

- Another option: Force the last $\tau_i = 0$. How SAS does it.

- These options give different numerical estimates for the parameters, but all conclusions about effects and contrasts will be the same.

Unbalanced Data

- Using the standard ANOVA formulas is easy, but it will give wrong results when data are unbalanced (different numbers of observations across cells).
- Dummy variable approach always gives correct answers.

Illustration: A unbalanced 2-factor factorial study. (Table 11.3 data, p. 585)

		C	
		1	2
A	1	4,5,6	8
	2	5	7,9

- Question: Does factor A have a significant effect on the response? (For simplicity, ignore any interaction between A and C for this example).

Recall: Our F-statistic formula for this type of test was:

$$F^* = \frac{MSA}{MSW} = \frac{SSA / (a-1)}{SSW / ac(n-1)}$$

and $SSA = cn \sum_i (\bar{Y}_{i..} - \bar{Y}_{...})^2$

- This formula is based on the variation between the marginal means $\bar{Y}_{1..}$ and $\bar{Y}_{2..}$

- For the Table 11.3 data:

$$\bar{Y}_{1..} = \frac{4+5+6+8}{4} = 5.75$$

$$\bar{Y}_{2..} = \frac{5+7+9}{3} = 7$$

→ Based on this, there is some sample variation between the means for levels 1 and 2 of factor A.

- However, let's look at the sample means for levels 1 and 2 of A, separately at each level of C:

For level 1 of C:

$$\bar{Y}_{11.} = \frac{4+5+6}{3} = 5$$

$$\bar{Y}_{21.} = 5$$

For level 2 of C:

$$\bar{Y}_{12.} = 8$$

$$\bar{Y}_{22.} = \frac{7+9}{2} = 8$$

- These results imply that (at each level of C) there is no sample variation between the means for levels 1 and 2 of factor A.

• Which conclusion is correct?

• Our model is (recall there is no interaction term):

$$Y_{ijk} = \mu + \alpha_i + \gamma_j + \varepsilon_{ijk}$$

Note: $\bar{Y}_{11\cdot} - \bar{Y}_{21\cdot}$ is an estimate of: $E(Y_{11k}) - E(Y_{21k}) = (\mu + \alpha_1 + \gamma_1) - (\mu + \alpha_2 + \gamma_1) = \alpha_1 - \alpha_2$

Also, $\bar{Y}_{12\cdot} - \bar{Y}_{22\cdot}$ is an estimate of: $E(Y_{12k}) - E(Y_{22k}) = (\mu + \alpha_1 + \gamma_2) - (\mu + \alpha_2 + \gamma_2) = \alpha_1 - \alpha_2$

• So these do estimate the true difference in the means for levels 1 and 2 of factor A.

But ... $\bar{Y}_{1\cdot\cdot} - \bar{Y}_{2\cdot\cdot}$, for these data, is:

$$\frac{3\bar{Y}_{11\cdot} + \bar{Y}_{12\cdot}}{4} - \frac{\bar{Y}_{21\cdot} + 2\bar{Y}_{22\cdot}}{3}$$

which estimates:

$$\frac{3}{4}(\mu + \alpha_1 + \gamma_1) + \frac{1}{4}(\mu + \alpha_1 + \gamma_2) - \frac{1}{3}(\mu + \alpha_2 + \gamma_1) - \frac{2}{3}(\mu + \alpha_2 + \gamma_2)$$

$$= \alpha_1 - \alpha_2 + \frac{5}{12}\gamma_1 - \frac{5}{12}\gamma_2$$

- This is not the true difference in factor A's level means that we wanted to estimate.
- For balanced data, the magnitudes of all the coefficients would be the same and everything would cancel out properly.
- With unbalanced data, we need to adjust for the fact that the various cell means are based on different numbers of observations per cell.
- Using a dummy variable regression model implies the effect of factor A is estimated holding factor C constant → produces correct results.
- Analysis for unbalanced data involves the least squares means, not the ordinary factor level means.

LS mean
for level 1
of A is:
 $\frac{5+8}{2} = 6.5$

LS mean
for level 2
of A is:
 $\frac{5+8}{2} = 6.5$

- The least squares mean (for, say, level 1 of factor A) is the unweighted average of the cell sample means corresponding to level 1 of factor A. With unbalanced data, this is different than simply averaging all response values for level 1 of factor A. (see example)
- With unbalanced data in the two-way ANOVA, our F-tests about the factors use the Type III sums of squares, rather than the ordinary (Type I) ANOVA SS.
- See example for calculating these F-statistics correctly.

Example: (Table 11.3 data)

• **Least squares means:**

for Level 1 of Factor A: $\frac{5+8}{2} = 6.5$

for Level 2 of Factor A: $\frac{5+8}{2} = 6.5$

for Level 1 of Factor C: $\frac{5+5}{2} = 5$

for Level 2 of Factor C: $\frac{8+8}{2} = 8$

• **Correct F-tests about factor effects:**

- Done based on Type III SS which uses a dummy variable regression model and full vs. reduced model F-tests.

- No significant A x C interaction.

- Factor A main-effects F-test had a P-value of 1.

- Factor C main-effect F-test had a P-value of .054

• **More complicated example: Suppose A has 3 levels and C has 2 levels.**

based on Type III SS.

• **Now we need to use 3 - 1 = 2 dummy variables for A and 2 - 1 = 1 dummy variable for C.**

Example:

		Factor C	
		1	2
Factor A	1	4,5,6	8
	2	5	7,9
	3	11	12

Again, Type III SS does the dummy variable regression approach;

A x C interaction was not significant. Factor A main-effects F-test had P-value .0305.

Factor C main-effects F-test had P-value .0703.