

HW9

Sam Olson

Q1

A plant scientist was interested in comparing two plant genotypes (1 and 2). An experiment was conducted in a greenhouse with one table, eight trays, and sixteen pots. The table in the greenhouse held the eight trays with two pots on each tray. For each of the eight trays, two genotype 1 seeds were planted in one pot, and two genotype 2 seeds were planted in the other pot. The assignment of genotypes 1 and 2 to the two pots within each tray was determined by flipping a fair coin. The response of interest is a quantitative measurement of overall plant health that was calculated for each plant 42 days after planting. These quantitative measurements of overall plant health are presented as integers in Table 1 to make calculations easier, but please answer all questions as if each measurement is a realization from a normal distribution.

Table 1. Measurements of overall plant health for each plant.

Tray	Genotype 1 Pot		Genotype 2 Pot	
	Plant 1	Plant 2	Plant 1	Plant 2
1	8	7	6	7
2	8	9	4	5
3	8	8	7	7
4	5	7	4	2
5	5	6	4	3
6	9	10	7	9
7	5	7	1	4
8	4	6	5	5

Let i index genotypes ($i = 1, 2$), j index trays ($j = 1, \dots, 8$), and k index plants within pots ($k = 1, 2$). Let y_{ijk} denote the response corresponding to genotype i , tray j , and plant k . Suppose:

$$y_{ijk} = \mu_i + t_j + p_{ij} + e_{ijk} \quad \forall i, j, k,$$

where μ_1 and μ_2 are unknown real-valued parameters, $t_j \sim \mathcal{N}(0, \sigma_t^2) \quad \forall j$, $p_{ij} \sim \mathcal{N}(0, \sigma_p^2) \quad \forall i, j$, $e_{ijk} \sim \mathcal{N}(0, \sigma_e^2) \quad \forall i, j, k$, and all t_j , p_{ij} , e_{ijk} terms are mutually independent.

a)

Explain what the p_{ij} terms represent and provide one reason for including them in the model.

Answer

The p_{ij} terms are random effects corresponding to pots.

Pots are the experimental units in this experiment because the levels of the factor genotype were randomly assigned to pots. Because there are multiple observations per experimental unit, it is important to explicitly include a random effect for each experimental unit in the model.

These random effects account for pot-to-pot variation in the response that may occur due to differences among pots. The pot random effects allow for the correlation between the two observations from a given pot to be larger than the correlation between two observations from different pots.

b)

Let $\bar{y}_{ij} = \frac{1}{2} \sum_{k=1}^2 y_{ijk} \quad \forall i, j$. Determine the distribution of $\bar{y}_{11} - \bar{y}_{21}$.

Answer

Linear combinations of normal random variables are normal, so $\bar{y}_{11} - \bar{y}_{21}$ will be normally distributed. We then need only find the mean and variance to uniquely characterize the distribution.

To that end, note:

$$\bar{y}_{11} - \bar{y}_{21} = \mu_1 - \mu_2 + p_{11} - p_{21} + \bar{e}_{11} - \bar{e}_{21}.$$

Thus,

$$E(\bar{y}_{11} - \bar{y}_{21}) = \mu_1 - \mu_2$$

$$\text{Var}(\bar{y}_{11} - \bar{y}_{21}) = \sigma_p^2 + \sigma_p^2 + \sigma_e^2/2 + \sigma_e^2/2 = 2\sigma_p^2 + \sigma_e^2.$$

$$\bar{y}_{11} - \bar{y}_{21} \sim N(\mu_1 - \mu_2, 2\sigma_p^2 + \sigma_e^2)$$

c)

Compute the value of an unbiased estimator of the variance of $\bar{y}_{11} - \bar{y}_{21}$.

Answer

Let $d_j = \bar{y}_{1j} - \bar{y}_{2j} \cdot \forall j$.

From part b), it follows:

$$d_1, \dots, d_8 \stackrel{iid}{\sim} N(\mu_1 - \mu_2, 2\sigma_p^2 + \sigma_e^2)$$

Thus,

$$s_d^2 = \sum_{j=1}^8 (d_j - \bar{d})^2 / (8 - 1)$$

is an unbiased estimator of $2\sigma_p^2 + \sigma_e^2$.

Calculating explicitly:

```
plant_data_wide <- data.frame(
  Tray = 1:8,
  G1_P1 = c(8, 8, 8, 5, 5, 9, 5, 4),
  G1_P2 = c(7, 9, 8, 7, 6, 10, 7, 6),
  G2_P1 = c(6, 4, 7, 4, 4, 7, 1, 5),
  G2_P2 = c(7, 5, 7, 2, 3, 9, 4, 5)
)

plant_data_wide$d_j <- rowMeans(plant_data_wide[, c("G1_P1", "G1_P2")]) -
  rowMeans(plant_data_wide[, c("G2_P1", "G2_P2")])

d_bar <- mean(plant_data_wide$d_j)

s_d_squared <- sum((plant_data_wide$d_j - d_bar)^2) / (length(plant_data_wide$d_j) - 1)
s_d_squared

## [1] 1.928571
```

$$s_d^2 = 13.5/7 \approx 1.929$$

d)

Provide a 95% confidence interval for $\mu_1 - \mu_2$.

