

PSY 3307

# Analysis of Variance (ANOVA)

Jonathan A. Pedroza PhD

Cal Poly Pomona

2021-10-12

# Agenda

- Review t-tests
- Overview of Analysis of Variance (ANOVA)
  - Quick Review of all the different types of ANOVA
  - One-way ANOVA
  - Two-way ANOVA
  - Factorial ANOVA
  - Repeated-measures ANOVA
  - Mixed-effect ANOVA
- Components of ANOVA
- Performing ANOVA
- Post-hoc Tests (Tukey's HSD Test)
- Effect Size and  $\eta^2$
- JP is Including More

# Review t-tests

- Things we need to remember
  - How to read a table (z, t, and now F)
  - Differences between within- and between-designs
  - What is an IV, DV, and conditions

# Analysis of Variance (ANOVA)

- **Factor** is just another word for IV
- A **level** is the same thing as a condition (from t-tests) and similar to a t-test, ANOVA has to do with differences between or among sample means
  - however ANOVA does not have restrictions on the number of groups we can test
- **k** is the symbol for the number of levels in a factor
  - otherwise known as the number of conditions in an IV

*k = number of levels in factor*

# Example of the Absurdity

```
## Warning: Missing column names filled in: 'X1' [1]
```

```
##
```

```
## -- Column specification -----
```

```
## cols(
```

```
##   .default = col_double(),
```

```
##   state_fips_code = col_character(),
```

```
##   county_fips_code = col_character(),
```

```
##   fips_code = col_character(),
```

```
##   state_abbreviation = col_character(),
```

```
##   county_name = col_character(),
```

```
##   reading_scores_aian = col_logical(),
```

```
##   math_scores_aian = col_logical(),
```

```
##   communicable_disease = col_logical(),
```

```
##   cancer_incidence = col_logical(),
```

```
##   coronary_heart_disease_hospitalizations = col_logical(),
```

```
##   cerebrovascular_disease_hospitalizations = col_logical(),
```

```
##   smoking_during_pregnancy = col_logical(),
```

```
##   opioid_hospital_visits = col_logical(),
```

```
##   alcohol_related_hospitalizations = col_logical(),
```

```
##   motor_vehicle_crash_occupancy_rate = col_logical(),
```

```
##   on_road_motor_vehicle_crash_related_er_visits = col_logical(),
```

```
aov_find <- aov(adult_smoking ~ state_abbreviation, data = data)
summary(aov_find)
```

```
##              Df Sum Sq Mean Sq F value           Pr(>F)
## state_abbreviation    50   2.338  0.04676    80.92 <0.00000000000000002 ***
## Residuals           3142   1.815  0.00058
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
tukey_find <- TukeyHSD(aov_find)
tukey_find$`state_abbreviation`[1:10, 1:4]
```

##		diff	lwr	upr	p adj
##	AK-CA	0.087343654	0.065752343	0.10893497	0.000000000000000
##	AL-CA	0.079242299	0.062110931	0.09637367	0.000000000000000
##	AR-CA	0.083350172	0.066642931	0.10005741	0.000000000000000
##	AZ-CA	0.043470569	0.016330229	0.07061091	0.0000002016392
##	CO-CA	0.022238079	0.004924024	0.03955213	0.0003319958281
##	CT-CA	-0.001099152	-0.035556182	0.03335788	1.000000000000000
##	DC-CA	0.038867947	-0.030362089	0.10809798	0.9909882820559
##	DE-CA	0.042165472	-0.007583593	0.09191454	0.3079515102058
##	FL-CA	0.070667490	0.053536122	0.08779886	0.000000000000000
##	GA-CA	0.061787805	0.047121974	0.07645364	0.000000000000000

