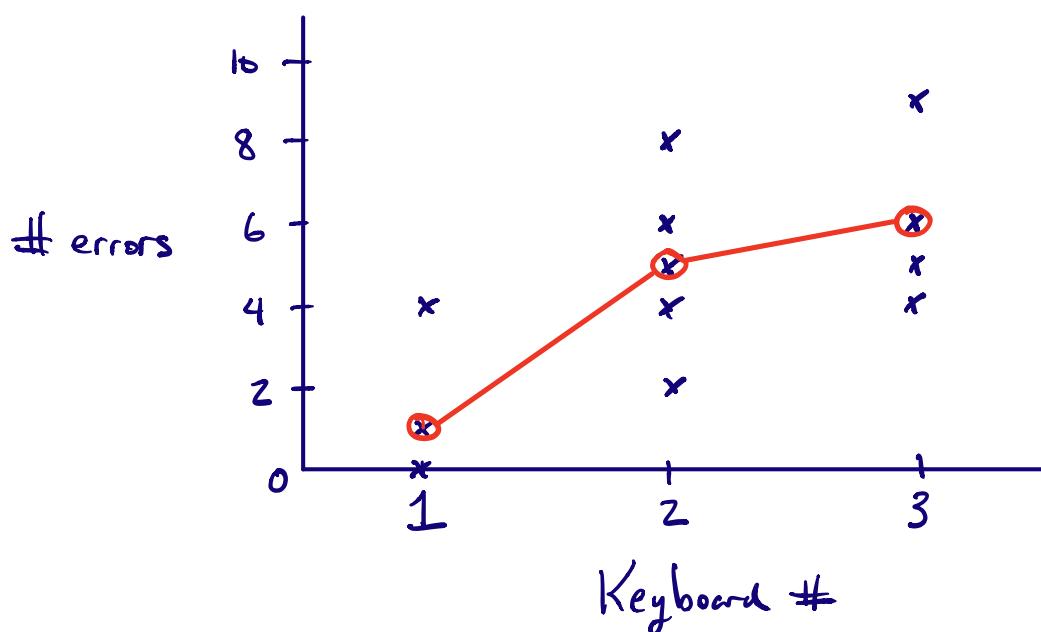


Lecture 3 - Inference in Single Factor Designs

Guiding example: An I/O psychologist studied three computer keyboard designs. Subjects were randomly assigned to type on one of the three designs, and the # of errors was recorded.

<u>Keyboard 1</u>	<u>Keyboard 2</u>	<u>Keyboard 3</u>
0	6	6
4	8	5
0	5	9
1	4	4
0	2	6
$\overline{x}_1 = 1$	$\overline{x}_2 = 5$	$\overline{x}_3 = 6$

To answer questions about these data, we need to put some structure on the data



Structure = statistical model

- * write a formula that generates each observed data point
- * "data story"

Linear model:

$$X_{ij} = \mu_j + \varepsilon_{ij}$$

where: X_{ij} = score of subject # i ($i=1, \dots, 5$)
in condition # j ($j=1, 2, 3$)

μ_j = population mean for condition # j ($j=1, 2, 3$)

ε_{ij} = deviation of subject i in condition j
from μ_j (i.e., model error)

In words, this means that each observed score can be written as the sum of

- (1) the condition mean
- (2) individual variation around the condition mean.

Inference goals:

