

Q1 10
Q2 15

Q3 10 Q5 15
Q4 10 Q6 30

2101 // 80

01	6	
23	6	22
45	6	455
67	6	66
89	6	888
01	7	0100001
23	7	3223333
45	7	454
67	7	66
89	7	
	8	

$\bar{x} = 69.8$

442 // 45

2	
2	8755
3	
3	9885
4	4210110
4	8

$\bar{x} = 36.9$

Ex J - handout 2 typos so far: x_3 is in grams
not mm.
equation (J2) not labelled

Designed experiments / factorial treatment structure

adv. - because we set x 's, we can be more sure that
they "cause" the observed change in y .

(see § 9.1.2)

- we can balance the exp't by observing well
chosen values for x 's

§ 9.1

on randomization

e.g. Ex J 3 levels of x_1, x_2, x_3
equal # of runs at ea. = 1
combination.

variables in data set are orthogonal,
so interpretation of $\hat{\beta}$ is clearer

i.e. $\hat{\beta}_1$ same in $y_i = \beta_0 + \beta_1 x_{1i} + \varepsilon_i$

or $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \varepsilon_i$

- we can randomize the assignment
of treatments to experimental units
(balanced on unknowns, on average)
(we can sometimes 'blind' subjects,
so other confounders are also balanced)

Treatments are very often structured in a factorial design

"x" 's are set at different levels

factors can be quantitative or qualitative

Ex J 3 quant. factors each at 3 levels
in HO $x \rightarrow \text{len, amp, load}$
 $\log() \log() \log()$

HW 2 : factorial analysis regression 1m
using -1, 0, 1 etc.

Table 9.3 teaching methods 1 qual. factor A B C D E levels

Seedrate HO
1 quant. factor 5 levels

Table 8.2 3 factors seat height 2 levels } quant.
tire pressure 2 levels }
dynamo on/off qual

Example 9.6 poisons (Table 8.10)
2 factors "treatment" 4 levels
"poison" 3 levels

For all of these, basic starting point is a linear model.

Ex J $y_i = \beta_0 + \beta_1 x_{1i} + \dots + \epsilon_i$

SR HO $y_{tr} = \alpha + \gamma_t + \epsilon_{tr}$

linear structure for $E(y)$; additive errors $\sim (0, \sigma^2)$

parametrization changes

even w Ex G $y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_{10} x_{10i} + \epsilon_i$

models are all $y = X\beta + \epsilon$

fitting is always LS

$$\min_{\beta} \sum_i (y_i - E(y_i))^2 = SS(\beta)$$

R is using lm in all cases.

make an $n \times 1$ vector of y 's

$n \times p$ matrix X — model.matrix

(even if you ask for aov)

When factors are qualitative, usually summarize the analysis with

- 1) anova table, F-tests (composite H_0)
e.g. $\gamma_1 = \dots = \gamma_T = 0$
- 2) group means, e.g. \bar{y}_t + est'd std. errors
- 3) compare various group means, as relevant for the problem

Ex Teaching methods A B C D E

p. 428 1) $\bar{y}_A - \bar{y}_B \pm 1.96 \cdot \widehat{S.E.}(\bar{y}_A - \bar{y}_B)$
 $\searrow \quad \swarrow$
 (uncorrelated) This suggests

2) $\bar{y}_u = \frac{1}{2}(\bar{y}_A + \bar{y}_B)$

$\bar{y}_C - \bar{y}_u$: CI for $\mu_C - \frac{1}{2}(\mu_A + \mu_B)$ is consistent w data
 $\bar{y}_D - \bar{y}_u$:
 $\bar{y}_E - \bar{y}_u$:

Ex Seedrate (H_0) \leftarrow 1 factor quantitative

on p. 1. : fit $y_{tj} = \alpha + \gamma_t + \varepsilon_{tj}$

under 2 constraints 1) $\gamma_1 = 0$

2) $\sum_{t=1}^T \gamma_t = 0$

1) is the default param = for factors that aren't ordered

> barley \$ SeedRate = factor(...)

most usual in software packages \rightarrow > options ("contrasts") # checking what's in place
 $\$contrast$
 unordered ordered
 "contr.treatment" "contr.poly" $\gamma_1 = 0$

*) \rightarrow > options (contrast = c("contr.sum", "contr.poly")) # $\sum_{t=1}^T \gamma_t = 0$
 ☺ \rightarrow > model.matrix(my.lm)
 \rightarrow refit the same model $\hat{\gamma}$'s change

→ refit the same model $\hat{\gamma}$'s change
anova table is the same.

