

53.483 DATA ANALYSIS I

Tuesday 22nd January 2008
10.00a.m. - 12noon

All questions may be attempted.
Credit will be given for the best **THREE** answers only.

1.(a) The mathematical model for the response variable X_{ij} in a completely randomised design with k treatments and n replications per treatment can be described by

$$X_{ij} = \mu + \tau_j + \varepsilon_{ij}, \quad 1 \leq i \leq n, \quad 1 \leq j \leq k,$$

where μ is the grand mean, τ_j the effect of treatment j and ε_{ij} the experimental error. Assume that treatments are fixed effects, $\sum \tau_j = 0$. Assume also that ε_{ij} are independent and follow a $N(0, \sigma_e^2)$ distribution. Give the definitions for SS_{treat} and MS_{treat} . Show that

$$E(MS_{treat}) = \sigma_e^2 + n\sigma_T^2,$$

where $\sigma_T^2 = (\tau_1^2 + \dots + \tau_k^2)/(k-1)$. (You may use the following results without proof. Let Y_i ($1 \leq i \leq m$) be independent and normally distributed $N(0, \sigma^2)$ variables. Define $\bar{Y} = (Y_1 + \dots + Y_m)/m$. Then $\bar{Y} \sim N(0, \sigma^2/m)$ and

$$E\left[\frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y})^2\right] = \sigma^2.)$$

(13 marks)

(b) The partially completed ANOVA table for a $p \times q$ design is shown below:

Source	SS	d.f.	MS	F
A	3.06		1.02	
B	159.50	3		
A × B	11.79			
Error				
Total	180.44	31		

Assume A is a fixed effect but B is a random effect so this is a mixed model. Using a suitable notation, give a clear description of the model, including any assumptions necessary to undertake the analysis of variance. Fill in the blanks in the ANOVA table. Carry out the F -tests for interaction $A \times B$ and main effects A and B . Test using $\alpha = 0.05$.

(12 marks)

2. The following data were generated from a randomised block design with 4 treatments and 4 blocks:

	Block 1	Block 2	Block 3	Block 4
Treat 1	1.25	1.43	1.32	1.31
Treat 2	2.05	1.56	1.68	1.69
Treat 3	1.95	2.00	1.83	1.81
Treat 4	1.75	1.93	1.70	1.59

Assume that treatments and blocks are fixed effects. The classical mathematical model for the design is

$$X_{ij} = \mu + \beta_i + \tau_j + \varepsilon_{ij}, \quad 1 \leq j \leq 4, 1 \leq i \leq 4,$$

where μ is the grand mean, τ_j is the effect of treatment j , β_i is the effect of block i while ε_{ij} is the experimental error. The GLIM mathematical model for the design is

$$X_{ij} = \tilde{\mu} + \tilde{\beta}_i + \tilde{\tau}_j + \varepsilon_{ij}, \quad 1 \leq j \leq 4, 1 \leq i \leq 4.$$

(a) State the classical assumptions and the GLIM assumptions. State the relationships between the classical parameters and the GLIM parameters by expressing $\tilde{\mu}$, $\tilde{\beta}_i$ and $\tilde{\tau}_j$ in terms of μ , β_i and τ_j .

(3 marks)

(b) State the GLIM commands you would use to input the data in order to produce the following results:

(3 marks)

[i] ? \$fit\$

[o] deviance - 0.95584

[o] residual df = 15

[i] ? \$fit +block\$

[o] deviance = 0.89167 (change = -0.06417)

[o] residual df = 9 (change = -3)

[o]

[i] ? \$fit + treat\$

[o] deviance = 0.17296 (change = -0.7187)

[o] residual df = 9 (change = -3)

[o]

[i] ? \$ disp e s v r\$

```

[o] estimate      s.e.   parameter
[o] 1      1.399  0.09169      1
[o] 2     -0.02000 0.09802  BLOCK(2)
[o] 3     -0.1175  0.09802  BLOCK(3)
[o] 4     -0.1500  0.09802  BLOCK(4)
[o] 5     -0.4175  0.09802  TREAT(2)
[o] 6      0.5700  0.09802  TREAT(3)
[o] 7      0.4150  0.09802  TREAT(4)
[o] scale parameter 0.01922
[i] ? $tab the yield mean for treat$
[o] TREAT      1      2      3
[o] MEAN    1.327  1.745  1.898  1.742

```

(c) Based on the GLIM output, draw up the ANOVA table and state what conclusions can be drawn from it (tests using $\alpha = 0.05$).