Answer

```
est <- mean(plant_data_wide$d_j)

moe <- qt(p = 0.975, df = 7)*(sqrt(s_d_squared/8))

lb <- est - moe
ub <- est + moe

cat("Lower bound:", lb, "\nUpper bound:", ub, "\n")

## Lower bound: 0.8389925
## Upper bound: 3.161007
```

$$\widehat{\mu_1 - \mu_2} \pm t_{0.975,7} \sqrt{s_d^2/8} \rightarrow 2 \pm (2.36 \cdot 0.49) \rightarrow (0.84, 3.16)$$

e)

The model can be written in the form $y = X\beta + Zu + e$. Provide X , β , Z , and u .

Answer

If we order the y vector by working our way across trays with the genotype 1 pot first and the genotype 2 pot second within each tray, we have:

$$X = \mathbf{1}_{8 \times 1} \otimes \mathbf{I}_{2 \times 2} \otimes \mathbf{1}_{2 \times 1}$$

$$\beta = (\mu_1, \mu_2)^\top$$

$$Z = [\mathbf{I}_{8 \times 8} \otimes \mathbf{1}_{4 \times 1}, \mathbf{I}_{16 \times 16} \otimes \mathbf{1}_{2 \times 1}]$$

and

$$u = [t_1, t_2, \dots, t_8, p_{11}, p_{21}, p_{12}, p_{22}, \dots, p_{18}, p_{28}]^\top$$

f)

Suppose the researchers would like to repeat their experiment again, using the same basic resources: eight trays, two pots per tray, sixteen seeds of genotype 1, and sixteen seeds of genotype 2. Would you recommend any changes to their experimental design? Explain why or why not.

Answer

It would be better to put one plant of genotype 1 and one plant of genotype 2 in each of the 16 pots. Then the variance of the genotype 1 average minus the genotype 2 average would be $\sigma_e^2/8$ because both tray and pot random effects would cancel in the difference of averages. This variance ($\sigma_e^2/8$) is less than the variance for the estimated difference in genotype means for the original design: $(2\sigma_p^2 + \sigma_e^2)/8 = \sigma_p^2/4 + \sigma_e^2/8$. Thus, we have a more precise unbiased estimator of the difference in genotype means if we put one plant of each genotype in each pot.

Q2

This is a continuation of Problem 1. Suppose the experiment actually involved a second factor—bacterial infection with levels 1=present and 2=absent—in addition to the factor genotype, randomly assigned to pots as discussed previously. Within each pot, one of the two plants was randomly selected for infection with a bacteria, which was applied by rubbing a gel containing the bacteria on the top leaf of the plant. The other plant in each pot was rubbed with the same gel but with the bacteria absent.

Let y_{ijk} be the response for the plant of genotype i on tray j that received bacterial infection k ($i = 1, 2; j = 1, \dots, 8; k = 1, 2$). Suppose the data are the same as in Table 1 and arranged so that in each pot, Plant 1 corresponds to the plant infected with the bacteria and Plant 2 corresponds to the plant not infected with the bacteria. Suppose:

$$y_{ijk} = \mu_{ik} + t_j + p_{ij} + e_{ijk} \quad \forall i, j, k,$$

where, as in model (1), $t_j \sim \mathcal{N}(0, \sigma_t^2) \forall j$, $p_{ij} \sim \mathcal{N}(0, \sigma_p^2) \forall i, j$, $e_{ijk} \sim \mathcal{N}(0, \sigma_e^2) \forall i, j, k$, and all t_j , p_{ij} , and e_{ijk} terms are mutually independent.

a)

This is a split-plot experiment. What are the whole-plot experimental units?

Answer

Pots

b)

What are the split-plot experimental units?

Answer

Plants

c)

What is the whole-plot treatment factor?

Answer

Genotype

d)

What is the split-plot treatment factor?

Answer

Bacterial infection

e)

Create an ANOVA table with columns Source and Degrees of Freedom.

