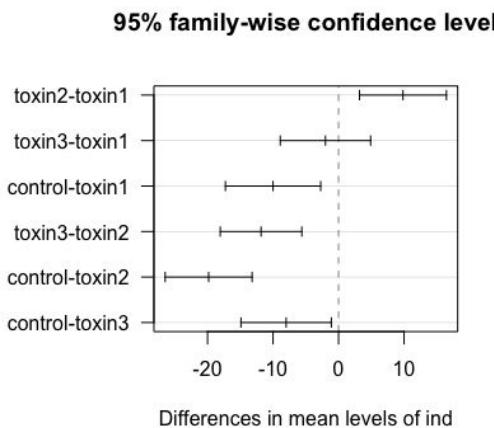


# Two-way ANOVA (univariate)

*ANOVA table*

07/29/2021

# In the previous lecture,



- Simultaneous pairwise comparisons:
  - Bonferroni method: adjusting with  $1 - \alpha/m$ .
  - Tukey's method: studentized range distribution.  
 $\text{Tukey's CI} \subseteq \text{Bonferroni CI}$ .
- Non-parametric Kruskal-Wallis test:
  - It extends Mann-Whitney test and focuses on medians;
  - The test statistic is:
$$H = \frac{12}{n(n+1)} \sum_{i=1}^k n_i \left( \frac{R_i}{n_i} - \frac{n+1}{2} \right)^2.$$
  - The rejection region is  $\{H \geq c\}$ .  
 $SS_B$
- Two-way factorial design:
  - Treatment means plot: check non-parallel lines to see whether there are interaction effects.
  - Model assumption:

$$Y_{ijl} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijl} .$$

$\stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

$$\overbrace{\mu + \alpha_i + \beta_j + \delta_{ij}} = \bar{Y}_{ij}.$$

## Two-way ANOVA: LRT

**Proposition D.** Assume  $Y_{ijl} = \underbrace{\mu + \alpha_i + \beta_j + \delta_{ij}} + \epsilon_{ijl}$ , where  $\epsilon_{ijl} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ . We can derive the maximum likelihood in  $\Theta = \Theta_0 \cup \Theta_1$ :

$$\sup_{\Theta} L = \left( \frac{1}{\sqrt{2\pi\hat{\sigma}_n^2}} \right)^n e^{-\frac{n}{2}},$$

where  $\hat{\sigma}_n^2 = n^{-1} \sum_i \sum_j \sum_{l=1}^{n_{ij}} (Y_{ijl} - \bar{Y}_{ij.})^2$ .

$SS_W$

Sum of squares within cells:

$$SS_E = \sum_{i=1}^I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} (Y_{ijl} - \bar{Y}_{ij.})^2.$$

$$\begin{aligned} & \underbrace{\bar{Y}_{ij1}, \dots, \bar{Y}_{ijn_{ij}}}_{\text{iid } N(\mu + \alpha_i + \beta_j + \delta_{ij}, \sigma^2)} \\ & \sum_{i=1}^I \sum_{j=1}^J \frac{1}{6} \sum_{l=1}^{n_{ij}} (\bar{Y}_{ijl} - \bar{Y}_{ij.})^2 \sim \chi^2_{n_{ij}-1} \\ & \perp \bar{Y}_{ij.} \end{aligned}$$

**Corollary D.** Under the two-way ANOVA assumptions,

$$\frac{SS_E}{\sigma^2} \sim \chi^2_{n-IJ}, \text{ and it's independent of each cell mean } \bar{Y}_{ij.},$$

where  $i = 1, \dots, I, j = 1, \dots, J$ .

$$\begin{aligned} & \sum_{i=1}^I \sum_{j=1}^J (n_{ij}-1) \\ & = n - \sum_{i=1}^I \sum_{j=1}^J 1 \\ & = n - IJ. \end{aligned}$$

$$F = \left\{ \lambda(Y_n) \leq c \right\} = \left\{ \frac{\hat{\sigma}_A^2}{\hat{\sigma}_n^2} \geq c' \right\}$$

## Two-way ANOVA: LRT for factor A

**Proposition E1.** Assume  $Y_{ijl} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijl}$ , where  $\epsilon_{ijl} \stackrel{iid}{\sim} N(0, \sigma^2)$ . The likelihood ratio for testing  $H_0 : \alpha_1 = \dots = \alpha_I = 0$  is equal to:

$$\lambda(Y_n) = \left( \frac{\hat{\sigma}_n^2}{\hat{\sigma}_n^2 + \hat{\sigma}_A^2} \right)^{n/2},$$

where  $\hat{\sigma}_A^2 = n^{-1} \sum_i n_{i..} (\bar{Y}_{i..} - \bar{Y}_{...})^2$ .