(5 marks)

(d) Carry out the Newman-Keuls-MRT for the 4 treatments (test using $\alpha = 0.05$) and state your conclusions.

(7 marks)

(e) Based on the GLIM output, the standard error (s.e.) for the GLIM parameter $\tilde{\tau}_2$ is 0.09802. State (without proof) the least-squares estimator for the $\tilde{\tau}_2$ in terms of the observations, and show its variance from first principles and hence verify its s.e. = 0.09802.

(7 marks)

3. Data are provided on content uniformity of film-coated tablets produced for a cardiovascular drug used to lower blood pressure. A random sample of 3 batches was chosen at each of two manufacturing sites, and from each batch 5 tablets were examined. We wish to know if there are differences between the two sites (fixed effects) and between batches within sites (random effects) i.e. mixed model.

	Site					
	1			2		
	Batch		3	Batch		3
1	2	1		2		
Rep1	5.03	4.64	5.10	5.05	5.46	4.90
	5.10	4.73	5.15	4.96	5.15	4.95
	5.25	4.82	5.20	5.12	5.18	4.86
	4.98	4.95	5.08	5.12	5.18	4.86
	5.05	5.06	5.14	5.05	5.11	5.07

(a) State a linear model that may be used to represent the data and state the assumptions of the model.

(6 marks)

(b) State the GLIM commands you would use to input the data in order to produce the following results.

[i] ? \$fit\$

[o] deviance = 0.76247

[o] d.f. = 29

[o]

[i] ? \$fit + site\$

[o] deviance = 0.74421 (change = -0.01825)

[o] d.f. = 28 (change = -1)

[o]

[i] ? \$fit + site.batch\$

[o] deviance = 0.29020 (change = -0.4540)

[o] d.f. = 24 (change = -4)

(5 marks)

(c) Using the GLIM results provided, give the ANOVA table and undertake hypothesis tests to see if there are significant differences between the two sites and between batches within sites. Test using $\alpha = 0.05$.

(7 marks)

(d) Compute the 90% confidence interval for σ_e^2 (the variance of the experimental error).

(7 marks)

4.(a) In a 2^4 design with treatment A, B, C, D and 8 blocks and 3 replicates, the design is arranged so that AB, BC and CD are confounded with block effects in 3 replicates. Specify the units to go into each block and identify the composition of the ANOVA table.

(12 marks)

(b) In a 2^4 -design, 3 replications are carried out for each of the following half of the 16 treatment combinations:

(1) $b \ c \ bc \ ad \ abd \ acd \ abcd$

Identify the defining contrast and indicate aliases.

(9 marks)

(c) Identify the composition of the ANOVA table for part (b).

(4 marks)

Solutions to DA1 Exam 2007/8

$$1(a) \quad T_j = \sum_{i=1}^n X_{ij}, \quad \bar{T}_j = \frac{T_j}{n}$$

$$G = \sum_{i=1}^n \sum_{j=1}^k X_{ij}, \quad \bar{G} = \frac{G}{nk}$$

$$SS_{\text{treat}} = n \sum_{j=1}^k (\bar{T}_j - \bar{G})^2$$

$$MS_{\text{treat}} = \frac{SS_{\text{treat}}}{k-1}$$

$$\text{Note } \bar{T}_j = \frac{1}{n} \sum_{i=1}^n (\mu + \tau_j + \varepsilon_{ij}) \\ = \mu + \tau_j + \bar{e}_j$$

$$\text{where } \bar{e}_j = \frac{1}{n} \sum_{i=1}^n \varepsilon_{ij} \sim N(0, \frac{\sigma_e^2}{n})$$

$$\bar{G} = \frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k (\mu + \tau_j + \varepsilon_{ij}) \\ = \mu + \frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k \varepsilon_{ij} \\ = \mu + \bar{e}$$

$$\text{where } \bar{e} = \frac{1}{k} \sum_{j=1}^k \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_{ij} \right) = \frac{1}{k} \sum_{j=1}^k \bar{e}_j \\ \sim N(0, \frac{\sigma_e^2}{nk})$$