Answer

Source	DF	SS
Trays	7	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{.j.} - \bar{y}_{...})^2$
Genotypes	1	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{i..} - \bar{y}_{...})^2$
Trays \times Genotypes	7	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2$
Infections	1	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{..k} - \bar{y}_{...})^2$
Genotypes \times Infections	1	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{i.k} - \bar{y}_{i..} - \bar{y}_{..k} + \bar{y}_{...})^2$
Error	14	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (y_{ijk} - \bar{y}_{ij.} - \bar{y}_{i.k} + \bar{y}_{i..})^2$
Corrected Total	31	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (y_{ijk} - \bar{y}_{...})^2$

f)

Give formulas for each of the Sums of Squares of the ANOVA table. (Shortcut formulas for degrees of freedom and sums of squares work in this case because of the balanced experimental design.)

Answer

Source	SS
Trays	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{.j.} - \bar{y}_{...})^2$
Genotypes	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{i..} - \bar{y}_{...})^2$
Trays \times Genotypes	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2$
Infections	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{..k} - \bar{y}_{...})^2$
Genotypes \times Infections	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (\bar{y}_{i.k} - \bar{y}_{i..} - \bar{y}_{..k} + \bar{y}_{...})^2$
Error	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (y_{ijk} - \bar{y}_{ij.} - \bar{y}_{i.k} + \bar{y}_{i..})^2$
Corrected Total	$\sum_{i=1}^2 \sum_{j=1}^8 \sum_{k=1}^2 (y_{ijk} - \bar{y}_{...})^2$

g)

Derive the expected mean square for the second to last line of the ANOVA table (the line right before corrected total). This line is typically called error or split-plot error.

Answer

$$E(MS_{error}) = \sigma_e^2$$

h)

Compute the value of the best linear unbiased estimator of $\mu_{11} - \mu_{12}$.

Answer

$$\bar{y}_{1.1} - \bar{y}_{1.2} = -1$$

i)

Derive an expression for the variance of the best linear unbiased estimator of $\mu_{11} - \mu_{12}$ in terms of model (2) parameters.

Answer

$$\text{Var}(\widehat{\mu_{11} - \mu_{12}}) = \frac{1}{4}\sigma_e^2$$

j)

Compute a 95% confidence interval for $\mu_{11} - \mu_{12}$.

HERE OFFICER

Answer

```
# est <- -plant_data_wide$d_j[1]
est <- plant_data_wide$G1_P1 - plant_data_wide$G1_P2

moe <- qt(p = 0.975, df = 14)*(sqrt(s_d_squared/8))

lb <- est - moe
ub <- est + moe

cat("Lower bound:", lb, "\nUpper bound:", ub, "\n")
```

```
## Lower bound: -0.05306936 -2.053069 -1.053069 -3.053069 -2.053069 -2.053069 -3.053069 -3.053069
## Upper bound: 2.053069 0.05306936 1.053069 -0.9469306 0.05306936 0.05306936 -0.9469306 -0.9469306
```

$$\widehat{\mu_{11} - \mu_{12}} \pm t_{0.975, 14} \cdot \left(\frac{1}{4}\right) \sigma_e^2 \rightarrow -1 \pm (2.1448 \cdot 0.4818) = (-2.0334, 0.0334)$$

k)

Determine the distribution of $y_{111} - y_{112} - y_{211} + y_{212}$.

Answer

$$N(\mu_{11} - \mu_{12} - \mu_{21} + \mu_{22}, 4\sigma_e^2)$$

l)

Compute the value of the best linear unbiased estimator of $\mu_{11} - \mu_{12} - \mu_{21} + \mu_{22}$.

Answer

$$\bar{y}_{1.1} - \bar{y}_{1.2} - \bar{y}_{2.1} + \bar{y}_{2.2} = -0.5.$$

m)

Derive an expression for the variance of the best linear unbiased estimator of $\mu_{11} - \mu_{12} - \mu_{21} + \mu_{22}$ in terms of model (2) parameters.

Answer

$$\text{Var}(\bar{y}_{1.1} - \bar{y}_{1.2} - \bar{y}_{2.1} + \bar{y}_{2.2}) = \frac{1}{2}\sigma_e^2$$

Q3

Researchers created a device to test the effectiveness of helmets at reducing the stress caused by head impacts. The device includes a head-shaped sensor on which a helmet can be placed, as well as a striking weight that can produce impacts to the front or side of a helmet placed on the sensor. The intensity of each impact can be controlled by the researchers. When an impact is delivered, a measurement of the amount of stress experienced by the head-shaped sensor is recorded. A measurement of 0 indicates no stress, while a measurement of 100 indicates stress high enough to cause serious brain injury.