$$= \left\{ \frac{SS_A}{SSE} \geq c' \right\}$$

Sum of squares across rows (due to factor A):

$$SS_A = \sum_{i=1}^I n_{i..} (\bar{Y}_{i..} - \bar{Y}_{...})^2.$$

$$\mathbb{H}_0 = \left\{ \underbrace{\alpha_i = 0}_{I-1}, \underbrace{\sum_j \beta_j = 0}_{J-1}, \underbrace{\sum_i \delta_{ij} = \sum_j \delta_{ij} = 0}_{2}, \underbrace{\mu \in \mathbb{R}, \sigma^2 > 0}_{2} \right\}$$

$$\dim \mathbb{H} - \dim \mathbb{H}_0 = I - 1,$$

$$-2 \log \lambda(Y_n) \xrightarrow{d} \chi^2_{I-1} \text{ under } H_0.$$

Proof. We want to calculate  $\sup_{\text{H}_0} L$  under  $H_0: \alpha_1 = \dots = \alpha_I = 0$ .

$$\sup_{\text{H}_0} l(\vec{\alpha}, \vec{\beta}, \vec{\delta}, \mu, \sigma^2) = \sup_{\text{H}_0} \left[ \frac{n}{2} \log(2\pi b^2) - \frac{1}{2b^2} \sum_i \sum_j \sum_{l=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij})^2 \right]$$

$$\frac{\partial l}{\partial \mu} = \boxed{\frac{1}{b^2} \sum_i \sum_j \sum_{l=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij}) = 0}$$

$$\frac{\partial l}{\partial \beta_j} = \frac{1}{b^2} \sum_{i=1}^I \sum_{l=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij}) - \frac{1}{b^2} \sum_{i=1}^I \sum_{l=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij}) = 0$$

↑  
 $j = 1, \dots, J-1$

Sum all the likelihood equations for  $\frac{\partial l}{\partial \beta_j} = 0$ :

$$\frac{1}{b^2} \sum_{i=1}^I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij}) - \frac{J \cancel{\beta_j}}{b^2} = 0$$

$$\boxed{\sum_{i=1}^I \sum_{j=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij}) = 0}$$

Assume balanced design  
 $n_{ij} = M$

$$Y_{\cdot j} - M \sum_{i=1}^I \frac{1}{M} n_{ij} (\mu + \beta_j) - M \sum_{i=1}^I \frac{1}{M} n_{ij} \delta_{ij} = 0 \Rightarrow \mu + \beta_j = \bar{Y}_{\cdot j}$$

$$\frac{\partial l}{\partial \delta_{ij}} = \sum_{l=1}^{n_{ij}} (Y_{ijl} - \mu - \beta_j - \delta_{ij}) + \sum_{l=1}^{n_{ij}} (Y_{iJl} - \mu - \beta_J - \delta_{iJ})$$

↓

$$\sum_{i=1, \dots, I-1} \left( Y_{iJl} - \mu - \beta_j - \delta_{iJ} \right) - \sum_{j=1, \dots, J-1} \left( Y_{iJl} - \mu - \beta_J - \delta_{iJ} \right)$$

Sum all  $\frac{\partial l}{\partial \delta_{ij}} = 0$  to get:

$$0 = \sum_{I=1}^I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} + (I \boxed{J}) (J \boxed{I}) \sum_{l=1}^{n_{ij}} - (I \boxed{J}) \sum_{j=1}^J \sum_{l=1}^{n_{ij}} - (J \boxed{I}) \sum_{i=1}^I \sum_{l=1}^{n_{ij}} = 0$$

$$0 = \sum_I \sum_{l=1}^{n_{ij}} (Y_{iJl} - \mu - \beta_J - \delta_{iJ}) - \sum_I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} (Y_{iJl} - \mu - \beta_J - \delta_{iJ})$$

$$IM \bar{Y}_{iJ} = \frac{Y_{iJ}}{IM}$$

$$\frac{1}{JM} \left( Y_{IJ.} - \mu \bar{Y}_{J.} - \mu \delta_{iJ} \right) - \left( \frac{Y_{I..}}{JM} - \frac{Y_{...}}{IJM} \right) = 0$$

$$\bar{Y}_{IJ.} - \bar{Y}_{J.} - \delta_{iJ} - \bar{Y}_{I..} + \bar{Y}_{...} = 0 \Rightarrow \delta_{iJ} = \bar{Y}_{IJ.} - \bar{Y}_{J.} - \bar{Y}_{I..} + \bar{Y}_{...}$$