# Another Reason Why I Like Regression

- In a real-life scenario, you would already have a hypothesis where you would be interested in whether or not a state is different from the rest
- Therefore, you'd already have a reference group to compare everyone to and just need to run one test
- Below is the same test, just through a regression framework



```
lm_find <- lm(adult_smoking ~ state_abbreviation, data = data)
summary(lm_find)
```

```
##
## Call:
## lm(formula = adult_smoking ~ state_abbreviation, data = data)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-0.07620	-0.01314	-0.00115	0.01065	0.24411

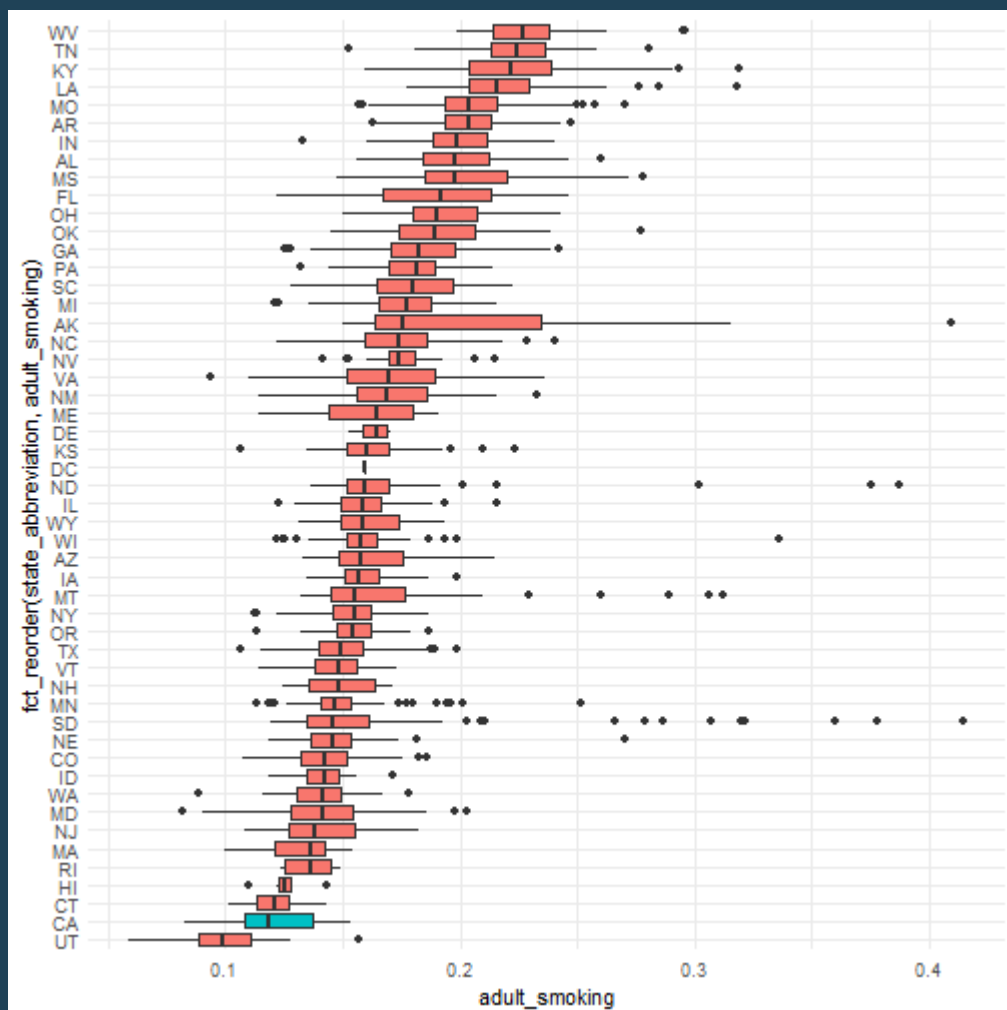
```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.120788	0.003129	38.598	< 0.000000000000000002 ***
state_abbreviationAK	0.087344	0.005390	16.205	< 0.000000000000000002 ***
state_abbreviationAL	0.079242	0.004277	18.529	< 0.000000000000000002 ***
state_abbreviationAR	0.083350	0.004171	19.984	< 0.000000000000000002 ***
state_abbreviationAZ	0.043471	0.006775	6.416	0.00000000001609179 ***
state_abbreviationCO	0.022238	0.004322	5.145	0.00000002839612803 ***
state_abbreviationCT	-0.001099	0.008602	-0.128	0.898331
state_abbreviationDC	0.038868	0.017283	2.249	0.024584 *
state_abbreviationDE	0.042165	0.012419	3.395	0.000694 ***
state_abbreviationFL	0.070667	0.004277	16.524	< 0.000000000000000002 ***
state_abbreviationGA	0.061788	0.003661	16.876	< 0.000000000000000002 ***
state_abbreviationHI	0.004873	0.010300	0.473	0.636138
state_abbreviationIA	0.037774	0.003946	9.573	< 0.000000000000000002 ***
state_abbreviationID	0.020741	0.004757	4.360	0.0000134333660949 ***

```
broom::tidy(lm_find)
```

```
## # A tibble: 51 x 5
```

