SS</u>	<u>MS</u>	<u>F-val</u>	<u>p-val</u>
	model	$p-1 = 4$	66308.4	16577.1	146.7	0.0001
	error	$n-p = 395$	44635	113		
	C-total	$n-1 = 399$				

$$F = \frac{MSM}{MSE}$$

d) Add ht, wt to model.

① ANOVA test compares (1) and (2) shows wt/ht are important  $\rightarrow$   
 $\Rightarrow$  means the larger model explains more variation in LDL levels

② Insignificant Wald tests

$\rightarrow$  both coefficients are not significantly diff from 0 in added last tests.

$\Rightarrow$  No, they don't necessarily contradict each other. The  $t$ -tests are assessing whether height and weight <sup>each</sup> add significant info to the model in the presence of all other variables, while the overall  $F$ -test is assessing whether the two variables jointly improve the overall model fit.

# Question 3

2010, MS2

Drug for LDL reduction.

$n = 400$

$\Rightarrow 200$  men /  $200$  women

$\Rightarrow$   $\begin{matrix} 100 & 100 \\ \swarrow & \searrow \\ \text{drug} & \text{placebo} \end{matrix}$   $\begin{matrix} 100 & 100 \\ \swarrow & \searrow \\ \text{drug} & \text{placebo} \end{matrix}$

$$y = \beta_0 + \overset{1}{\text{age}} + \overset{2}{\text{trt}} + \overset{3}{\text{gender}} + \overset{4}{\text{trt} \times \text{gender}}$$

$\downarrow$   $\downarrow$   $\downarrow$   
 $= 1$  if  $= 1$  if  
 drug man

a)

<u>drug</u>	<u>gender</u>	<u><math>E(y_i)</math></u>
plac.	f	$\beta_0 + \beta_1$
plac.	m	$\beta_0 + \beta_1 + \beta_3$
drug.	f	$\beta_0 + \beta_1 + \beta_2$
drug.	m	$\beta_0 + \beta_1 + \beta_2 + \beta_3 + \beta_4$

b) i.  $y = X\beta + \epsilon$   $y = 400 \times 1$   $X = 400 \times 5$   $\beta = 5 \times 1$   $\epsilon = 400 \times 1$

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

iii.  $\hat{b} \sim N(b, \sigma^2(X'X)^{-1})$

		<u>SS</u>	<u>MS</u>	<u>SS/df</u>	<u>MSM</u> <u>MSE</u>	<u>p-val</u>
c)	<u>source</u>	<u>df</u>				
	model	4	66308.4	16577.1	146.7	<0.0001
p-1	error	395	44635	113		
n-p	C-total	399	110943.4			
	n-1					



d) Now add weight and height to the model

$$y = \beta_0 + \beta_1 \text{age} + \beta_2 \text{trt} + \beta_3 g + \beta_4 \text{trt} * g + \beta_5 \text{wt} + \beta_6 \text{ht}$$

test comparing 2 models  $\Rightarrow p = 0.0098$  (wt/ht is important)

but wt and ht are insignificant in T3 testing

$\Rightarrow$  Does not contradict

the overall ANOVA test indicates that there is more info explained by the larger model overall.

wt and ht are not significant, after all the other vars are in the model.

wt is not very insignif. should be careful.

e) i. Is drug effect diff for m/f

$\Rightarrow$  yes. interaction is significant. (effect of drug differs by m/f)

ii. 'Drug has no effect in both males and females.'  $\Rightarrow$  INTERPRETING as: no effect of the drug overall

$$H_0: \beta_2 = \beta_4 = 0$$

$$C = (0 \ 0 \ 1 \ 0 \ 1 \ 0 \ 0) ; \theta_0 = (0, 0)^T$$

$$\text{iii. } E(y) = \beta_0 + \beta_1(\text{age}) + \beta_2(\text{trt}) + \beta_5(\text{wt}) + \beta_6(\text{ht}) \Leftarrow \text{fem } (g=0)$$

$$E(y) = \beta_0 + \beta_1(\text{age}) + \beta_2(\text{trt}) + \beta_3(g) + \beta_4(\text{trt} * g) + \beta_5(\text{wt}) + \beta_6(\text{ht})$$

f) No. conditional distr.  $y|x=x$  matters.

$\uparrow$   
males ( $g=1$ )

e)

i. Drug effect diff btwn males / females?

⇒ yes, the interaction is significant ( $p = 0.000619$ )

ii. Drug has no effect in both males and females

⇒ drug has no effect

⇒  $H_0: \beta_{\text{drug}} = \beta_{\text{drug} \times \text{gender}} = 0$

$$C = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad q = (0 \ 0)^T$$

iii. linear model for males / females

$$\text{MALES} \Rightarrow E(Y_i) = (\beta_0 + \beta_3) + \beta_1(\text{age}) + (\beta_2 + \beta_4)(\text{trt}) + \beta_5(\text{wt}) + \beta_6(\text{wt})$$

$$\text{Female } E(Y_i) = (\beta_0) + \beta_1(\text{age}) + \beta_2(\text{trt}) + \beta_5(\text{wt}) + \beta_6(\text{wt})$$

---

f) No. linear regression does not require any assumptions about the distribution of  $Y$ . It is important for the distribution of  $Y$  conditional on all the parameters to be normal, but that is not what is shown here.

Question 1

$$\vec{y} = \vec{X}\vec{\beta} + \vec{\varepsilon}$$

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

a)

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & -1 \\ 1 & 3 & -2 \\ 1 & -1 & 2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \end{bmatrix}$$

b) 2 (column rank).

Not full rank:  $C1 - C2 = C3$ c) Re-express the model as  $y = X^* \beta^* + \varepsilon$  where  $X^*$  is full ranksquare  $X_2$ 

$$X = \begin{pmatrix} 1 & 1^0 & 0 \\ 1 & 2 & 1 \\ 1 & 3 & 4 \\ 1 & -1 & 4 \end{pmatrix}$$

~~center around means?~~  
~~around a random~~  
~~value?~~

d)  $\theta_1 = \beta_1 - \beta_2$  estimable (?) To be estimable  $\rightarrow C = TX$

$$C = (0 \ 1 \ -1)_{1 \times p} = T_{1 \times n} \begin{pmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 1 & 3 & 4 \\ 1 & -1 & 4 \end{pmatrix}_{n \times p}$$

$$T = (a_1 \ a_2 \ a_3 \ a_4)$$

$$(0 \ 1 \ -1) = (a_1 + a_2 + a_3 + a_4, a_1 + 2a_2 + 3a_3 - a_4, a_2 + 4a_3 + 4a_4)$$

$$0 = a_1 + a_2 + a_3 + a_4$$

$$1 = a_1 + 2a_2 + 3a_3 - a_4$$

$$-1 = a_2 + 4a_3 + 4a_4$$

$$0 = \sum_{i=1}^4 a_i$$

come up w/ some combo / matrix  
for  $a$  that works.

$$\text{need } \Rightarrow X\%t = c$$

$$c = t \% X$$

$$\text{e.g. } T = (1 \ -1 \ 0 \ 0) \Rightarrow \theta_1 \text{ estimable}$$

$$\hat{\theta}_1 = c\hat{\beta} \quad \hat{\beta} = (X'X)^{-1}X'y$$
$$\text{Var } \hat{\theta}_1 = \sigma^2 c(X'X)^{-1}c' \rightarrow \text{MSE}$$

how would you do  
this by hand?