

## Stat 421, Test 2, Fall, Nov. 16, 2015; Marzban

9 + 25

ONLY a half-size "cheat sheet" is allowed

Multiple choice: Circle all the correct answers; there is wrong-answer penalty

For rest, SHOW answer &amp; work; NO CREDIT for correct answer without explanation

Points

1. Which of the following statements is/are true regarding contrasts?

- 1
- a) All comparisons of  $\mu_i$  can be written in terms of zero-sum contrasts constructed from the  $\mu_i$ .  
 b) A specific comparison between several  $\mu_i$ , constructed from a zero-sum contrast, can be written in terms of a zero-sum contrast constructed from the corresponding effects  $\alpha_i = \mu - \mu_i$ .  
 c) For  $a$  treatment levels, there exists a **unique** set of  $(a - 1)$  orthogonal contrasts.  
 d) If an anova F-test has led to the rejection of the null hypothesis of equal means, then testing a set of orthogonal contrasts can help in identifying a specific combination of means that is "responsible" for the rejection.
- 0.5 each

~ 1

2. In using the maximum likelihood criterion for estimating model parameters, if we change the constraints, then \_\_\_\_\_ of the parameter estimates will change.

- a) none    b) some    c) all

~ 1

3. In using the maximum likelihood criterion for estimating model parameters, if we change the constraints, then \_\_\_\_\_ functions of the parameters will change.

- a) no    b) some    c) all

1

4. Suppose we have developed the model  $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$   $i = 1 \dots a$ ,  $j = 1 \dots n$  based on a CRD. Circle all of the following quantities that will be different in an RCBD model  $y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$  of the same data.

- a) SST    b) SSA    c) SSE
- 0.5 each     $SSA = \sum (\bar{y}_{i.} - \bar{y}_{..})^2$  in both models.

~ 1

5. In the model  $y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$   $i = 1 \dots a$ ,  $j = 1 \dots b$  based on an RCBD design, how does  $F_A$  generally change with increasing  $b$ ?

- a) Generally decreases    b) Generally remains constant    c) Generally increases    d) None of the above.
- \* c)

6. In which of the following models should we be concerned if the "grand mean" of the residuals turns out to be nonzero?

- 2
- a)  $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$     b)  $y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$     c)  $y_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ijk}$     d)  $y_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \epsilon_{ijk}$     e)  $y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_k + \epsilon_{ijk}$
- all models have  $e_{....} = 0$ .  
0.5 each, max = 2

7. In which of the following models should we not be surprised if some "conditional mean" of the residuals turns out to be nonzero?

- 1
- a)  $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$     b)  $y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$     c)  $y_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ijk}$     d)  $y_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \epsilon_{ijk}$     e)  $y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_k + \epsilon_{ijk}$
- 0.5 each     $e_{i.} = 0$      $e_{i.} \neq 0$      $e_{i.} \neq 0$  (let 13)

~ 1

8. In a problem with quantitative factors, one can develop a regression model that has a comparable number of parameters as the anova model. This statement is

- a) always true, period.  
 b) true, but the highest possible order depends on the number of factors in the problem.  
 c) true, but the highest possible order depends on the number of levels in the factors.  
 d) never true, period.

5)  $F = \frac{MSA}{MSE} = \frac{\sum_i (\bar{y}_{i.} - \bar{y}_{..})^2 / (a-1)}{\sum_{i,j} (y_{ij} - \dots)^2 / (ab(n-1))}$      $\sim b(n-1)$

\* The more levels in B, the easier to find A effect!

Vadav example

9. We are preparing to conduct a study to find if IQ depends on gender and on whether or not one wears glasses. Due to practical limitations, the data must be collected on 4 different days, and in 4 different cities, and it will be impossible to randomize across the days and the cities.

a) In words, what kind of design is most appropriate if we have sufficient resources to conduct  $4^3$  runs? Make sure you also define the treatment factors, and the block factors, if any.

RCBD factorial, with Gender (B/G) and Glasses (Yes/No) as 2 treatment factors, and Day and City as 2 block factors.

b) What if we have sufficient resources to conduct only 16 runs?

LSD with Day and City as 4-level row and col. factors, and combination of Gender and Glasses (B,Y), (B,N), (G,Y), (G,N) as

The 4-level treatment factor.

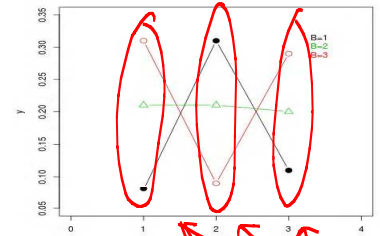
10. Consider the data (shown in the adjacent figure) on the response  $y$  and two 3-level factors A and B.

a) Suppose, initially we are not given the information on B at all, i.e., we have no idea that it even exists. Do we expect to find that A has an effect on the response?