##	term	estimate	std.error	statistic	p.value
##	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
##	1 (Intercept)	0.121	0.00313	38.6	4.13e-267
##	2 state_abbreviationAK	0.0873	0.00539	16.2	8.82e- 57
##	3 state_abbreviationAL	0.0792	0.00428	18.5	7.99e- 73
##	4 state_abbreviationAR	0.0834	0.00417	20.0	9.75e- 84
##	5 state_abbreviationAZ	0.0435	0.00678	6.42	1.61e- 10
##	6 state_abbreviationCO	0.0222	0.00432	5.14	2.84e- 7
##	7 state_abbreviationCT	-0.00110	0.00860	-0.128	8.98e- 1
##	8 state_abbreviationDC	0.0389	0.0173	2.25	2.46e- 2
##	9 state_abbreviationDE	0.0422	0.0124	3.40	6.94e- 4
##	10 state_abbreviationFL	0.0707	0.00428	16.5	7.02e- 59
##	# ... with 41 more rows				



# Breaking Down A One-way ANOVA

- **One-way ANOVA** is when we have an IV that has multiple levels (3+)
  - **NOTE** if you were to go on to SPSS and run a one-way ANOVA with the **sex** variable you would get the same answer
  - Essentially the same test being run
- Similar to a t-test, this also has within- and between-subjects designs
- Now instead of a t-table, we will be using a F-table

# We Are Now Working With Modeling

- There's just one problem. We have to work with ANOVA modeling
- **ANOVA** is a parametric procedure for determining whether significant differences occur in an experiment containing two or more sample means

$$X_{ij} = \mu + \gamma_j + \epsilon_{ij}$$

- $\mu$  is the grand mean
- $\gamma_j$  is the specific treatment effect for group  $j$  (which group you are interested in looking at)
- $\epsilon_{ij}$  is the error/residual of a specific individual (how much an individual deviates from the group's mean)

# Assumptions

- **homogeneity of variance** is the assumption that each population has the same variance

$$\sigma_1^2 = \sigma_2^2 = \sigma_3^2 = \dots$$

- **error variance** variance unrelated to any treatment differences
- **heterogeneity of variance** is when populations have different variances
- **normality** DV values are normally distributed
- **independence** observations are independent of one another
  - it really is that the residual/error is independent but for now we'll keep it as observations are different from one another
- The **n** in each level doesn't have to be exactly the same but they should not be drastically different