Compute

$$E(MS_{\text{treat}}) = \frac{n}{k-1} E \left[\sum_{j=1}^k (\mu + \tau_j + \bar{e}_j - \mu - \bar{e})^2 \right] \\ = \frac{n}{k-1} E \left[\sum_{j=1}^k (\tau_j^2 + 2\tau_j(\bar{e}_j - \bar{e}) + (\bar{e}_j - \bar{e})^2) \right]$$

$$= \frac{n}{k-1} \left[\sum_{j=1}^k \tau_j^2 + \sum_{j=1}^k 2\tau_j (\bar{E}e_j - \bar{E}\bar{e}) + E \sum_{j=1}^k (e_j - \bar{e})^2 \right]$$

$$= n \frac{\tau_1^2 + \dots + \tau_k^2}{k-1} + 0 + n E \left[\frac{1}{k-1} \sum_{j=1}^k (e_j - \bar{e})^2 \right]$$

$$= n \sigma_T^2 + n \cdot \frac{\sigma_e^2}{n} = n \sigma_T^2 + \sigma_e^2$$

1(b) This is a $p \times q$ design with $p = 4$ (A) and $q = 4$ (B)

$$X_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + E_{ijk}$$

$$1 \leq i \leq 4, 1 \leq j \leq 4, 1 \leq k \leq 2$$

where μ - grand mean

α_i - effect of factor A level i

β_j - effect of factor B level j

$\alpha\beta_{ij}$ - effect of interaction between factor A level i & factor B level j

E_{ijk} - experimental error

\therefore A is a fix effect while B is a random effect

\therefore This a mixed-effects model

Assumptions:

$$\alpha_i \text{ 's constants, } \sum_{i=1}^4 \alpha_i = 0$$

$$\beta_j \sim N(0, \sigma_B^2)$$

$$\alpha\beta_{ij} \sim N(0, \sigma_{AB}^2)$$

$$E_{ijk} \sim N(0, \sigma_e^2)$$

} independent

ANOVA

Source	SS	d.f.	MS	F
A	3.06	3	1.02	$\frac{1.02}{1.31} = 0.77$
B	159.50	3	53.17	$\frac{53.17}{0.38} = 139.92$
AB	11.79	9 (=3x3)	1.31	$\frac{1.31}{0.38} = 3.44$
Error	6.09	16	0.38	
Total	180.44	31		

Hypothesis tests

AB: $H_0: \sigma_{AB}^2 = 0$
 $H_a: \sigma_{AB}^2 > 0$

$$F = \frac{MS_{AB}}{MS_{Error}} = \frac{1.31}{0.38} = 3.44$$

$$F_{0.05}(9, 16) = 2.55 < F$$

\therefore reject H_0 , i.e. AB is significant

A: $H_0: \alpha_i = 0$
 $H_a: \alpha_i$'s differ

$$F = \frac{MS_A}{MS_{AB}} = \frac{1.02}{1.31} = 0.77$$

$$F_{0.05}(3, 9) = 3.86 > 0.77$$

\therefore don't reject H_0 , i.e. A is not sign.

B: $H_0: \sigma_B^2 = 0$
 $H_a: \sigma_B^2 > 0$

$$F = \frac{MS_B}{MS_{Error}} = \frac{53.17}{0.38} = 139.92$$

$$F_{0.05}(3, 16) = 3.24 < 139.92$$

\therefore reject H_0 , i.e. B is significant

2 @ The classical assumptions.