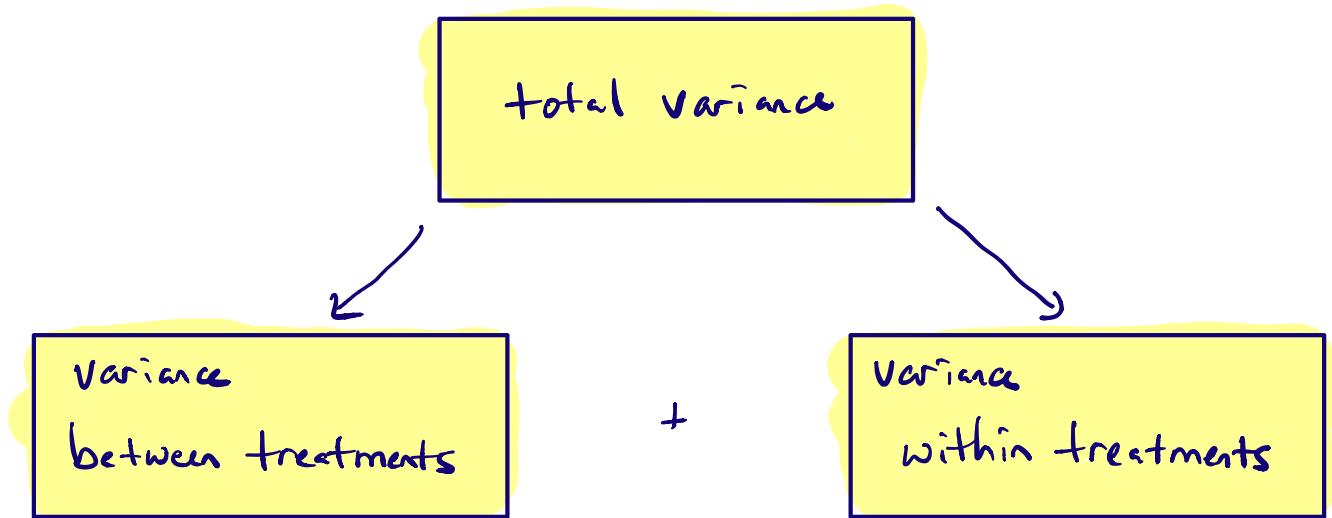
- (1) test hypotheses about the condition means μ_j
- (2) estimate the condition means μ_j

Hypothesis testing:

Consider $\begin{cases} H_0: \mu_1 = \mu_2 = \mu_3 \\ H_1: \text{not all } \mu_i \text{'s equal.} \end{cases}$. Which better predicts the observed data?

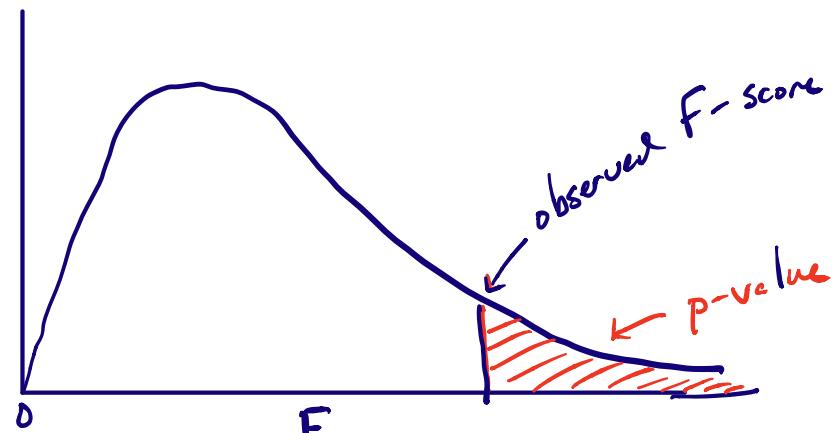
Recall: frequentist approach \rightarrow compute $p(\text{data} | H_0)$.

R.A. Fisher (1925) developed a creative method for computing this likelihood - partitioning the variance in the observed data.



$$F = \frac{\text{variance between treatments}}{\text{variance within treatments}}$$

Under H_0 , this statistic forms an F-distribution



So how do we compute these two different variances?

Recall: Variance = $\frac{SS}{df}$ ← need to compute
(1) SS between group means
(2) SS within groups

① To compute SS between group means, we sum the squared deviations of each group mean ($\bar{x}_1, \bar{x}_2, \bar{x}_3$) from the grand mean (\bar{x})

$$SS = (\bar{x}_1 - \bar{x})^2 + (\bar{x}_2 - \bar{x})^2 + (\bar{x}_3 - \bar{x})^2 \\ = \sum_{j=1}^3 (\bar{x}_j - \bar{x})^2$$

Technical note: this SS, when divided by the appropriate df, will estimate the variance of the distribution of sample means, σ_{means}^2 .

