## **ANOVA** Tests

Prof. Wells

STA 209, 5/5/23

### Outline

In this lecture, we will...

- Construct a statistic to measure the differences in mean among several groups
- Discuss the theoretical and simulation-based distribution of the F statistic
- Use ANOVA to test for a difference in means among several groups

# Section 1

# Differences Among Several Populations

• Suppose we want to determine in a population whether a quantitative variable is independent of another variable.

- Suppose we want to determine in a population whether a quantitative variable is independent of another variable.
- Previously, we...

- Suppose we want to determine in a population whether a quantitative variable is independent of another variable.
- Previously, we...
  - Used a *t*-test for difference in means to determine if a *quantitative* variable and a *categorical* variable (with only 2 levels) were independent.

- Suppose we want to determine in a population whether a quantitative variable is independent of another variable.
- Previously, we...
  - Used a *t*-test for difference in means to determine if a *quantitative* variable and a *categorical* variable (with only 2 levels) were independent.
- The Analysis of Variance (ANOVA) test will allow us to assess whether the mean values of a quantitative variable differ across the levels of a categorical variable.

- Suppose we want to determine in a population whether a quantitative variable is independent of another variable.
- Previously, we...
  - Used a *t*-test for difference in means to determine if a *quantitative* variable and a *categorical* variable (with only 2 levels) were independent.
- The Analysis of Variance (ANOVA) test will allow us to assess whether the mean values of a quantitative variable differ across the levels of a categorical variable.
  - This gives a test to determine if a *quantitative* variable and a *categorical* variable (with more than 2 levels) are independent.

#### There's No Accounting For Taste

*Research Question*: Certainly, individual tastes in movie genres vary. But in aggregate, do audience ratings of movies depend on genre? To answer, we assess the Rotten Tomatoes audience rating for 132 films from 2011 spread across 7 genera.

#### There's No Accounting For Taste

*Research Question*: Certainly, individual tastes in movie genres vary. But in aggregate, do audience ratings of movies depend on genre? To answer, we assess the Rotten Tomatoes audience rating for 132 films from 2011 spread across 7 genera.

Movie	AudienceScore	Genre
Insidious	65	Horror
Paranormal Activity 3 Bad Teacher	58 38	Horror Comedy
Bridesmaids	77	Comedy
Midnight in Paris	84	Romance
The Help	91	Drama

### There's No Accounting For Taste

*Research Question*: Certainly, individual tastes in movie genres vary. But in aggregate, do audience ratings of movies depend on genre? To answer, we assess the Rotten Tomatoes audience rating for 132 films from 2011 spread across 7 genera.

Movie	AudienceScore	Genre
Insidious	65	Horror
Paranormal Activity 3	58	Horror
Bad Teacher	38	Comedy
Bridesmaids	77	Comedy
Midnight in Paris	84	Romance
The Help	91	Drama

- Observational unit: a single film
- Sample: 132 films from 2011
- Population: All films (maybe from last 20 years?)
- Variables: Audience Rating and Genre
- **Parameters**: Average audience rating for each genre,  $\mu_1, \ldots, \mu_7$ .
- Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
- Alternative Hypothesis: At least one  $\mu$  is not equal to the others

#### Data Exploration

Do ratings differ by genre?

#### Data Exploration

#### Do ratings differ by genre?



#### Data Exploration

#### Do ratings differ by genre?



```
movies %>% group_by(Genre) %>%
summarize(number = n(), avg_rating = mean(AudienceScore), st_dev = sd(AudienceScore))
```

##	#	A tibble:	7 x 4		
##		Genre	number	avg_rating	st_dev
##		<fct></fct>	<int></int>	<dbl></dbl>	<dbl></dbl>
##	1	Action	32	58.6	18.4
##	2	Animation	12	64.1	13.9
##	3	Comedy	27	59.1	15.7
##	4	Drama	21	72.1	14.5
##	5	Horror	17	48.6	15.9
##	6	Romance	10	64.8	12.9
##	7	Thriller	12	67.7	9.01

##			Movie	number	avg_rating	st_dev
##	1	A11	Films	131	61.38931	16.61351

• We saw a clear visual difference in mean scores for different genres.