μ, β_i, τ_j are all constants,

$$\sum_{i=1}^4 \beta_i = 0, \quad \sum_{j=1}^4 \tau_j = 0$$

$E_{ij} \sim N(0, \sigma_e^2)$ independently

The GLIM assumptions:

$\tilde{\mu}, \tilde{\beta}_i, \tilde{\tau}_j$ are all constants,

$$\tilde{\beta}_1 = 0, \quad \tilde{\tau}_1 = 0$$

$$\tilde{\mu} = \mu + \beta_1 + \tau_1$$

$$\tilde{\beta}_i = \beta_i - \beta_1 \quad (2 \leq i \leq 4)$$

$$\tilde{\tau}_j = \tau_j - \tau_1 \quad (2 \leq j \leq 4)$$

6

```

[i] ? $uni 16$
[i] ? $data yields$
[i] ? $read
[i] $REA? 1.25 1.43 1.32 1.31
[i] $REA? 2.05 1.56 1.68 1.69
[i] $REA? 1.95 2.00 1.83 1.81
[i] $REA? 1.75 1.93 1.70 1.59
[i] ? $ca block=%gl(4,1)$
[i] ? $ca treat=%gl(4,4)$
[i] ? $factor block 4 treat 4$
[i] ? $look yield block treat$
[O]      YIELD  BLOCK  TREAT
[O] 1    1.250   1.000   1.000
[O] 2    1.430   2.000   1.000
[O] 3    1.320   3.000   1.000
[O] 4    1.310   4.000   1.000
[O] 5    2.050   1.000   2.000
[O] 6    1.560   2.000   2.000
[O] 7    1.680   3.000   2.000
[O] 8    1.690   4.000   2.000
[O] 9    1.950   1.000   3.000
[O] 10   2.000   2.000   3.000
[O] 11   1.830   3.000   3.000
[O] 12   1.810   4.000   3.000
[O] 13   1.750   1.000   4.000
[O] 14   1.930   2.000   4.000
[O] 15   1.700   3.000   4.000
[O] 16   1.590   4.000   4.000
[i] ? $yvar yields$

```


2(c)

ANOVA

Source	SS	df	MS	F
Block	0.06417	3	0.02139	1.11290
Treat	0.7187	3	0.23957	12.4646
Error	0.17296	9	0.01922	
Total	0.95583	12		

Hypothesis test for treatments.

H₀: τ₁ = τ₂ = τ₃ = τ₄ = 0

H_a: τ_j's differ

Given α = 0.05, F_{3,9(0.05)} = 3.86

∴ F = 12.4646 > F_{3,9(0.05)}

∴ reject H₀ and accept H₁, i.e. there is a significant difference among treatment means.

Hypothesis test for blocks.

H₀: β₁ = β₂ = β₃ = β₄ = 0

H_a: β_i's differ

∴ F = 1.1129 < F_{3,9(0.05)}

∴ don't reject H₀

3(d) MRT:

Treat	1	4	2	3
Mean	1.327	1.742	1.745	1.898

MSerror = 0.01922, d.f. = 9, s.e. of mean = sqrt(0.01922/4) = 0.06932

P	2	3	4
Range	3.2	3.95	4.42

P	2	3	4
LS range	0.2218	0.2738	0.3064

Tests

- 3 vs 1: $1.898 - 1.327 = 0.571 > 0.3064 \Rightarrow$ significant
- 3 vs 4: $1.898 - 1.742 = 0.156 < 0.2738 \Rightarrow$ not significant
- 3 vs 2: $1.898 - 1.745 = 0.153 < 0.2218 \Rightarrow$ not significant
- 2 vs 1: $1.745 - 1.327 = 0.418 > 0.2738 \Rightarrow$ significant
- 2 vs 4: $1.745 - 1.742 = 0.003 < 0.2218 \Rightarrow$ not significant
- 4 vs 1: $1.742 - 1.327 = 0.415 > 0.2218 \Rightarrow$ significant