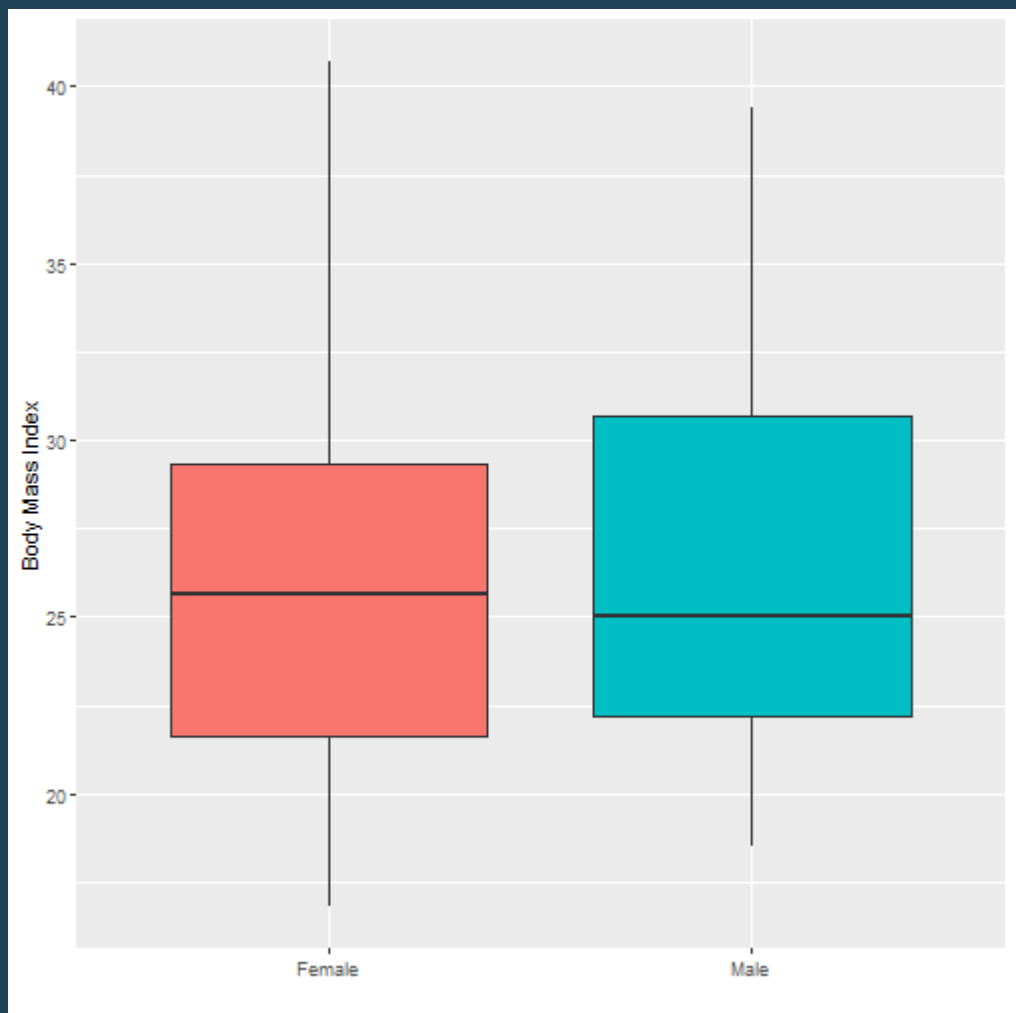
# Controlling Experiment-Wise Error Rate

- I had mentioned this previously that multiple independent-samples t-tests could do the same thing as a one-way ANOVA
- **experiment-wise error rate** is the probability of making a Type I error when comparing all means in an experiment
- with an F-test we are less likely to commit a type I error because we are not running all the tests possible

# Example

```
jp <- rio::import(here::here("jp_thesis_1.sav")) %>%  
  janitor::clean_names() %>%  
  rowid_to_column() %>%  
  rename(sex = ccc_gender)
```





```
bmi_aov <- aov(ccc_bmi ~ factor(sex), data = jp)
summary(bmi_aov)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## factor(sex)   1      3    2.515   0.086  0.769
## Residuals    370  10793   29.169
```

```
TukeyHSD(bmi_aov)
```

```
##      Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = ccc_bmi ~ factor(sex), data = jp)
##
## $`factor(sex)`
##              diff              lwr              upr              p adj
## 2-1 -0.1775135 -1.366163  1.011136  0.7691805
```

```
t.test(ccc_bmi ~ factor(sex), data = jp, var.equal = TRUE)
```

```
##  
##      Two Sample t-test  
##  
## data:  ccc_bmi by factor(sex)  
## t = 0.29366, df = 370, p-value = 0.7692  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
##  -1.011136  1.366163  
## sample estimates:  
## mean in group 1 mean in group 2  
##      26.32931      26.15180
```

# Steps to Conduct ANOVA

- Hypotheses
  - Similar to the independent-samples t-test, a one-way ANOVA's null hypothesis states that there is no difference between the levels/conditions
  - We just have more levels in a one-way ANOVA.

