

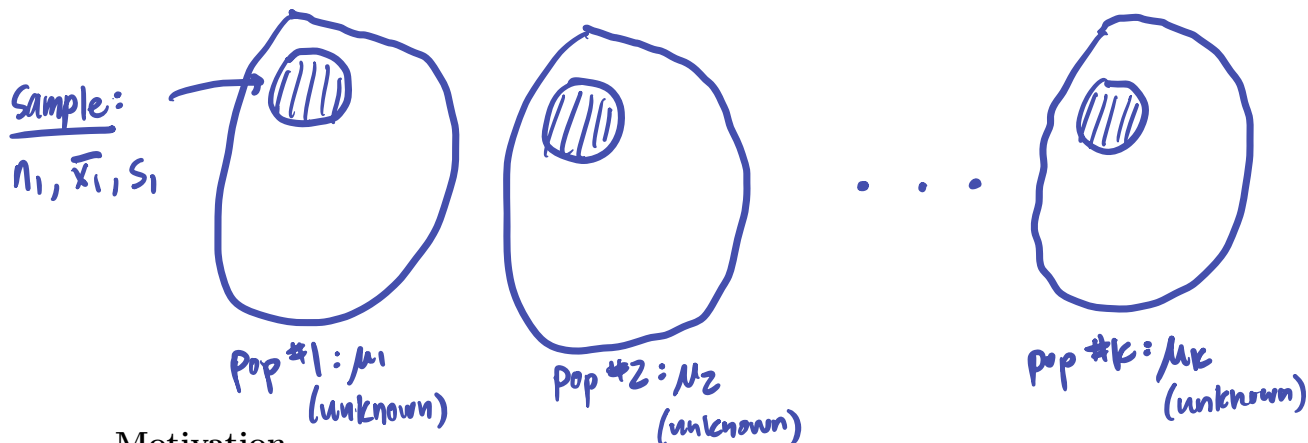
10 Analysis of Variance (ANOVA)

Chapter Preview

In chapter 9, we considered comparing the population means or population proportions of two populations. In this chapter, we will extend this framework to consider comparing the population means or proportions of three or more populations.

10.1: Single-Factor ANOVA

Illustration of ANOVA



Motivation

Q: Are the means the same in each population? Does $\mu_1 = \mu_2 = \dots = \mu_k$?

A: Collect a sample from each population and use ANOVA to make the determination.

Our hypotheses will be the following:

$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ (all pop. means are equal)

H_a : at least one μ_i differs

↳ this is not saying $\mu_1 \neq \mu_2 \neq \dots \neq \mu_k$

- Do the average GPAs differ by class level at this university?

Sample statistics:

First year:

$$n_1 = 5$$

$$\bar{x}_1 = 2.3$$

$$s_1 = 0.318$$

Second year:

$$n_2 = 5$$

$$\bar{x}_2 = 3.3$$

$$s_2 = 0.412$$

Third year:

$$n_3 = 5$$

$$\bar{x}_3 = 3.0$$

$$s_3 = 0.425$$

Fourth year:

$$n_4 = 5$$

$$\bar{x}_4 = 3.3$$

$$s_4 = 0.453$$

Is there enough evidence to suggest that the average GPA is not the same for all four levels?
Our hypotheses for this example will be:

$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$$

$$H_a: \text{at least one } \mu_i \text{ differs}$$

Recall that a two-sample t-test will compare the means of two populations (or treatments). ANOVA extends this to 3 or more populations (or treatments).

Definitions and Notation

Def: A **factor** is the characteristic that distinguishes the populations/treatments.

Def: A **level** of a factor is one of the specific populations or treatments.

Example: We want to test the average effectiveness of four drugs.

Factor: four drugs

Level: one of the four drugs (e.g., Ibuprofen)

Example: We want to compare the average growth of sunflowers for six growing conditions.