$$\underbrace{\mu + \beta_j + \delta_{ij}}_{\text{constant}} = \overline{Y}_{I\cdot} + \overline{Y}_{IJ\cdot} - \overline{Y}_{I\cdot\cdot} - \overline{Y}_{I\cdot\cdot\cdot}$$

$$= \overline{Y}_{IJ\cdot} - \overline{Y}_{I\cdot\cdot} + \overline{Y}_{I\cdot\cdot\cdot}$$

$$\Rightarrow \mu + \beta_j + \delta_{ij} = \underbrace{\overline{Y}_{IJ\cdot} - \overline{Y}_{I\cdot\cdot}}_{\text{constant}} + \overline{Y}_{I\cdot\cdot\cdot}$$

$$\Rightarrow \sup_{\mathbb{H}^n} L = \left( \frac{1}{\sqrt{2\pi b_0^2}} \right)^n e^{-\frac{n}{2}}$$

$$\begin{aligned} \Rightarrow \lambda(Y_n) &= \left( \frac{\hat{b}_n^2}{\hat{b}_0^2} \right)^{\frac{n}{2}} \\ &= \left( \frac{\hat{b}_n^2}{\hat{b}_A^2 + \hat{b}_n^2} \right)^{\frac{n}{2}} \quad \blacksquare \end{aligned}$$

where  $\hat{b}_0^2 = \sum_i \sum_j \sum_l \frac{n_{ij}}{n} (Y_{ijl} - \overline{Y}_{ij\cdot} + \overline{Y}_{i\cdot\cdot} - \overline{Y}_{\cdot\cdot\cdot})^2$

$$\begin{aligned} &= \frac{1}{n} \sum_i \sum_j \sum_l \left[ (Y_{ijl} - \overline{Y}_{ij\cdot})^2 + (\overline{Y}_{i\cdot\cdot} - \overline{Y}_{\cdot\cdot\cdot})^2 \right. \\ &\quad \left. + 2(Y_{ijl} - \overline{Y}_{ij\cdot})(\overline{Y}_{i\cdot\cdot} - \overline{Y}_{\cdot\cdot\cdot}) \right] \\ &= 0 \end{aligned}$$

$$\begin{aligned} &= \frac{1}{n} \sum_i \sum_j \sum_l \frac{n_{ij}}{n} (Y_{ijl} - \overline{Y}_{ij\cdot})^2 \\ &+ \frac{1}{n} \sum_i n_{ii} (\overline{Y}_{i\cdot\cdot} - \overline{Y}_{\cdot\cdot\cdot})^2 \end{aligned}$$

$$\frac{\hat{b}_n^2}{\hat{b}_A^2}$$



## Two-way ANOVA: LRT for factor B

$$P = \left\{ \frac{SS_B}{SSE} \geq c' \right\}$$

**Proposition E2.** Assume  $Y_{ijl} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijl}$ , where  $\epsilon_{ijl} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ . The likelihood ratio for testing  $H_0 : \beta_1 = \dots = \beta_J = 0$  is equal to:

$$\lambda(\mathbf{Y}_n) = \left( \frac{\hat{\sigma}_n^2}{\hat{\sigma}_n^2 + \hat{\sigma}_B^2} \right)^{n/2}.$$

where  $\hat{\sigma}_B^2 = n^{-1} \sum_j n_{\cdot j} (\bar{Y}_{\cdot j} - \bar{Y}_{\dots})^2$ .

Sum of squares across columns (due to factor B):

$$SS_B = \sum_{j=1}^J n_{\cdot j} (\bar{Y}_{\cdot j} - \bar{Y}_{\dots})^2.$$

↑    ↑

$$\dim(\mathbb{H}) - \dim(\mathbb{H}_0) = J - 1$$

## Two-way ANOVA: LRT for A\*B

$$\begin{aligned}\mu_{ij} &= \mu + \alpha_i + \beta_j + \delta_{ij} \\ \delta_{ij} &= \underline{\mu_{ij}} - \underline{\mu} - \underline{\alpha_i} - \underline{\beta_j} = \bar{Y}_{ij\cdot} - \bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdot j\cdot} + \bar{Y}_{\cdot\cdot\cdot}\end{aligned}$$