Conclusions: treat 1 differs from treatments 2, 3, 4 while treatments 2, 3, 4 do not differ

2e) From 2a, we know $\tilde{\tau}_0 = \tau_2 - \tau_1$, while the l.s. estimators for τ_2 and τ_1 are \bar{T}_2 and \bar{T}_1 , respectively, so the l.s. estimator for $\tilde{\tau}_0$ is

$$\bar{T}_2 - \bar{T}_1$$

where $\bar{T}_2 = \frac{1}{4} \sum_{i=1}^4 X_{i2}$, $\bar{T}_1 = \frac{1}{4} \sum_{i=1}^4 X_{i1}$

Since X_{ij} 's are independent random variables with variance σ_e^2 ,

$$\begin{aligned} \text{Var}(\bar{T}_2 - \bar{T}_1) &= \text{Var}(\bar{T}_2) + \text{Var}(\bar{T}_1) \\ &= \frac{1}{16} \sum_{i=1}^4 \text{Var}(X_{i2}) + \frac{1}{16} \sum_{i=1}^4 \text{Var}(X_{i1}) \\ &= \frac{1}{16} \sum_{i=1}^4 \sigma_e^2 + \frac{1}{16} \sum_{i=1}^4 \sigma_e^2 \\ &= \frac{\sigma_e^2}{2} \end{aligned}$$

The s.e. of the l.s. estimator for $\tilde{\tau}_0$ is then $\sigma_e^2/\sqrt{2}$. Estimate σ_e by $s_e = \sqrt{MS_{\text{error}}} = 0.1386$.

Then the s.e. of the l.s. estimator for $\tilde{\tau}_0$ is

$$\frac{0.1386}{\sqrt{2}} = 0.09802$$

which is the same as the GLM output

3 a) This is a nested designs, since different batches are used at each site. The mathematical model is:

$$X_{ijk} = \mu + \alpha_i + (\beta_{j(i)} + \epsilon_{ijk},$$

$$1 \leq i \leq 2, \quad 1 \leq j \leq 3, \quad 1 \leq k \leq 5$$

where μ - grand mean

α_i - effect of site i

$\beta_{j(i)}$ - effect of batch j nested in ~~treatment~~ site i

ϵ_{ijk} - experimental error.

From the given conditions, we know that sites are fixed effects but batches are random effects. The corresponding assumptions are:

α_1 and α_2 are constants, $\alpha_1 + \alpha_2 = 0$

$$\left. \begin{aligned} \beta_{j(i)} &\sim N(0, \sigma_B^2) \\ \epsilon_{ijk} &\sim N(0, \sigma_e^2) \end{aligned} \right\} \text{independently.}$$

b)

```

o] ? $units 30
i] ? $data content$
i] ? $read
i] $REA? 5.03 5.10 5.25 4.98 5.05
i] $REA? 4.64 4.73 4.82 4.95 5.06
i] $REA? 5.10 5.15 5.20 5.08 5.14
i] $REA? 5.05 4.96 5.12 5.12 5.05
i] $REA? 5.46 5.15 5.18 5.18 5.11
i] $REA? 4.90 4.95 4.86 4.86 5.07
i] ? $ca site=%gl(2,15) :batch=%gl(3,5)$
i] ? $fac site 2 batch 3$
i] ? $yvar content$

```

}
 may use

\$ ca site = %gl(2,15) : batch = %gl(3,5) \$
 \$ fac site 2 batch 6 \$

5

c) ANOVA

Source	SS	d.f.	MS	F
Sites	0.01825	1	0.01825	$\frac{0.01825}{0.01209} = 0.16$
Batches within sites	0.4540	4	0.11350	$\frac{0.11350}{0.01209} = 9.39$
Error	0.2902	24	0.01209	
Total	0.76247	29		