Yes/No, explain in words.

No.

Without any info on B, there is too much overlap over the data in the 3 categories/populations.



b) Now, suppose it is revealed to us that B affects the data as shown in the figure. If we fit a model that continues to ignore B, do we expect to find that A has an effect? Yes/No, explain in words.

No.

A model without the B term is the same as the model above.

c) Write down the appropriate model if all the information in the figure is provided to you, but our main focus is the effect of A on the response. Explain your reasoning for the proposed model.

$$y_{ij} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ij}$$

more indices if repl.

There is clearly an interaction between A, B because the effect of A depends on the level of B.

11. Consider a CRD and the model  $y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_k + \epsilon_{ijk}$ , where  $i = 1 \dots a$ ,  $j = 1 \dots b$ ,  $k = 1 \dots c$ . Show that the sum of the predictions is equal to the sum of the observed values of  $y_{ijk}$ .

$$\hat{y}_{ijk} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + \hat{\gamma}_k = \bar{y}_{...} + (\bar{y}_{i..} - \bar{y}_{...}) + (\bar{y}_{.j.} - \bar{y}_{...}) + (\bar{y}_{...k} - \bar{y}_{...})$$

$$\frac{1}{bc} y_{i..} = \bar{y}_{i..} + \bar{y}_{.j.} + \bar{y}_{...k} - 2\bar{y}_{...}$$

$$\sum_{i,j,k} \hat{y}_{ijk} = \sum_{i,j,k} \bar{y}_{i..} + \sum_{i,j,k} \bar{y}_{.j.} + \sum_{i,j,k} \bar{y}_{...k} - 2\bar{y}_{...} \sum_{i,j,k} 1$$

$$= \frac{1}{bc} \sum_{i,j,k} y_{i..} + \frac{1}{ac} \sum_{i,j,k} y_{.j.} + \frac{1}{ab} \sum_{i,j,k} y_{...k} - 2 \frac{1}{abc} y_{...} \sum_{i,j,k} 1$$

$$= y_{...} + y_{...} + y_{...} - 2 y_{...} = y_{...}$$

~ 2

hom and  
Sapienter

12. Show that  $y_{1..} + y_{2..} + y_{3..} = y_{...}$  where all "dots" refer to restricted sums over the LSD given by

$$y_{1..} = y_{111} + y_{122} + y_{133}$$

$$y_{2..} = y_{221} + y_{232} + y_{213}$$

$$y_{3..} = y_{331} + y_{312} + y_{323}$$

$$\begin{pmatrix} 111 & 122 & 133 \\ 221 & 232 & 213 \\ 331 & 312 & 323 \end{pmatrix}$$

	$B_1$	$B_2$	$B_3$
$A_1$	$C_1$	$C_2$	$C_3$
$A_2$	$C_2$	$C_3$	$C_1$
$A_3$	$C_3$	$C_1$	$C_2$

$$\therefore y_{1..} + y_{2..} + y_{3..} = \text{sum of all 9 entries} \\ = y_{...}$$

~ 3

25

Let us

13. Consider the model  $y_{ij} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \epsilon_{ij}$ , where  $i = 1 \dots a$ ,  $j = 1 \dots b$ , and  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ .

a) Starting from the expression for the Likelihood of data, compute the maximum-likelihood (ML) estimate of the interaction term. You may use the fact that the ML estimates of  $\mu, \alpha_i, \beta_j$  are, respectively,  $\bar{y}_{..}$ ,  $(\bar{y}_{i.} - \bar{y}_{..})$  and  $(\bar{y}_{.j} - \bar{y}_{..})$ .

$$L = \prod_{i,j} \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} e^{-\frac{1}{2} \left( \frac{y_{ij} - \mu - \alpha_i - \beta_j - (\alpha\beta)_{ij}}{\sigma_\epsilon} \right)^2} = e^{-\frac{1}{2} \left( \sum_{i,j} \frac{1}{2\sigma_\epsilon^2} (y_{ij} - \mu - \alpha_i - \beta_j - (\alpha\beta)_{ij})^2 - \sum_{i,j} \frac{1}{2} \ln(2\pi\sigma_\epsilon^2) \right)}$$

$$\frac{\partial}{\partial (\alpha\beta)_{kl}} : \frac{\partial}{\partial (\alpha\beta)_{kl}} \sum_{i,j} (y_{ij} - \mu - \alpha_i - \beta_j - (\alpha\beta)_{ij})^2 = -2 (y_{kl} - \mu - \alpha_k - \beta_l - (\alpha\beta)_{kl})$$

$$\therefore (\hat{\alpha\beta})_{ij} = y_{ij} - \hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j = y_{ij} - \bar{y}_{..} - (\bar{y}_{i.} - \bar{y}_{..}) - (\bar{y}_{.j} - \bar{y}_{..}) \\ = \underline{y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..}}$$