**Proposition E3.** Assume  $Y_{ijl} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijl}$ , where  $\epsilon_{ijl} \stackrel{iid}{\sim} N(0, \sigma^2)$ . The likelihood ratio for testing  $H_0$ : all  $\delta_{ij}$ 's are zero is equal to:

$$R = \left\{ \frac{SS_{AB}}{SSE} \geq c' \right\}$$

$$\lambda(\mathbf{Y}_n) = \left( \frac{\hat{\sigma}_n^2}{\hat{\sigma}_n^2 + \hat{\sigma}_{AB}^2} \right)^{n/2}.$$

$$\text{where } \hat{\sigma}_{AB}^2 = \sum_i \sum_j n_{ij} (\bar{Y}_{ij\cdot} - \bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdot j\cdot} + \bar{Y}_{\cdot\cdot\cdot})^2.$$

Sum of squares among cells (due to interaction between A and B):

$$SS_{AB} = \sum_{i=1}^I \sum_{j=1}^J n_{ij} \left( \bar{Y}_{ij\cdot} - \bar{Y}_{i\cdot\cdot} - \bar{Y}_{\cdot j\cdot} + \bar{Y}_{\cdot\cdot\cdot} \right)^2.$$

$$\mathbb{H}_0 = \left\{ \begin{array}{l} \sum_i \alpha_i = 0, \sum_j \beta_j = 0 \\ \delta_{ij} = 0, \mu \in \mathbb{R}, b^2 \geq 0 \end{array} \right\}$$

$$\dim \mathbb{H} - \dim \mathbb{H}_0 = (I-1)(J-1)$$

## Two-way ANOVA: Summary

**Theorem F.** Under the assumptions of two-way ANOVA model,

- 1) Denote the total sum of squares as

$$SS_{\text{Tot}} = \sum_{i=1}^I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} (Y_{ijk} - \bar{Y}_{...})^2.$$

- 2) Then  $SS_{\text{Tot}} = \underbrace{SS_A}_{\text{Under } H_0: \alpha_1 = \dots = \alpha_I = 0} + \underbrace{SS_B}_{\text{Under } H_0: \beta_1 = \dots = \beta_J = 0} + \underbrace{SS_{AB}}_{\text{Under } H_0: \text{all } \delta_{ij}\text{'s are zero,}} + \underbrace{SS_E}.$

$$\frac{SS_A}{\sigma^2} \sim \chi_{I-1}^2.$$

- 3) Under  $H_0: \beta_1 = \dots = \beta_J = 0$

$$\frac{SS_B}{\sigma^2} \sim \chi_{J-1}^2.$$

- 4) Under  $H_0: \text{all } \delta_{ij}\text{'s are zero,}$

$$\frac{SS_{AB}}{\sigma^2} \sim \chi_{(I-1)(J-1)}^2.$$

- 5)  $SS_A, SS_B, SS_{AB}$  and  $SS_E$  are mutually independent.

$$\sum_{i,j} (Y_{ijk} - \bar{Y}_{...})^2 \quad \text{a function of } Y_{ij}.$$

## Two-way ANOVA: Summary

$$R = \left\{ \begin{array}{l} \frac{SS_A}{SSE} \geq c' \\ \end{array} \right\}$$

$$= \left\{ \frac{SS_A / (I-1)}{SSE / (n - IJ)} \geq c'' \right\}$$

Source	df	Sum Sq	Mean Sq	F value	p-value
Factor A	$I - 1$	$SS_A = \sum_{i=1}^I n_i (\bar{Y}_{i..} - \bar{Y}_{...})^2$	$MS_A = SS_A / (I-1)$	$MS_A / MS_E$	$pf(MS_A / MS_E, I-1, n-I*J, lower.tail = F)$
Factor B	$J - 1$	$SS_B = \sum_{j=1}^J n_{.j} (\bar{Y}_{.j.} - \bar{Y}_{...})^2$	$MS_B = SS_B / (J-1)$	$MS_B / MS_E$	$pf(MS_B / MS_E, J-1, n-I*J, lower.tail = F)$
Interaction A*B	$(I-1)(J-1)$	$SS_{AB} = \sum_{i=1}^I \sum_{j=1}^J n_{ij} (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...})^2$	$MS_{AB} = SS_{AB} / ((I-1)(J-1))$	$MS_{AB} / MS_E$	$pf(MS_{AB} / MS_E, (I-1)*(J-1), n-I*J, lower.tail = F)$
Error	$n - IJ$	$SS_E = \sum_{i=1}^I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} (Y_{ijl} - \bar{Y}_{ij.})^2$	$MS_E = SS_E / (n - IJ)$		
Total	$n - 1$	$SS_{Tot} = \sum_{i=1}^I \sum_{j=1}^J \sum_{l=1}^{n_{ij}} (Y_{ijk} - \bar{Y}_{...})^2$			