- We saw a clear visual difference in mean scores for different genres.
  - But maybe we would see a similar difference just by separating into 7 arbitrary groups

- We saw a clear visual difference in mean scores for different genres.
  - But maybe we would see a similar difference just by separating into 7 arbitrary groups



7 Arbitrary Groups

- We saw a clear visual difference in mean scores for different genres.
  - But maybe we would see a similar difference just by separating into 7 arbitrary groups







#### There's No Accounting for Taste... But There is Accounting for Variance

Which of the following experiments gives *strongest* evidence of a difference in population means? Which gives the *weakest* evidence?

#### There's No Accounting for Taste... But There is Accounting for Variance

Which of the following experiments gives *strongest* evidence of a difference in population means? Which gives the *weakest* evidence?



### There's No Accounting for Taste... But There is Accounting for Variance

Which of the following experiments gives *strongest* evidence of a difference in population means? Which gives the *weakest* evidence?



- Strongest: Experiment 2
- Weakest: Experiment 3

• To assess whether a collection of population means are equal, we treat the sample means as a data set.

- To assess whether a collection of population means are equal, we treat the sample means as a data set.
  - If the population means are equal, the sample means should be close.
  - Otherwise, the sample means should be spread out.

- To assess whether a collection of population means are equal, we treat the sample means as a data set.
  - If the population means are equal, the sample means should be close.
  - Otherwise, the sample means should be spread out.
- But how spread out do sample means need to be to give good evidence that population means are different?

- To assess whether a collection of population means are equal, we treat the sample means as a data set.
  - If the population means are equal, the sample means should be close.
  - Otherwise, the sample means should be spread out.
- But how spread out do sample means need to be to give good evidence that population means are different?
  - If the populations have large standard deviation, we would expect the sample means to exhibit greater spread (even if the population means are equal)

- To assess whether a collection of population means are equal, we treat the sample means as a data set.
  - If the population means are equal, the sample means should be close.
  - Otherwise, the sample means should be spread out.
- But how spread out do sample means need to be to give good evidence that population means are different?
  - If the populations have large standard deviation, we would expect the sample means to exhibit greater spread (even if the population means are equal)
- Is the variation observed among sample means greater than what can be explained by variability in observations within each group alone?



• The Total Variability among all observations is the sum of Variability Between Groups and Variability Within Groups



- The Total Variability among all observations is the sum of Variability Between Groups and Variability Within Groups
- Variability Between Groups: How much do means vary?
  - Compare red dots



- The Total Variability among all observations is the sum of Variability Between Groups and Variability Within Groups
- Variability Between Groups: How much do means vary?
  - Compare red dots
- Variability Within Groups: How much do observations in groups vary from mean?
  - Within each group, compare black dots to red dot

Prof. Wells

• TotalVariability = Variability Between Groups + Variability Within Groups

- TotalVariability = Variability Between Groups + Variability Within Groups
- Variability Between Groups: How much do means vary?

Variability Between Groups (SSG) = 
$$\sum n_i (\bar{x}_i - \bar{x})^2$$

• The sum is across the k different groups

- TotalVariability = Variability Between Groups + Variability Within Groups
- Variability Between Groups: How much do means vary?

Variability Between Groups (SSG) = 
$$\sum n_i (\bar{x}_i - \bar{x})^2$$

- The sum is across the k different groups
- Variability Within Groups: How much do observations in groups vary from mean?

Variability Within Groups (SSE) = 
$$\sum (x_i - \bar{x}_i)^2$$

• The sum is across the n different observations, but each uses its group's mean

- TotalVariability = Variability Between Groups + Variability Within Groups
- Variability Between Groups: How much do means vary?

Variability Between Groups (SSG) 
$$= \sum n_i (\bar{x}_i - \bar{x})^2$$

- The sum is across the k different groups
- Variability Within Groups: How much do observations in groups vary from mean?

Variability Within Groups (SSE) =  $\sum (x_i - \bar{x}_i)^2$ 

- The sum is across the n different observations, but each uses its group's mean
- Total Variability: How much do observations vary from overall mean?

