

## Unbalanced data in factorial designs

Stat 346  
5-27-25

①

		Temp (°F)		
		15	70	125
Material	Type	1	34 40	70 58
		2	159 26	136 115
3			138 160	150 139 96

②

Look at the  
cell counts

4	4	2	10
2	2	1	5
2	2	1	5
8	8	4	20

Notice that the rows are proportional

and the columns are proportional

$$\text{and } n_{ij} = \frac{n_i \cdot n_{\cdot j}}{n_{\cdot \cdot}}$$

This is called Proportional data

(3)

You may be able to run this using the balanced ANOVA module in the stat package

Exact Analysis

$$\text{Model: } Y_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + \varepsilon_{(ijk)}$$

$$i = 1, 2, 3 \quad j = 1, 2, 3 \quad k = 1, \dots, n_{ij}$$

$$\sum_i \tau_i = 0 \quad \sum_j \beta_j = 0 \quad \sum_i \tau \beta_{ij} = \sum_j \tau \beta_{ij} = 0$$

(4)

1W

$$\begin{bmatrix} \mu & \tau_1 & \tau_2 & \beta_1 & \beta_2 & \tau\beta_{11} & \tau\beta_{12} & \tau\beta_{21} & \tau\beta_{22} \\ & Qx_1 & & & & & & & \end{bmatrix}$$

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$x_{15}$	$x_{16}$	$x_{17}$	$x_{18}$	$x_{19}$	$x_{20}$
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

11

$y$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$x_{15}$	$x_{16}$	$x_{17}$	$x_{18}$	$x_{19}$	$x_{20}$
130	155	74	180	34	40	28	15	70	58	59	126	136	115	45	38	160	50	139	96	

(5)

$$\text{Note: } \hat{\beta}_{13} = -\hat{\beta}_{11} - \hat{\beta}_{12}$$

$$\hat{\beta}_{23} = -\hat{\beta}_{21} - \hat{\beta}_{22}$$

$$\hat{\beta}_{31} = -\hat{\beta}_{11} - \hat{\beta}_{21}$$

$$\hat{\beta}_{32} = -\hat{\beta}_{12} - \hat{\beta}_{22}$$

$$\hat{\beta}_{33} = -\hat{\beta}_{31} - \hat{\beta}_{32}$$

$$= \hat{\beta}_{11} + \hat{\beta}_{21} + \hat{\beta}_{12} + \hat{\beta}_{22}$$

(6)

To get an exact analysis, run this as a regression, entering the variables  $Y, X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8$  as your data set.

Look at the sequential sums of squares.

$$SS_A = SS_{X_1} + SS_{X_2|X_1}$$

$$SS_{B|A} = SS_{X_3|X_1, X_2} + SS_{X_4|X_1, X_2, X_3}$$

$$SS_{AB|A,B} = SS_{X_5|X_1, X_2, X_3, X_4} + \dots$$

(7)

### Approximate methods:

Imputation: In Stat 525, we considered the case of an empty cell.

If a cell has some values, but fewer than the other cells, you can use  $\bar{y}_{ij}$  as the imputed value.

You must manually subtract an error of for each imputed value.

(8)

### Setting aside a data value:

Randomly set aside values from the over-populated cells.

Yates' Method : Compute the cell means + treat those as if they were the data set.

(9)

	Temp			
	15	70	125	
Type	1	134.75	57.23	64
	2	142.5	125.5	45
	3	149	144.5	96

Run this as a 2-way ANOVA with interaction. There will 0 df for error.

Compute  $\hat{\sigma}^2 = \frac{\sum \sum \sum (y_{ijk} - \bar{y}_{ij.})^2}{n_{..} - ab}$

This estimates  $\sigma^2 = V(y_{ijk})$  (10)

However, Yates' Method is an ANOVA on the cell means, i.e. the  $\bar{y}_{ij.}$ , not the original  $y_{ijk}$ .

$$V(\bar{y}_{ij.}) = \frac{\sigma^2}{n_{ij}} \quad (\text{not the same for every cell})$$

Their average is  $\bar{V}(\bar{y}_{ij.}) = \frac{\sum \sum \frac{\sigma^2}{n_{ij}}}{ab}$

(11)

This can be estimated by

$$\frac{\hat{\sigma}^2}{ab} \sum_{i,j} \frac{1}{n_{ij}}$$

Yates said to use this as the manually computed MSE, with  $n_{..} - ab$  df.

(12)

Cell-means ANOVA

Source	df
A	a-1
B	b-1
AB	(a-1)(b-1)
Error	0
Total	ab - 1

Then, manually compute the MSE + reconstruct the ANOVA table

Source	df
A	a-1
B	b-1
AB	(a-1)(b-1)
Error	$n_{..} - ab$
Total	$n_{..} - 1$