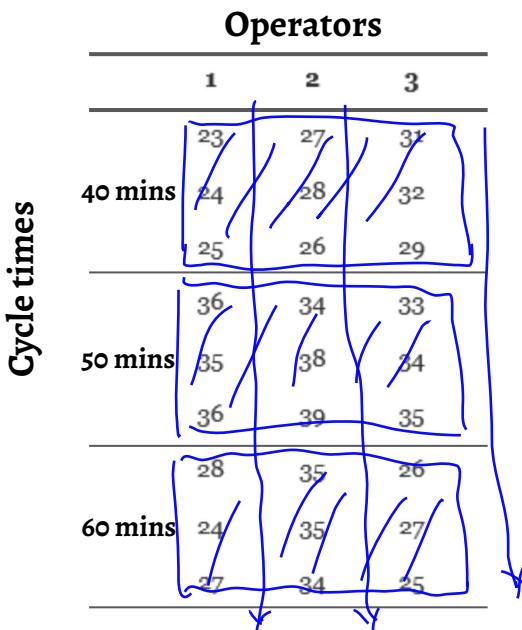
↑

↑

↑

## Two treatment factors

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected. The finished cloth was compared to a standard, and a numerical score was assigned.



$$SS_A = \sum_{i=1}^I n_i (\bar{Y}_{i..} - \bar{Y}_{...})^2$$

↓      ↓  
cycle-times    operator obs

```

> ## Bookkeeping
> I = 3; J = 3; n = nrow(Obs_w_names)
> all_mean = mean(Obs_w_names$Obs)

## Calculate SS_A
> S_A = aggregate(Obs ~ cycle_times, data = Obs_w_names,
+                   FUN = function(vec) length(vec)*(mean(vec)-all_mean)^2)
> S_A
   cycle_times      Obs
1        40 102.23457
2        50 221.67901
3        60 22.82716
}
> SS_A = sum(S_A$Obs)
> SS_A
[1] 346.7407

```

## Two treatment factors

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected. The finished cloth was compared to a standard, and a numerical score was assigned.

Operators			
	1	2	3
40 mins	23 24 25	27 28 26	31 32 29
50 mins	36 35 36	34 38 39	33 34 35
60 mins	28 24 27	35 35 34	26 27 25

$$SS_B = \sum_{j=1}^J n_{\cdot j} (\bar{Y}_{\cdot j} - \bar{Y}_{\dots})^2$$

```
## Calculate SS_B
> S_B = aggregate(Obs ~ operator, data = Obs_w_names,
   FUN = function(vec) length(vec)*(mean(vec)-all_mean)^2)

> SS_B = sum(S_B$Obs)
> SS_B
[1] 82.07407
```

## Two treatment factors

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected. The finished cloth was compared to a standard, and a numerical score was assigned.

Operators			
	1	2	3
40 mins	23 24 25	27 28 26	31 32 29
50 mins	36 35 36	64 38 39	33 34 35
60 mins	28 24 27	35 35 34	26 27 25

$$SS_E = \sum_{i=1}^I \sum_j^J \sum_{l=1}^{n_{ij}} (Y_{ijl} - \bar{Y}_{ij})^2$$

```
> ## Calculate SS_E
> S_E = aggregate(0bs ~ cycle_times + operator, data = 0bs_w_names,
+                 FUN = function(vec) sum((vec-mean(vec))^2))
> head(S_E, 6)
   cycle_times operator    0bs
1        40         1  2.000
2        50         1  0.667
3        60         1  8.667
4        40         2  2.000
5        50         2 14.000
6        60         2  0.667

> SS_E = sum(S_E$0bs)
> SS_E
[1] 36.66667
```

## Two treatment factors

$a = c(1,2,3)$

$\xrightarrow{\text{rep(a, times = 3)}}$        $\xrightarrow{\text{times}}$        $\xrightarrow{\text{each}}$

$1, 2, 3, 1, 2, 3, 1, 2, 3$

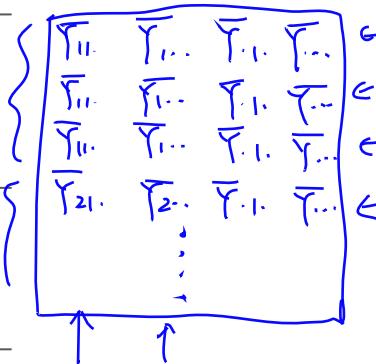
$1, 1, 1, 2, 2, 2, 3, 3, 3$

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected. The finished cloth was compared to a standard, and a numerical score was assigned.