- our null hypothesis would be

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

$$H_1 : \text{not all } \mu \text{ are equal}$$

- our alternative/research hypothesis is not that all groups would be significantly different from one another
- **JP Note** another option is to state that one group will be significantly different from the other groups

# F statistic

- Steps to ANOVA
  1. Compute a F-statistic
  2. Conduct a post-hoc test

# F-statistic

- we compute a F-statistic to see if any means are different
  - Significant F-statistic then there are differences somewhere between the multiple levels
  - Non-significant F-statistic means there are no differences between any levels
- The F-statistic only tells us whether or not a significant difference is found between any of the levels
- F-obtained is compared to a F-critical value to find statistical significance

# Post-hoc Tests OR Planned Comparisons

- Post-hoc often refers to after the fact
- Post-hoc tests are often considered after finding a statistically significant F-statistic
- like t-tests when comparing all combinations of each levels to see if there differences between two specific levels
- Planned comparisons are when you are interested in having a specific levels being different from the other levels

# Post-hoc Tests

- only look at post-hoc findings if there is a significant F-statistic
- no post-hoc tests are run when you only have two levels



# Different Types of Post-hoc Tests

- Tukey Test or Tukey's HSD (Honestly Significant Difference)
- Fisher's Least Significant Difference (LSD) Procedure
- Newman-Keuls Test
- Scheffe Test
- Dunnett's Test
- Benjamini-Hochberg Test
- Bonferroni Test
- We'll get to these in the upcoming weeks (if only shortly)

# Componentenets of ANOVA

$$S_x^2 = \frac{\Sigma(X - \bar{X})^2}{N - 1}$$

$$S_x^2 = \frac{\text{Sum of Squares (SS)}}{\text{Degrees of Freedom (df)}}$$

- Sum of squares (SS)/degrees of freedom (df) is equal to mean square
- mean square is often seen as MS
- when referring to ANOVA, sum of the squared deviations is called sum of squares
- this calculation of sum of squares divided by the degrees of freedom is called mean square or MS
- so our variance calculation is really our mean square in ANOVA
  - this is because we are calculating the mean square within groups and mean square between groups

# Mean Square Within Groups

- **MS within groups** describes the variability in scores within the conditions/levels of an experiment

$$MS_{wn}$$

- we find differences among values in each level/condition and "pool" them together
- **MS within groups** is the "average" variability of scores within each level
- it is essentially a measure of variability of individual scores

# Mean Square Between Groups

- **MS between groups** is the differences between the means of each condition/level in a factor/IV

$$MS_{bn}$$

- measure difference between the means by treating them as scores, with an "average" amount they deviate from their mean, which in this case is the overall mean of the experiment
- similar to how scores deviate from a mean, this is a measure of the deviations of sample means from the overall mean

# The F-ratio

- MS between groups tells us whether or not the levels differ from one another and support our null/alternative hypotheses
- MS within groups estimates variability of individual scores
- if working with one population, our MS between should equal our MS within
- so if our null is true the MS between should be the same answer as the MS within
- **F-ratio** is a fraction consisting of MS between divided by the MS within or

$$F_{obt} = \frac{MS_{bn}}{MS_{wn}}$$

# F-ratio

- while it is unlikely that an exact value of 1 would be a F-obtained value, a value around 1 should be supportive of the null hypothesis
- the larger the F-ratio the more likely that the result is from sampling error or from the IV
- if our F-obtained value is larger than the F-critical value then we reject the null and accept the alternative/research hypothesis

# Performing an ANOVA

- before beginning, its important to note that for many of the calculations symbols will be similar like **MS** but it is important to note the subscripts

$$MS_{bn} \neq MS_{wn}$$

<i>Source</i>	<i>Sum of Squares</i>	<i>df</i>	<i>MeanSquare</i>	<i>F – ratio</i>
<i>Between</i>	$SS_{bn}$	$df_{bn}$	$MS_{bn}$	$F_{obt}$
<i>Within</i>	$SS_{wn}$	$df_{wn}$	$MS_{wn}$	
<i>Total</i>	$SS_{total}$	$df_{total}$		