Variability Within Groups (TSS) 
$$=\sum (x - \bar{x})^2$$

• The sum is across the *n* different observations and each uses the mean for all data.

### Mean Squares

- Suppose we have a single population of 120 people which we divide randomly into...
  - a groups
  - 6 20 groups

### Mean Squares

- Suppose we have a single population of 120 people which we divide randomly into...
  - 3 groups
  - 6 20 groups
- All else equal, which of these divisions do we expect to have higher SSG?

$$\mathrm{SSG} \;=\; \mathrm{Variability\;Between\;Groups\;} = \sum n_i (\bar{x}_i - \bar{x})^2$$

### Mean Squares

- Suppose we have a single population of 120 people which we divide randomly into...
  - 3 groups
  - 6 20 groups
- All else equal, which of these divisions do we expect to have higher SSG?

$$\mathrm{SSG} ~=~ \mathrm{Variability~Between~Groups} ~= \sum n_i (\bar{x}_i - \bar{x})^2$$

• We standardize sum of squares to compare SSG to SSE

Mean Variability Between Groups (MSG) = 
$$\frac{\text{SSG}}{k-1}$$
  
Mean Variability Within Groups (MSE) =  $\frac{\text{SSE}}{n-k}$
#### Mean Squares

- Suppose we have a single population of 120 people which we divide randomly into...
  - 3 groups
  - 6 20 groups
- All else equal, which of these divisions do we expect to have higher SSG?

SSG = Variability Between Groups = 
$$\sum n_i (\bar{x}_i - \bar{x})^2$$

• We standardize sum of squares to compare SSG to SSE

Mean Variability Between Groups (MSG) = 
$$\frac{\text{SSG}}{k-1}$$
  
Mean Variability Within Groups (MSE) =  $\frac{\text{SSE}}{n-k}$ 

• Our goal is to use MSG and MSE to build a test statistic which measures when variability between groups is much greater than variability within groups

Mean Variability Between Groups 
$$=\frac{SSG}{k-1} = MSG$$
  
Mean Variability Within Groups  $=\frac{SSE}{n-k} = MSE$ 

Mean Variability Between Groups 
$$=\frac{SSG}{k-1} = MSG$$
  
Mean Variability Within Groups  $=\frac{SSE}{n-k} = MSE$ 

• The F statistic is

$$F = \frac{\text{MSG}}{\text{MSE}} = \frac{\frac{1}{k-1}\sum n_i(\bar{x}_i - \bar{x})^2}{\frac{1}{n-k}\sum (x_i - \bar{x}_i)^2}$$

Mean Variability Between Groups 
$$=\frac{SSG}{k-1} = MSG$$
  
Mean Variability Within Groups  $=\frac{SSE}{n-k} = MSE$ 

- The F statistic is  $F = \frac{\text{MSG}}{\text{MSE}} = \frac{\frac{1}{k-1}\sum n_i(\bar{x}_i \bar{x})^2}{\frac{1}{k-1}\sum (x_i \bar{x}_i)^2}$
- If all observations come from the same population, what is a typical value for F?

Mean Variability Between Groups 
$$=\frac{SSG}{k-1} = MSG$$
  
Mean Variability Within Groups  $=\frac{SSE}{n-k} = MSE$ 

- The F statistic is  $F = \frac{\text{MSG}}{\text{MSE}} = \frac{\frac{1}{k-1}\sum n_i(\bar{x}_i \bar{x})^2}{\frac{1}{k-1}\sum (x_i \bar{x}_i)^2}$
- If all observations come from the same population, what is a typical value for F?

Mean Variability Between Groups 
$$=\frac{SSG}{k-1} = MSG$$
  
Mean Variability Within Groups  $=\frac{SSE}{n-k} = MSE$ 

- The F statistic is  $F = \frac{\text{MSG}}{\text{MSE}} = \frac{\frac{1}{k-1}\sum n_i(\bar{x}_i \bar{x})^2}{\frac{1}{k-1}\sum (x_i \bar{x}_i)^2}$
- If all observations come from the same population, what is a typical value for F?

• If mean of groups are not equal, what values of F are typical?