Cycle times

Operators		
	1	2
1	23	27
40 mins	24	28
	25	26
36	34	33
50 mins	35	38
	36	39
28	35	26
60 mins	24	35
	27	34

	cycle_times	operator	Obs
1	40	1	23
2	40	1	24
3	40	1	25
4	50	1	36
5	50	1	35
6	50	1	36



$$SS_{AB} = \sum_{i=1}^I \sum_{j=1}^J n_{ij} (\bar{Y}_{ij} - \bar{Y}_{i..} - \bar{Y}_{.j} + \bar{Y}...)^2$$

## calculate SS\_AB

```

> cell_mean = aggregate(Obs ~ cycle_times + operator, data = Obs_w_names,
  FUN = mean)
> cell_mean_vec = rep(cell_mean$Obs, each = p)
>
> A_means = aggregate(Obs ~ cycle_times, data = Obs_w_names,
  FUN = mean)
> A_means_vec = rep(rep(A_means$Obs, times = J), each = p)
>
> B_means = aggregate(Obs ~ operator, data = Obs_w_names,
  FUN = mean)
> B_means_vec = rep(rep(B_means$Obs, each = I), each = p)
>
> Tmp_data = cbind(cell_mean_vec, A_means_vec, B_means_vec, all_mean)
> S_AB = apply(Tmp_data, 1, function(vec) (vec[1]-vec[2]-vec[3]+vec[4])^2)
> SS_AB = sum(S_AB)
> SS_AB
[1] 143.037

```

## Two treatment factors

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected. The finished cloth was compared to a standard, and a numerical score was assigned.

$$H_0 : \alpha_1 = \cdots = \alpha_I = 0$$

Reject  $H_0$

```
> ## Test main effect of A
> F_stat = {SS_A/(I-1)}/{SS_E/(n-I*J)}
> F_stat
[1] 85.10909

> # Critical value ←
> qf(0.05, I-1, n-I*J, lower.tail = FALSE)
[1] 3.554557

> # p-value
> pf(F_stat, I-1, n-I*J, lower.tail = FALSE)
[1] 6.691113e-10
```

$$H_0 : \beta_1 = \cdots = \beta_J = 0$$

Reject  $H_0$

```
> ## Test main effect of B
> F_stat = {SS_B/(J-1)}/{SS_E/(n-I*J)}
> F_stat
[1] 20.14545 ←

> # Critical value
> qf(0.05, J-1, n-I*J, lower.tail = FALSE)
[1] 3.554557 ←

> # p-value
> pf(F_stat, J-1, n-I*J, lower.tail = FALSE)
[1] 2.552966e-05
```

$$H_0 : \text{all } \delta_{ij} \text{'s are zero}$$

Reject  $H_0$

```
> ## Test interaction A*B
> F_stat = {SS_AB/((I-1)*(J-1))}/{SS_E/(n-I*J)}
> F_stat
[1] 17.55455 ←

> # Critical value
> qf(0.05, (I-1)*(J-1), n-I*J, lower.tail = FALSE)
[1] 2.927744 ←

> # p-value
> pf(F_stat, (I-1)*(J-1), n-I*J, lower.tail = FALSE)
[1] 5.003958e-06 ←
```

## Two treatment factors

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected. The finished cloth was compared to a standard, and a numerical score was assigned.

```
> fit <- lm(Obs ~ cycle_times*operator, data = Obs_w_names)
> anova(fit)
Analysis of Variance Table

Response: Obs
            Df  Sum Sq Mean Sq F value    Pr(>F)
cycle_times       2  346.74 173.370  85.109 6.691e-10 ***
operator          2   82.07  41.037  20.145 2.553e-05 ***
cycle_times:operator 4  143.04  35.759  17.555 5.004e-06 ***
Residuals        18   36.67   2.037
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

$Obs \sim cycle\_times + operator + cycle\_times : operator$



$$\text{Model 1: } y_{ijl} = \mu + \alpha_i + \beta_j + (\delta_{ij}) + \epsilon_{ijl} \quad \left. \right\} H_0: \text{Model 1 is better than Model 2}$$

Two treatment factors

$$\text{Model 2: } y_{ijl} = \mu + \alpha_i + \beta_j + \epsilon_{ijl}$$

$$\text{Model 2: } y_{ijl} = \mu + \alpha_i + \epsilon_{ijl}$$

**Example 4 cont'd.** A fabric plant is studying the effect of several factors on the dyeing of cotton-synthetic cloth. At 300°C, three operators and three cycle times were selected.