7

Tests for

Sites: $F_{0.05}(1, 4) = 7.71 > 0.16$

⇒ Insignificant

Batches within sites: $F_{0.05}(4, 24) = 2.75 < 9.39$

⇒ Significant

(9)

d) Note from the ANOVA that

$$s_e^2 = MS_{\text{error}} = 0.01209$$

$$\text{d.f. for error} = 24$$

\therefore the 90% C.I. for σ_e^2 is

$$\frac{24 s_e^2}{\chi_{24}^2(0.05)} < \sigma_e^2 < \frac{24 s_e^2}{\chi_{24}^2(0.95)}$$

From the χ^2 -table,

$$\chi_{24}^2(0.05) = 36.42$$

$$\chi_{24}^2(0.95) = 13.85$$

\therefore The C.I. is

$$\frac{24 \times 0.01209}{36.42} < \sigma_e^2 < \frac{24 \times 0.01209}{13.85}$$

$$0.00796 < \sigma_e^2 < 0.0209$$

(7)

4 (a) From the def. contrasts AB, BC, CD we derive the following interactions are also confounded.

$$AB \times BC \pmod{2} = AC$$

$$AB \times CD \pmod{2} = ABCD$$

$$BC \times CD \pmod{2} = BD$$

$$AB \times BC \times CD \pmod{2} = AD$$

For all 3 replicates - the 16 treat. comb's are assigned into 8 blocks as follows.

$$(1) \quad \begin{array}{cccccccc} a & b & ab & c & ac & bc & abc & \\ d & ad & bd & abd & cd & acd & bcd & abcd \end{array}$$

↓ By AB

$$\left\{ \begin{array}{l} \text{Gp 1: } (1) \quad a b \quad c \quad abc \quad d \quad abd \quad cd \quad abcd \\ \text{Gp 2: } \quad a \quad b \quad ac \quad bc \quad ad \quad bd \quad acd \quad bcd \end{array} \right.$$

↓ By BC

$$\left\{ \begin{array}{l} \text{Gp 1: } (1) \quad abc \quad d \quad abcd \\ \text{Gp 2: } \quad ab \quad c \quad abd \quad cd \\ \text{Gp 3: } \quad a \quad bc \quad ad \quad bcd \\ \text{Gp 4: } \quad b \quad ac \quad bd \quad acd \end{array} \right.$$

↓ By CD

$$\text{Block 1: } (1) \quad abcd$$

$$2: \quad abc \quad d$$

$$3: \quad ab \quad cd$$

$$4: \quad c \quad abd$$

$$5: \quad a \quad bcd$$

$$6: \quad bc \quad ad$$

$$7: \quad b \quad acd$$

$$8: \quad ac \quad bd$$

ANOVA

Source	d.f.
A	1
B	1
C	1
D	1
ABC	1
ABD	1
ACD	1
BCD	1
Blocks	7
Replicates	2
Blocks x Repl	14
Error	16
Total	47

(12)

4① The $\frac{1}{2}$ -replicate is

- (i) b c bc ad abd acd abcd

A x (the $\frac{1}{2}$ -repl.) (mod 2) yields

a ab ac abc d bd ed bed
 = the other $\frac{1}{2}$ -repl.

so A \in the def. contrast

B x (the $\frac{1}{2}$ -repl.) (mod 2) yields

b (i) bc b abd ad abcd acd
 = the same $\frac{1}{2}$ -repl.

so B \notin the def. contrast

C x (the $\frac{1}{2}$ -repl.) (mod 2) yields

c bc cd b acd abcd ad abd
= the same $\frac{1}{2}$ -repl.

So C \notin the def. contrast

D x (the $\frac{1}{2}$ -repl) (mod. 2) yields

d bd cd bcd a ab ac abc
= the other $\frac{1}{2}$ -repl.

So D \in the defining contrast

Hence the def. contrast is AD

Factor/interaction	Alias
A	D
B	ABD
C	ACD
AB	BD
AC	CD
BC	ABCD
ABC	BCD

9

4C ANOVA

Source	d.f.
Repl.	2
A or D	1
B or ABD	1
C or ACD	1
AB or BD	1
AC or CD	1
BC or ABCD	1
ABC or BCD	1
Error	14
Total	23

4