The researchers used the device to test a total of 10 helmets consisting of 5 helmets of type 1 and 5 helmets of type 2. The 10 helmets were tested in random order. When each helmet was tested, it was struck a total of 4 times: once with low impact to the front, once with high impact to the front, once with low impact to the side, and once with high impact to the side. The order of the 4 impacts was determined separately for each helmet using the following procedure. A fair coin was flipped. If the result of the flip was heads, the first two impacts were front impacts and the last two impacts were side impacts. If the result of the flip was tails, the first two impacts were side impacts and the last two impacts were front impacts. For the first two impacts, the coin was flipped again. If the result of the flip was heads, the first impact was at low intensity and the second was at high intensity. If the result of the flip was tails, the first impact was at high intensity and the second at low intensity. A coin was flipped a third time to determine the order of the impact intensities for the third and fourth impacts so that each order (low and then high vs. high and then low) was equally likely.

Let $i = 1, 2$ index helmet types 1 and 2. Let $j = 1, \dots, 5$ index helmets nested within helmet types. Let $k = 1, 2$ index the direction of impact, with $k = 1$ for front and $k = 2$ for side. Let $\ell = 1, 2$ index the intensity of the impact, with $\ell = 1$ for low and $\ell = 2$ for high. Let y_{ijkl} be the stress measurement for the corresponding values of i, j, k , and ℓ . For $i = 1, 2, j = 1, \dots, 5, k = 1, 2$, and $\ell = 1, 2$, consider the model:

$$y_{ijkl} = \mu_{ik\ell} + a_{ij} + b_{ijk} + e_{ijkl},$$

where the $\mu_{ik\ell}$ values are unknown parameters, $a_{ij} \sim \mathcal{N}(0, \sigma_a^2)$, $b_{ijk} \sim \mathcal{N}(0, \sigma_b^2)$, $e_{ijkl} \sim \mathcal{N}(0, \sigma_e^2)$, and all random terms are independent. Model (3) was fit to the dataset, and the following ANOVA table was obtained. Because we have a balanced experimental design, the type I and type III sums of squares are the same, and the lines of the ANOVA table can be reordered in a variety of ways without changing the results.

Source	Sum of Squares
Type	226
Direction	255
Intensity	8910
Type \times Direction	207
Type \times Intensity	2
Direction \times Intensity	7
Type \times Direction \times Intensity	9
Helmet(Type)	254
Direction \times Helmet(Type)	114
Error	59
C. Total	10043

a)

We learned a shortcut for expressing sums of squares in summation notation that works for balanced designs like the one considered here. Use that shortcut to express the sum of squares for Direction \times Intensity using summation notation.

Answer

$$(k-1)(l-1) = kl - k - l + 1 \sum_{i=1}^2 \sum_{j=1}^5 \sum_{k=1}^2 \sum_{l=1}^2 (\bar{y}_{-kl} - \bar{y}_{-k} - \bar{y}_{-l} + \bar{y}_{-m})^2 = 10 \sum_{k=1}^2 \sum_{l=1}^2 (\bar{y}_{-kl} - \bar{y}_{-k} - \bar{y}_{-l} + \bar{y}_{-m})^2$$

b)

Compute a t statistic that can be used to test $H_0 : \bar{\mu}_{1...} = \bar{\mu}_{2...}$.

Answer

$$t = \sqrt{\frac{MS_{Type}}{MS_{Helmet(Type)}}} = \sqrt{\frac{226/1}{254/8}}$$

c)

The statistic in part b) has a noncentral t distribution. Provide an expression for the noncentrality parameter in terms of model (3) parameters.

Answer

$$\frac{\bar{\mu}_{1...} - \bar{\mu}_{2...}}{\sqrt{(4\sigma_a^2 + 2\sigma_b^2 + \sigma_e^2)/10}}$$

d)

Compute the value of an unbiased estimator for σ_a^2 .

Answer

$$\frac{MS_{Helmet(Type)} - MS_{Direction \times Helmet(Type)}}{4} = \frac{254/8 - 114/8}{4} = \frac{140}{32} = 4.375$$

e)

The best linear unbiased estimator of $\bar{\mu}_{12...} - \bar{\mu}_{11...}$ is equal to 0.5 for this dataset. Provide a 95% confidence interval for $\bar{\mu}_{12...} - \bar{\mu}_{11...}$.

ALSO HERE

Answer

$$0.5 \pm 2.306\sqrt{2.85}$$

where $t_{0.975,8} = 2.306$ and the degrees of freedom (8) come from $Direction \times Helmet(Type)$.

f)

Compute a standard error for the best linear unbiased estimator of $\mu_{121} - \mu_{111}$.

Answer

The standard error for the BLUE of $\mu_{121} - \mu_{111}$ is:

$$\widehat{Var}(\bar{y}_{1.21} - \bar{y}_{1.11}) = \frac{1}{5} \left(\frac{114}{8} + \frac{59}{16} \right) = 3.5875$$

Giving us:

$$SE(\bar{y}_{1.21} - \bar{y}_{1.11}) = \sqrt{3.5875} = 1.89407$$