Mean Variability Between Groups 
$$=\frac{SSG}{k-1} = MSG$$
  
Mean Variability Within Groups  $=\frac{SSE}{n-k} = MSE$ 

The F statistic is

$$F = \frac{\text{MSG}}{\text{MSE}} = \frac{\frac{1}{k-1}\sum n_i(\bar{x}_i - \bar{x})^2}{\frac{1}{n-k}\sum (x_i - \bar{x}_i)^2}$$

• If all observations come from the same population, what is a typical value for F?

• If mean of groups are not equal, what values of F are typical?

F > 1

Prof. Wells

#### F is for Films

• We could use the previous formulas to calculate the F statistic by hand...

#### F is for Films

- We could use the previous formulas to calculate the F statistic by hand...
  - But let's use technology!

#### F is for Films

- We could use the previous formulas to calculate the F statistic by hand...
  - But let's use technology!

```
movies_F<- movies %>%
   specify(AudienceScore ~ Genre) %>%
   calculate(stat = "F")
movies_F
```

```
## Response: AudienceScore (numeric)
## Explanatory: Genre (factor)
## # A tibble: 1 x 1
## stat
## <dbl>
## 1 4.34
```

• Is this a large value of F?

# Section 2

# The Distribution of the F statistic

- Hypotheses
  - Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
  - Alternative Hypothesis: At least one  $\mu$  is not equal to the others

- Hypotheses
  - Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
  - Alternative Hypothesis: At least one  $\mu$  is not equal to the others
- The F statistic

$$F = \frac{MSG}{MSE} = \frac{Variability Between Groups}{Variability Within Groups}$$

- Hypotheses
  - Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
  - Alternative Hypothesis: At least one  $\mu$  is not equal to the others
- The *F* statistic

$$\label{eq:F} \textit{F} = \frac{\text{MSG}}{\text{MSE}} = \frac{\text{Variability Between Groups}}{\text{Variability Within Groups}}$$

• Extreme values of the F statistic should give evidence that the null is not true

- Hypotheses
  - Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
  - Alternative Hypothesis: At least one  $\mu$  is not equal to the others
- The F statistic

$$F = \frac{\text{MSG}}{\text{MSE}} = \frac{\text{Variability Between Groups}}{\text{Variability Within Groups}}$$

- Extreme values of the F statistic should give evidence that the null is not true
  - How do we know which values of F are extreme?

- Hypotheses
  - Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
  - Alternative Hypothesis: At least one  $\mu$  is not equal to the others
- The F statistic

$$F = \frac{\text{MSG}}{\text{MSE}} = \frac{\text{Variability Between Groups}}{\text{Variability Within Groups}}$$

- Extreme values of the F statistic should give evidence that the null is not true
  - How do we know which values of F are extreme?
- We can find the distribution F under the null hypothesis by...

- Hypotheses
  - Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \cdots = \mu_7$
  - Alternative Hypothesis: At least one  $\mu$  is not equal to the others
- The F statistic

$$F = \frac{\text{MSG}}{\text{MSE}} = \frac{\text{Variability Between Groups}}{\text{Variability Within Groups}}$$

- Extreme values of the F statistic should give evidence that the null is not true
  - How do we know which values of F are extreme?
- We can find the distribution F under the null hypothesis by...
  - Randomization
  - Theoretical Approximation.

• If we assume that the quantitative and categorical variable are independent, then all samples are actually drawn from the same population

- If we assume that the quantitative and categorical variable are independent, then all samples are actually drawn from the same population
- We can imitate drawing new samples from this population by permuting the group labels among observations

- If we assume that the quantitative and categorical variable are independent, then all samples are actually drawn from the same population
- We can imitate drawing new samples from this population by permuting the group labels among observations
  - i.e. we assume that the genre label on a movie is superfluous and shuffle those labels around, while preserving Audience Rating.