$$\begin{array}{ccc} \text{coef}(\dots \text{lm}) & \begin{array}{c} \hat{\alpha} \\ \hat{\gamma}_1 \\ \hat{\gamma}_2 \\ \hat{\gamma}_3 \\ \hat{\gamma}_4 \\ \hat{\gamma}_5 \end{array} & \begin{array}{c} \hat{\text{se}}(\hat{\alpha}) \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{array} \\ & & \hat{\gamma}_5 = -(\hat{\gamma}_1 + \dots + \hat{\gamma}_4) \end{array}$$

$H_0: \gamma_2 = 0$ doesn't make sense

$H_0: \gamma_2 = \gamma_1$ does make sense

"make sense" - depends on the problem

- another lm, insisting on treating suchdate as quantitative
now we have a 3rd possible param =
use orthogonal polynomials (yet another)

$$\begin{array}{lll} E(y_{tj}): \alpha + \gamma_t & \alpha + \gamma_t & \alpha + \underbrace{\beta_1 P_1(x)}_{\text{take 1 val}} + \beta_2 P_2(x) + \beta_3 P_3(x) + \beta_4 P_4(x) \\ \gamma_1 = 0 & \sum \gamma_t = 0 & -2, -1, 0, 1, 2 \\ (\alpha, \gamma_2, \dots, \gamma_5) & (\alpha, \gamma_1, \dots, \gamma_4) & (\alpha, \beta_1, \dots, \beta_4) \\ X & X' & X'' \end{array}$$

> SR = as.ordered(SR)

> contrasts(SR)

	.L	.Q	.C	.QQ
1	-2	-.63	2	
2	-1	-.31	-1	
3	0	0	-2	
4	1	-.31	-1	
5	2	.63	2	

If factors are quant., then the comparison of interest are usually related to regression.

That's why R was "contr. poly", because it gives coefficients relevant to

- linear relationship betw y & x
- quadratic " " "
- cubic " " "
- quartic " " "

very often of interest

See also Ex 9.12

> coef.lm or > summary (my.lm)

> anova (... split = list (...)) see last pg of H0

Sometimes, we have both, meaning

- one factor (or more)

HW #1

$z = \pm 1$

- quantitative variable (or more)

x cont^s

$$y_{tj} = \alpha + \gamma_t + \beta x_{tj} + \varepsilon_{tj} \quad \text{is one model}$$

(*) \rightarrow
§ analysis of
covariance
Ex 9.13

$$\text{OR } \boxed{\mu + \gamma_t + \beta(x_{tj} - \bar{x})} + \varepsilon_{tj} \quad \text{Ex 11}$$

OR (if x was set at 5 equally spaced levels)

$$y_{tj} = \beta_0 + \gamma_t + \beta_1 P_1(x) + \dots + \beta_4 P_4(x) + \varepsilon_{tj}$$

using orthog poly's

$$GR = \dots =$$

§ 9.3.2 for ~~an~~ a linear algebra explanation of contrasts

Example 9.6 has 2 factors . 4 obs^s at each condⁿ

Example 9.6 has 2 factors, 4 obsⁿ at each combⁿ
of factors §9.2.4 has the analysis

Example H

$$y_i = \begin{cases} 1 & \text{if there is a fault in sample } i \\ 0 & \text{if not} \end{cases}$$

$$x_i = \text{purity index for sample } i$$

$$z_i = \begin{cases} -1 & \text{if method = standard} \\ +1 & \text{if method = modified} \end{cases}$$

$$i = 1, \dots, 44$$

22 batches, divided in half

44 samples

- 1st analysis ignores x_i and uses pairing (paired t-test \rightarrow binary response)
- regression analysis ignores pairing & pretends 44 ind't obsⁿ

Not this model

$$y_i = \alpha + \beta_1 z_i + \beta_2 x_i + \varepsilon_i$$

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

$$P_1(y_i = 1) = p_i = \frac{e^{\alpha + \beta_1 z_i + \beta_2 x_i}}{1 + e^{\alpha + \beta_1 z_i + \beta_2 x_i}} = E y_i$$

random part

$$y_i = \begin{cases} 1 & \text{w prob } p_i \\ 0 & \text{w prob } 1 - p_i \end{cases} \quad \text{binomial dist.}$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \alpha + \beta_1 z_i + \beta_2 x_i \quad \text{linear on log-odds}$$

log-odds ratio (works for interpretation)

slight diff.

$$\log\left(\frac{p_i}{1-p_i}\right) = \begin{cases} \alpha + \Delta + \beta(x_i - \bar{x}) & \text{std.} \\ \alpha - \Delta + \beta(x_i - \bar{x}) & \text{mod.} \end{cases}$$

$$z_i = \begin{cases} +1 & \text{std} \\ -1 & \text{mod} \end{cases} \quad \Delta \text{ instead of } \beta_1$$