~ 2

b) Find the expression for the predictions  $\hat{y}_{ij}$ .

$$\hat{y}_{ij} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + (\hat{\alpha\beta})_{ij} \\ = \cancel{\hat{\mu}} + \cancel{\hat{\alpha}_i} + \cancel{\hat{\beta}_j} + y_{ij} - \cancel{\hat{\mu}} - \cancel{\hat{\alpha}_i} - \cancel{\hat{\beta}_j} \\ = \underline{y_{ij}}$$

14. Consider a CRD with a single qualitative factor A, and the model  $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$ ,  $i = 1 \dots a$ ,  $j = 1 \dots n$ . We want to test whether A has an effect on the response. We know that  $MS_{Tr} = \frac{\sum_{i,j} (\bar{y}_{i.} - \bar{y}_{..})^2}{a-1}$ . Compute  $SSE_{reduced}$  and  $SSE_{full}$ , and then show that  $MS_{Tr} = \frac{SSE_{reduced} - SSE_{full}}{a-1}$ . You may use the anova decomposition without deriving/proving it.

Reduced:  $y_{ij} = \mu + \epsilon_{ij} \Rightarrow \hat{\mu} = \bar{y}_{..}$

$\therefore \hat{y}_{ij} = \hat{\mu} = \bar{y}_{..} \Rightarrow SSE_{reduced} = \sum_{i,j} (y_{ij} - \bar{y}_{..})^2$

Full:  $y_{ij} = \mu + \alpha_i + \epsilon_{ij} \Rightarrow \hat{\mu} = \bar{y}_{..}, \hat{\alpha}_i = \bar{y}_{i.} - \bar{y}_{..}$

$\therefore \hat{y}_{ij} = \hat{\mu} + \hat{\alpha}_i = \bar{y}_{i.} \Rightarrow SSE_{full} = \sum_{i,j} (y_{ij} - \bar{y}_{i.})^2$

$\underline{SSE_{red} - SSE_{full}} = \underbrace{\sum_{i,j} (y_{ij} - \bar{y}_{..})^2}_{SST} - \underbrace{\sum_{i,j} (y_{ij} - \bar{y}_{i.})^2}_{SSE} = \sum_{i,j} (\bar{y}_{i.} - \bar{y}_{..})^2 = \underline{SS_{Tr}}$

ANOVA decomposition

15. Consider the following data on a response y (numbers in the table) from two treatment levels of a factor A.

A	Runs		
1	1	2	3
2	2	3	4

a) Suppose the data have been collected in a CRD. Use an appropriate model and report the numerical value of SSE. Hint: find  $\bar{y}_{i.}$  first.

$y_{ij} = \mu + \alpha_i + \epsilon_{ij} \quad a=2, n=3 \Rightarrow \hat{y}_{ij} = \bar{y}_{i.} = \begin{cases} \frac{1}{3}(1+2+3) = 2 & i=1 \\ \frac{1}{3}(2+3+4) = 3 & i=2 \end{cases}$

$SSE = \sum_{i,j} (y_{ij} - \bar{y}_{i.})^2 = (1-2)^2 + (2-2)^2 + (3-2)^2 \leftarrow i=1$   
 $+ (2-3)^2 + (3-3)^2 + (4-3)^2 \leftarrow i=2$   
 $= 1 + 0 + 1 + 1 + 0 + 1 = \boxed{4}$

b) Now, suppose the 3 Runs occurred on 3 different days, and there was no randomization between days. Use an appropriate RCBD model and report the numerical value of SSE. Hint: find  $\bar{y}_{.j}, \bar{y}_{..}$ .

$y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij} \Rightarrow \hat{y}_{ij} = \bar{y}_{i.} + \bar{y}_{.j} - \bar{y}_{..}, \quad \bar{y}_{.j} = \begin{cases} \frac{1}{2}(1+2) = 3/2 \\ \frac{1}{2}(2+3) = 5/2 \\ \frac{1}{2}(3+4) = 7/2 \end{cases}$

$\bar{y}_{..} = \frac{15}{6} = 5/2$

$SSE = \sum_{i,j} (y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2$

$= (1 - 2 - \frac{3}{2} + \frac{5}{2})^2 + (2 - 2 - \frac{5}{2} + \frac{5}{2})^2 + (3 - 2 - \frac{7}{2} + \frac{5}{2})^2 \leftarrow i=1$   
 $(2 - 3 - \frac{3}{2} + \frac{5}{2})^2 + (3 - 3 - \frac{5}{2} + \frac{5}{2})^2 + (4 - 3 - \frac{7}{2} + \frac{5}{2})^2 \leftarrow i=2$   
 $= \boxed{0}$