

## Lecture 34.

Recall (ANOVA for single Factor Experiments).

- we have a single variable (the factor) that takes on some small set of values (treatments) that the experimenter chooses.
- Given a many treatments, we assume a linear statistical model

$$Y_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad \begin{cases} i \in \{1, \dots, a\} \\ j \in \{1, \dots, n\} \end{cases}$$

random sample      common mean      treatment effects      Error term  
 $\varepsilon_{ij} \sim N(0, \sigma^2)$

- Under the assumption that  $\sum_{i=1}^a \tau_i = 0$ , we aim to test the hypotheses

$$H_0: \tau_1 = \tau_2 = \dots = \tau_a = 0, \quad H_1: \tau_i \neq 0 \text{ for at least one } i.$$

- Given a set of data  $\{y_{ij} : 1 \leq i \leq a, 1 \leq j \leq n\}$  we partition the total variation according

to the ANOVA Equation

$$\sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2 = n \sum_{i=1}^a (\bar{y}_{i..} - \bar{y}_{..})^2 + \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{i..})^2$$

$\underbrace{\hspace{10em}}$   $\underbrace{\hspace{10em}}$

$SS_{\text{treatments}}$   $SS_{\text{E}}$ .

$SS_T$

- From here, we form a test statistic:

$$F_0 = \frac{SS_{\text{treat}}/a-1}{SS_{\text{E}}/a(n-1)} = \frac{MS_{\text{treat}}}{MS_{\text{E}}}.$$

- The ANOVA f-test results in a rejection of  $H_0$  iff  $F_0 = f_0$  is too big, i.e.

$$f_0 > f_{\alpha, n-1, a(n-1)}$$

for a given significance level  $\alpha$ .

- The ANOVA calculations are usually summarized in an ANOVA table:

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F <sub>0</sub>
Treatments	$SS_{\text{Treatments}}$	$a - 1$	$MS_{\text{Treat}} = \frac{SS_{\text{Treat}}}{a-1}$	$MS_{\text{Treat}} / MS_E$
Error	$SS_E$	$a(n-1)$	$MS_E = \frac{SS_E}{a(n-1)}$	
Total	$SS_T$	$an - 1$		

Confidence Intervals on the treatment means:

- Recall that the treatment means are defined to be  $\mu_i = \mu + \tilde{\mu}_i$ ,  $1 \leq i \leq a$ .
- A point estimator for  $\mu_i$  is  $\hat{\mu}_i = \bar{Y}_{i \cdot} = \frac{\sum_{j=1}^n Y_{ij}}{n}$ .  
(i.e. just the usual sample mean)
- If we assume  $\epsilon_{ij} \sim N(\mu, \sigma^2)$ , then  $\bar{Y}_i \sim N(\mu_i, \frac{\sigma^2}{n})$
- If  $\sigma^2$  is known, we can construct a  $100(1-\alpha)\%$  CI based on the normal distribution.
- If  $\sigma^2$  is unknown, then we can use  $MS_E$  as an unbiased estimator of  $\sigma^2$ .

- In this case

$$T = \frac{\bar{Y}_{i\cdot} - \mu_i}{\sqrt{MSE/n}}$$

Follows a t-distribution with  $a(n-1)$  degrees of freedom, and so a  $100(1-\alpha)\%$  CI for  $\mu_i$  is

$$\bar{Y}_{i\cdot} - t_{\alpha/2, a(n-1)} \sqrt{\frac{MSE}{n}} \leq \mu_i \leq \bar{Y}_{i\cdot} + t_{\alpha/2, a(n-1)} \sqrt{\frac{MSE}{n}}$$

- similarly given any two means  $\mu_i, \mu_j$ ,  $\bar{Y}_{i\cdot} - \bar{Y}_{j\cdot}$  is a point estimator for  $\mu_i - \mu_j$ , and  $V(\bar{Y}_{i\cdot} - \bar{Y}_{j\cdot}) = \frac{\sigma^2}{n} + \frac{\sigma^2}{n} = \frac{2\sigma^2}{n}$  (by independence)

Therefore,

$$T = \frac{\bar{Y}_{i\cdot} - \bar{Y}_{j\cdot} - (\mu_i - \mu_j)}{\sqrt{2MSE/n}}$$

has a t-distribution with  $a(n-1)$  degrees of freedom, and so a  $100(1-\alpha)\%$  CI for  $\mu_i - \mu_j$  has critical points

$$(\bar{Y}_{i\cdot} - \bar{Y}_{j\cdot}) \pm t_{\alpha/2, a(n-1)} \sqrt{\frac{2MSE}{n}}$$

## A Note on Unequal Sample Sizes.

- Recall that we assumed that for each  $i \in \{1, \dots, a\}$ ,  $n_i = n$ . That is, for each treatment, the sample sizes are all the same.
- If we allow the  $n_i$ 's to be different, we can adjust the ANOVA equation slightly
- Let  $N = \sum_{i=1}^a n_i$ .

- Then

$$SS_T = \sum_{i=1}^a \sum_{j=1}^n y_{ij}^2 - \frac{y_{..}^2}{N}$$

$$SS_{\text{treat}} = \sum_{i=1}^a \frac{y_{i..}^2}{n_i} - \frac{y_{..}^2}{N}$$

and

$$SS_{\text{res}} = SS_T - SS_{\text{treat}}$$

- this is the unbalanced situation.

- balanced is better

↳ power is maximized.

↳ less sensitive to slightly varying  $\sigma_i^2$ .

## Last Topic Fisher's LSD test.

- Suppose we have performed an ANOVA test
- $H_0: \tau_i = \tau_j = 0 \forall i, j$ ,  $H_1: \tau_i \neq 0$  for some  $i$  and we have decided to reject  $H_0$ .
- ANOVA doesn't tell us which of the  $\tau_i \neq 0$ .
- we want a method for comparing  $\mu_i - \mu_j$
- there are many techniques, but we will discuss Fisher's "least significant differences" method.
- For each  $i \neq j$ , we create a statistic

$$t_{ij} = \frac{\bar{y}_{i\cdot} - \bar{y}_{j\cdot}}{\sqrt{\frac{2MSE}{n}}} \quad \left( \begin{array}{l} \text{since the order doesn't} \\ \text{matter, there are } \binom{a}{2} \\ \text{many such stats} \end{array} \right)$$

- Assuming  $n_i = n_j = n$  (same sample size for each treatment), we say that the difference between two means  $\mu_i$  and  $\mu_j$  is significant if

$$|\bar{y}_{i\cdot} - \bar{y}_{j\cdot}| > \underbrace{t_{\alpha/2, a(n-1)} \sqrt{\frac{2MSE}{n}}}_{\text{LSD.}}$$