By the central limit theorem, $\sigma_{\text{means}}^2 = \frac{\sigma_{\text{subjects}}^2}{N}$

Equivalently, $\sigma_{\text{subjects}}^2 = N \cdot \sigma_{\text{means}}^2$

So we need to multiply by N to "scale up" to the level of variation among subjects:

$$SS_{\text{between}} = N \sum_{j=1}^3 (\bar{x}_j - \bar{x})^2$$

For our data:

$$\bar{x}_1 = 1$$

$$\bar{x}_2 = 5 \quad \bar{X} = 4 \quad N = 5$$

$$\bar{x}_3 = 6$$

$$\begin{aligned} SS_{\text{between}} &= N \sum_{j=1}^3 (\bar{x}_j - \bar{X})^2 \\ &= 5 \left[(1-4)^2 + (5-4)^2 + (6-4)^2 \right] \\ &= 5 (9 + 1 + 4) \\ &= 70 \end{aligned}$$

Convert to variance (which we denote MS for "mean squared deviation")

$$MS_{\text{between}} = \frac{SS_{\text{between}}}{df_{\text{between}}} = \frac{70}{2} = 35$$

② to compute SS within groups, we pool together the SS for each group:

$$SS_{\text{within}} = \sum_{j=1}^3 \sum_{i=1}^5 (x_{ij} - \bar{x}_j)^2$$

$$\begin{aligned} SS_{\text{within}} &= (0-1)^2 + (4-1)^2 + (0-1)^2 + (1-1)^2 + (0-1)^2 \quad \text{group 1} \\ &\quad + (6-5)^2 + (8-5)^2 + (5-5)^2 + (4-5)^2 + (2-5)^2 \quad \text{group 2} \\ &\quad + (6-6)^2 + (5-6)^2 + (9-6)^2 + (4-6)^2 + (6-6)^2 \quad \text{group 3} \\ &= 1 + 9 + 1 + 0 + 1 \\ &\quad + 1 + 9 + 0 + 1 + 9 \\ &\quad + 0 + 1 + 9 + 4 + 0 \\ &= 46 \end{aligned}$$

Convert to variance: since there are $5-1=4$ degrees of freedom in each of the three groups, we have a total of $4 \times 3 = 12$ degrees of freedom within groups

$$\rightarrow MS_{\text{within}} = \frac{46}{12} = 3.833$$

Finally, we can compute F:

$$F = \frac{\text{Variance between groups}}{\text{Variance within groups}} = \frac{MS_{\text{between}}}{MS_{\text{within}}} = \frac{35}{3.833} = 9.14$$

How surprising is this if H_0 true?

$$p = P(F > 9.14 \mid H_0) = 0.00387$$

So, data are rare under $H_0 \rightarrow$ reject H_0 in favor of H_1 .

Bayes factor?

* a quick and dirty approximation comes from Faulkenberry (2018)

→ convert observed F to BF_{01}

$$BF_{01} \approx \sqrt{n^{\frac{df_1}{2}} \left(1 + \frac{F df_1}{df_2} \right)^{-n}}$$

OR: use online app, described in Faulkenberry (2018)

https://tomfaulkenberry.shinyapps.io/anova_BFcalc

Estimation

Let's get a confidence interval for μ_2

$$95\% \text{ CI} = \bar{x}_2 \pm t_{df}^* \cdot \frac{\hat{\sigma}_2}{\sqrt{n_2}}$$

$$\bar{x} = \bar{x}_2 = 5$$

$$t_{df}^* = 2.78 \quad (\text{use } df = 5-1 = 4)$$

$$\hat{\sigma}_2 = \sqrt{MS_{\text{within}}} = \sqrt{3.833} = 1.96$$

$$\text{so } 95\% \text{ CI} = 5 \pm 2.78 \cdot \frac{1.96}{\sqrt{5}}$$

$$= 5 \pm 2.44$$

$$= (2.56, 7.44)$$