```
> fit <- lm(Obs ~ cycle_times*operator, data = Obs_w_names)
> anova(fit)
```

Analysis of Variance Table

Response: Obs

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
cycle_times	2	346.74	173.370	85.109	6.691e-10 ***
operator	2	82.07	41.037	20.145	2.553e-05 ***
cycle_times:operator	4	143.04	35.759	17.555	5.004e-06 ***
Residuals	18	36.67	2.037		

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

*Model comparisons*

$$H_0: \beta_j = 0, \delta_{ij} = 0$$

```
> # F test as a model comparison
> fit_0 <- lm(Obs ~ cycle_times + operator, data = Obs_w_names)
> anova(fit_0, fit, test='F')
```

Analysis of Variance Table

Model 1: Obs ~ cycle\_times + operator

Model 2: Obs ~ cycle\_times \* operator

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	22	179.704			
2	18	36.667	4	143.04	17.555 5.004e-06 ***
---					

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

← Model 1  
← Model 2

```
> fit_0 <- lm(Obs ~ cycle_times, data = Obs_w_names)
> anova(fit_0, fit, test='F')
```

Analysis of Variance Table

Model 1: Obs ~ cycle\_times

Model 2: Obs ~ cycle\_times \* operator

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	24	261.778			
2	18	36.667	6	225.11	18.418 8.719e-07 ***
---					

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

# One-way MANOVA

*Some basics*

07/29/2021

# Multivariate analysis of variance (MANOVA)

The one-way layout involves **one factor** in the experimental design:

Treatments				
1	2	3	...	k
$\rightarrow \mathbf{y}_{11}$	$\mathbf{y}_{21}$	$\mathbf{y}_{31}$	$\cdots$	$\mathbf{y}_{k1}$
$\rightarrow \mathbf{y}_{12}$	$\mathbf{y}_{22}$	$\mathbf{y}_{32}$	$\cdots$	$\mathbf{y}_{k2}$
$\vdots$	$\vdots$	$\vdots$	$\dots$	$\mathbf{y}_{k3}$
		$\mathbf{y}_{3n_3}$		$\vdots$
$\rightarrow \mathbf{y}_{1n_1}$				
	$\mathbf{y}_{2n_2}$			$\mathbf{y}_{kn_k}$

$$\mathbf{y} = (y^1, \dots, y^p)$$

Model assumption ( $j = 1, \dots, n_i, i = 1, \dots, k$ ):

$$Y_{ij} = \underbrace{\mu}_{\text{common mean level}} + \underbrace{\alpha_i}_{\text{unique effect due to treatment } i} + \underbrace{\epsilon_{ij}}_{\substack{\text{iid } N(\mathbf{0}, \Sigma) \\ \text{Gaussian}}}.$$

$$H_0 : \alpha_1 = \cdots = \alpha_k = 0 \quad \text{versus} \quad H_1 : \alpha_i = \alpha_j \text{ for some } i \neq j.$$

## One-way MANOVA: LRT

$$\begin{pmatrix} Y_{ij1} - \bar{Y}_{..1} \\ \vdots \\ Y_{ijp} - \bar{Y}_{..p} \end{pmatrix} (Y_{ij1} - \bar{Y}_{..1}, \dots, Y_{ijp} - \bar{Y}_{..p})' = \begin{pmatrix} (\bar{Y}_{ij1} - \bar{Y}_{..1})^2 & \cdots & \cdots \\ (\bar{Y}_{ij2} - \bar{Y}_{..2}) & (\bar{Y}_{ij2} - \bar{Y}_{..2}) & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

**Lemma G.** The LRT rejection region is solely based on the following identity:

$$\underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{..}) (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{..})'}_{\text{Total sum of squares}} = \underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{i.}) (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{i.})'}_{\text{SS within groups}} + \underbrace{\sum_{i=1}^k n_i (\bar{\mathbf{Y}}_{i.} - \bar{\mathbf{Y}}_{..}) (\bar{\mathbf{Y}}_{i.} - \bar{\mathbf{Y}}_{..})'}_{\text{SS between groups}}$$

$$\mathbf{T} = \mathbf{E} + \mathbf{B}$$

Proof.

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{Y}_{ij} - \bar{Y}_{i.} + \bar{Y}_{i.} - \bar{Y}_{..}) (\bar{Y}_{ij} - \bar{Y}_{i.} + \bar{Y}_{i.} - \bar{Y}_{..})' \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{Y}_{ij} - \bar{Y}_{i.}) (\bar{Y}_{ij} - \bar{Y}_{i.})' \\ &+ \sum_{i=1}^k n_i (\bar{Y}_{i.} - \bar{Y}_{..}) (\bar{Y}_{i.} - \bar{Y}_{..})'. \\ & \cancel{+ 2 \sum_{i=1}^k \left[ \sum_{j=1}^{n_i} (\bar{Y}_{ij} - \bar{Y}_{i.}) \right] (\bar{Y}_{i.} - \bar{Y}_{..})'} = 0 \end{aligned}$$

Treatments				
1	2	3	...	k
$\mathbf{y}_{11}$	$\mathbf{y}_{21}$	$\mathbf{y}_{31}$	...	$\mathbf{y}_{k1}$
$\mathbf{y}_{12}$	$\mathbf{y}_{22}$	$\mathbf{y}_{32}$	...	$\mathbf{y}_{k2}$
$\vdots$	$\vdots$	$\vdots$	...	$\mathbf{y}_{k3}$
		$\mathbf{y}_{3n_3}$		$\vdots$
$\mathbf{y}_{1n_1}$				
	$\mathbf{y}_{2n_2}$			$\mathbf{y}_{kn_k}$

$$E/T \Rightarrow \underbrace{ET^{-1}}$$

## One-way MANOVA: LRT

**Proposition G.** There are many different test statistics:

1) Wilk's Lambda:  $\Lambda^* = \frac{|E|}{|\mathbf{B} + E|} = \frac{|E|}{|T|}$

2) Hotelling-Lawley trace:  $T_0^2 = \text{trace}(\mathbf{B}\mathbf{E}^{-1})$ .

3) Pillai trace:  $V = \text{trace}(\mathbf{B}(\mathbf{B}+\mathbf{E})^{-1})$ .

4) Roy's maximum root: the largest eigenvalue of  $\mathbf{B}\mathbf{E}^{-1}$ .

MANOVA		
Source	d.f.	SSP
Treatments	$k - 1$	$\mathbf{B}$
Error	$n - k$	$\mathbf{E}$
Total	$n - 1$	$\mathbf{T}$

In the univariate case,

$$\lambda(\mathbf{Y}_n) = \left( \frac{SS_W}{SS_{\text{Tot}}} \right)^{n/2},$$

$$R = \left\{ \frac{SS_B}{SS_W} > c'' \right\}.$$

The distribution of these statistics under  $H_0$  is not straightforward and can only be approximated except when  $p$  is small.

## One-way MANOVA: LRT

**Example 5.** A researcher randomly assigns 33 subjects to one of 3 groups that receive: 1) technical dietary information interactively from an online website; 2) the same information from a nurse practitioner; 3) the information from a video tape made by the same nurse practitioner. The researcher looks at three different ratings of the presentation, difficulty, usefulness and importance, to determine if there is a difference in the modes of presentation.

	Group	Useful	Difficulty	Importance
1	1	19.6	5.15	9.5
2	1	15.4	5.75	9.1
3	1	22.3	4.35	3.3
		:		
12	2	17.1	9.00	7.5
13	2	15.7	5.30	8.5
14	2	14.9	9.85	6.0
		:		