- If we assume that the quantitative and categorical variable are independent, then all samples are actually drawn from the same population
- We can imitate drawing new samples from this population by permuting the group labels among observations
  - i.e. we assume that the genre label on a movie is superfluous and shuffle those labels around, while preserving Audience Rating.
- This way, we can study how the size of the *F* statistic changes just due to random sampling

The Distribution of the F statistic 00000000

#### Randomization and Permutation II

```
null_dist<-movies %>%
specify(AudienceScore ~ Genre) %>%
hypothesize(null = "independence") %>%
generate(reps = 1000, type = "permute" ) %>%
calculate(stat = "F")
null_dist %>% visualize()
```



- Most F statistics are at most 3
  - i.e. Assuming independence, Variance BETWEEN groups is at most 3 times variance WITHIN groups

The Distribution of the F statistic 00000000

#### Randomization and Permutation III

How does the observed F statistic compare?

How does the observed *F* statistic compare?

```
movies_F<-movies %>% specify(AudienceScore ~ Genre) %>% calculate(stat = "F") movies_F
```

## stat ## 1 4.340672

How does the observed F statistic compare?

```
movies_F<-movies %>% specify(AudienceScore ~ Genre) %>% calculate(stat = "F") movies_F
```

## stat ## 1 4.340672

```
null_dist %>% visualize()+shade_p_value(obs_stat = movies_F, direction = "right")
```



#### Simulation-Based Null Distribution

How does the observed F statistic compare?

```
movies_F<-movies %>% specify(AudienceScore ~ Genre) %>% calculate(stat = "F") movies_F
```

## stat ## 1 4.340672

```
null_dist %>% visualize()+shade_p_value(obs_stat = movies_F, direction = "right")
```



## p\_value ## 1 0.001

Like other statistics, the F statistic also has a theoretical distribution

• Suppose we have a total of *n* observations among *k* groups and that the following conditions hold:

- Suppose we have a total of *n* observations among *k* groups and that the following conditions hold:
  - 1 Observations are independent
  - **2** Within each group, values are approximately Normal
  - **③** Standard deviation is relatively constant between groups

- Suppose we have a total of *n* observations among *k* groups and that the following conditions hold:
  - 1 Observations are independent
  - **2** Within each group, values are approximately Normal
  - **③** Standard deviation is relatively constant between groups
- Then the distribution for the F statistic under the null hypothesis is well approximated by the F-distribution with parameters  $df_1 = k 1$  and  $df_2 = n k$ .

- Suppose we have a total of *n* observations among *k* groups and that the following conditions hold:
  - 1 Observations are independent
  - **2** Within each group, values are approximately Normal
  - **3** Standard deviation is relatively constant between groups
- Then the distribution for the F statistic under the null hypothesis is well approximated by the F-distribution with parameters  $df_1 = k 1$  and  $df_2 = n k$ .
  - The p-value is the area in the right tail.

Like other statistics, the F statistic also has a theoretical distribution

- Suppose we have a total of *n* observations among *k* groups and that the following conditions hold:
  - Observations are independent
  - **2** Within each group, values are approximately Normal
  - 3 Standard deviation is relatively constant between groups
- Then the distribution for the F statistic under the null hypothesis is well approximated by the F-distribution with parameters  $df_1 = k 1$  and  $df_2 = n k$ .
  - The p-value is the area in the right tail.

```
p_value<- pf(q = 4.340672, df1 = 6, df2 = 125, lower.tail = FALSE)
p_value</pre>
```

## [1] 0.0005161513

## Theory-based and Simulation-based Distributions

#### Simulation-Based and Theoretical F Null Distributions



#### There is No Accounting for Taste ... Even on Average

• The observed F statistic had P-value less than  $\alpha = 0.001$ 

#### There is No Accounting for Taste ... Even on Average

- The observed F statistic had P-value less than  $\alpha = 0.001$ 
  - That is, such a large *F*-stat would occur less than 0.01% of the time if the means of all groups were equal.

#### There is No Accounting for Taste ... Even on Average

- The observed F statistic had P-value less than  $\alpha = 0.001$ 
  - That is, such a large *F*-stat would occur less than 0.01% of the time if the means of all groups were equal.
  - This gives extremely good evidence against the Null hypothesis.
## There is No Accounting for Taste ... Even on Average

- The observed F statistic had P-value less than  $\alpha = 0.001$ 
  - That is, such a large *F*-stat would occur less than 0.01% of the time if the means of all groups were equal.
  - This gives extremely good evidence against the Null hypothesis.
  - We conclude that Audience Rating does depend on genre.