# Computing for a F-obtained value

Steps for getting the F-obtained value

1. Calculate the sum of squares
2. Calculate the degrees of freedom
3. Calculate the mean squares
4. Calculate the F-obtained value

```
difficulty <- data.frame(easy = c(9, 12, 4, 8, 7),  
                          medium = c(4, 6, 8, 2, 10),  
                          hard = c(1, 3, 4, 5, 2))
```

```
easy_sum = 9+12+4+8+7  
med_sum = 4+6+8+2+10  
hard_sum = 1+3+4+5+2
```

```
easy_sum
```

```
## [1] 40
```

```
med_sum
```

```
## [1] 30
```

```
hard_sum
```

```
## [1] 15
```

```
easy_sum2 = 9^2+12^2+4^2+8^2+7^2  
med_sum2 = 4^2+6^2+8^2+2^2+10^2  
hard_sum2 = 1^2+3^2+4^2+5^2+2^2  
  
easy_sum2
```

```
## [1] 354
```

```
med_sum2
```

```
## [1] 220
```

```
hard_sum2
```

```
## [1] 55
```

```
easy_n = 5
med_n = 5
hard_n = 5

easy_mean = easy_sum/easy_n
med_mean = med_sum/med_n
hard_mean = hard_sum/hard_n

easy_mean
```

```
## [1] 8
```

```
med_mean
```

```
## [1] 6
```

```
hard_mean
```

```
## [1] 3
```

# Total Sum of all the Values

```
total_sum = easy_sum + med_sum + hard_sum  
total_sum
```

```
## [1] 85
```

# Total Sum of all the Squared Values

```
total_sum2 = easy_sum2 + med_sum2 + hard_sum2  
total_sum2
```

```
## [1] 629
```

# Total N

```
total_n = easy_n + med_n + hard_n  
total_n
```

```
## [1] 15
```



# k or the number of levels/conditions

```
k = 3  
k
```

```
## [1] 3
```

# Computing the Sums of Squares Total

$$SS_{total} = \Sigma X_{total}^2 - \frac{(\Sigma X_{total})^2}{N}$$

$$SS_{total} = 629 - \frac{(85)^2}{15}$$

```
85^2
```

```
## [1] 7225
```

$$SS_{total} = 629 - \frac{7225}{15}$$

```
## [1] 481.6667
```

$$SS_{total} = 629 - 481.67$$

```
629 - 481.67
```

```
## [1] 147.33
```

$$SS_{total} = 147.33$$

# Filling in the Table

$\left[ \begin{array}{cc} \text{Source} & \text{Between} & \text{Within} & \text{Total} \end{array} \right]$

$\left[ \begin{array}{cc} \text{Sum of Squares} & SS_{bn} & SS_{wn} & 147.33 \end{array} \right]$

$\left[ \begin{array}{cc} df & df_{bn} & df_{wn} & df_{total} \end{array} \right]$

$\left[ \begin{array}{cc} \text{Mean Square} & MS_{bn} & MS_{wn} & \end{array} \right]$

$\left[$

$F - ratio \ F_{obt}$

$\right]$

# Sums of Squares Between

$$SS_{bn} = \Sigma \left( \frac{(\Sigma X \text{ in each column})^2}{n \text{ in each column}} \right) - \frac{(\Sigma X_{total})^2}{N}$$

$$SS_{bn} = \Sigma \left( \frac{(40)^2}{5} + \frac{(30)^2}{5} + \frac{(15)^2}{5} \right) - \frac{(85)^2}{15}$$

```
40^2
```

```
## [1] 1600
```

```
30^2
```

```
## [1] 900
```

```
15^2
```

```
## [1] 225
```

```
85^2
```

```
## [1] 7225
```

$$SS_{bn} = \Sigma\left(\frac{1600}{5} + \frac{900}{5} + \frac{225}{5}\right) - \frac{7225}{15}$$



```
1600/5
```

```
## [1] 320
```

```
900/5
```

```
## [1] 180
```

```
225/5
```

```
## [1] 45
```

```
7225/15
```

```
## [1] 481.6667
```