```
# MANOVA test
res.man <- manova(cbind(Useful, Difficulty, Importance) ~ factor(Group),
                     data = dataset)
summary(res.man)
```

	Df	Pillai	approx F	num Df	den Df	Pr(>F)		
factor(Group)	2	0.47667	3.0248	6	58	0.01215 *		
Residuals	30							
---								
Signif. codes:	0	'***'	0.001	'**'	0.01	'*' 0.05	'. 0.1	' 1

# One-way MANOVA: LRT

```
> summary(res.man, 'Wilks')
      Df    Wilks approx F num Df den Df   Pr(>F)
factor(Group) 2 0.52579  3.5382       6     56 0.004859 ** 
Residuals     30
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> summary(res.man, 'Hotelling-Lawley')
      Df Hotelling-Lawley approx F num Df den Df   Pr(>F)
factor(Group) 2           0.89723  4.0375       6     54 0.002058 ** 
Residuals     30
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> summary(res.man, 'Roy')
      Df    Roy approx F num Df den Df   Pr(>F)
factor(Group) 2 0.89199  8.6225       3     29 0.0003023 ***
Residuals     30
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



```
# MANOVA test
res.man <- manova(cbind(Useful, Difficulty, Importance) ~ factor(Group),
                   data = dataset)
summary(res.man)
```

```
      Df    Pillai approx F num Df den Df   Pr(>F)
factor(Group) 2 0.47667  3.0248       6     58 0.01215 * 
Residuals     30
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Next Tuesday ...

Simple linear regression:

- Model assumptions;
- Geometric interpretations;
- Maximum likelihood estimation.