$$SS_{bn} = (320 + 180 + 45) - 481.67$$

```
(320 + 180 + 45) - 481.67
```

```
## [1] 63.33
```

$$SS_{bn} = 63.33$$

# Filling in the Table

Source Between Within Total

Sum of Squares

63.33  $SS_{wn}$  147.33

df

df  $df_{bn}$   $df_{wn}$   $df_{total}$

Mean Square  $MS_{bn}$   $MS_{wn}$

F-ratio

$F - ratio$   $F_{obt}$

# Sum of Squares Within Groups

$$SS_{wn} = SS_{total} - SS_{bn}$$

```
147.33 - 63.33
```

```
## [1] 84
```

# Filling in the Table

Source Between Within Total

*Sum of Squares* 63.33 84 147.33

df

Mean Square

*F – ratio*

# Calculating degrees of freedom

- df between groups is simply the number of groups/levels/conditions - 1

$$df_{bn} = k - 1$$

```
3 - 1
```

```
## [1] 2
```

- df within groups is  $N - k$

```
15 - 3
```

```
## [1] 12
```

- df total is still  $N - 1$

```
15 - 1
```

```
## [1] 14
```

# Filling in the Table

$$\begin{array}{cc} \text{Source} & \text{Between} & \text{Within} & \text{Total} \end{array}$$

$$\begin{array}{c} \text{Sum of Squares} \end{array}$$

$$63.33 \quad 84 \quad 147.33$$

$$\begin{array}{c} \text{df} \end{array}$$

$$\begin{array}{c} \text{Mean Square} \end{array}$$

$$2 \quad 12 \quad 14$$

$$\begin{array}{c} \text{MS}_{bn} \quad \text{MS}_{wn} \end{array}$$

$$\begin{array}{cc} \text{Mean Square} & \text{MS}_{bn} & \text{MS}_{wn} \end{array}$$

$$\begin{array}{c} F - \text{ratio} \end{array}$$

$$F_{obt}$$

$$\end{array}$$

# Computing the Mean Squares

$$MS_{bn} = \frac{SS_{bn}}{df_{bn}}$$

$$MS_{bn} = \frac{63.33}{2}$$



63.33/2

## [1] 31.665

$$MS_{bn} = 31.67$$

# Filling in the Table

Source Between Within Total

*Sum of Squares* 63.33 84 147.33

*df* 2 12 14

*Mean Square* 31.67  $MS_{wn}$

# Computing the Mean Square within Groups

$$MS_{wn} = \frac{SS_{wn}}{df_{wn}}$$

$$MS_{wn} = \frac{84}{12}$$

```
## [1] 7
```

$$MS_{wn} = 7$$

# Filling in the Table

Source Between Within Total

*Sum of Squares* 63.33 84 147.33

*df* 2 12 14

*Mean Square* 31.67 7

# Calculating the F-statistic

$$F_{obt} = \frac{MS_{bn}}{MS_{wn}}$$

$$F_{obt} = \frac{MS_{bn}}{MS_{wn}}$$

31.67/7

```
## [1] 4.524286
```

$$F_{obt} = 4.52$$

# Filling in the Table

Source Between Within Total

*Sum of Squares* 63.33 84 147.33

*df* 2 12 14

*Mean Square* 31.67 7



# Interpreting the F-obtained value

- **F-distribution** is the sampling distribution with values of F to represent when  $H_0$  is true and all conditions represent one population (no differences between groups/levels)
- F cannot be less than zero
- There is no limit to how large an F-obtained value can be
- the mean of the F-distribution is 1
- F-distribution shape also depends on the df

# F-table

- uses alpha, df within, df between
- line up the df within with the df between and choose the value based on the alpha decided on
- the F-table only tells us one thing, *is there a statistically significant difference between the means of the three groups/conditions/levels*
- in order to see which specific group comparisons are significantly different from one another, we will rely on post-hoc tests or examination of the contrasts