Factor: growing conditions

Level: one of the growing conditions (e.g., lots of irrigation, light fertilizer)

We will measure a quantitative response (Y) for a sample from each level of the factor.

Often, an ANOVA analysis has data that comes from a controlled experiment.

k = the number of levels/treatments/populations

Notation:

	True Population Parameters		Sample Statistics		
Group	Mean	Variance	Mean	Variance	Sample Size
Group #1	μ_1	σ_1^2	\bar{x}_1	s_1^2	n_1
Group #2	μ_2	σ_2^2	\bar{x}_2	s_2^2	n_2
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
Group #k	μ_k	σ_k^2	\bar{x}_k	s_k^2	n_k

$$N = n_1 + n_2 + \dots + n_k \quad (\text{total sample size})$$

$$T = n_1\bar{x}_1 + n_2\bar{x}_2 + \dots + n_k\bar{x}_k \quad (\text{sum of all observations})$$

$$\bar{\bar{x}} = T/N \quad (\text{overall/grand mean})$$

"x-double-bar"

Sum of Squares

Def: Sum of Squares (SS) quantify two types of variation.

- (1) Treatment Sum of Squares (SSTR) - measures the variability between groups.
How far apart are the sample means of the k populations?

$$\star \boxed{SSTR = n_1(\bar{x}_1 - \bar{\bar{x}})^2 + n_2(\bar{x}_2 - \bar{\bar{x}})^2 + \dots + n_k(\bar{x}_k - \bar{\bar{x}})^2}$$

$$\boxed{df_{TR} = k - 1} \quad (k = \# \text{ of groups})$$

- (2) Error Sum of Squares (SSE) - measures the variability within the k populations.
How spread out are the individual populations?

$$\star \boxed{SSE = (n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + \dots + (n_k - 1)s_k^2}$$

$$\boxed{df_E = N - k} \quad (N = \text{total } \# \text{ of observations})$$

- (3) Total Sum of Squares (SST) - measures the total variability in the data, ignoring populations and treatments.

"decomposition formula"

$$\star \boxed{SST = SSTR + SSE = \sum (x - \bar{\bar{x}})^2}$$

$$\boxed{df_T = N - 1}$$

Example: Calculate SSTR, SSE, and SST (and the associated degrees of freedom) for GPA data set.

$$k=4, N=20, \bar{\bar{x}} = \frac{T}{N} = \frac{1}{20}[5 \cdot 2.3 + 5 \cdot 3.3 + 5 \cdot 3.0 + 5 \cdot 3.3] = 2.975 = \bar{\bar{x}}$$

$$SSTR = \sum_{i=1}^k n_i (\bar{x}_i - \bar{\bar{x}})^2 = 5 \cdot (2.3 - 2.975)^2 + 5 \cdot (3.3 - 2.975)^2 + 5 \cdot (3.0 - 2.975)^2 + 5 \cdot (3.3 - 2.975)^2$$

$$\rightarrow \boxed{SSTR = 3.338, df_{TR} = k - 1 = 3}$$

$$SSE = \sum_{i=1}^k (n_i - 1) \cdot s_i^2 = 4 \cdot 0.318^2 + 4 \cdot 0.412^2 + 4 \cdot 0.425^2 + 4 \cdot 0.453^2$$

$$\rightarrow \boxed{SSE = 2.627, df_E = N - k = 16}$$

$$\boxed{SST = SSTR + SSE = 5.965}$$

$$df_T = N - 1 = 19$$

Note: ① If all pop. means are equal, $F_{test} \approx 1$.

② If one or more pop. means differ, $F_{test} \gg 1$.

Mean Squares ($MS = \frac{SS}{df}$)

Def: A mean square is the sum of squares divided by the degrees of freedom.

Measures "between" variability $\left\{ \begin{array}{l} MSTR = \frac{SSTR}{k-1} \end{array} \right.$

Measures "within" variability $\left\{ \begin{array}{l} MSE = \frac{SSE}{N-k} \end{array} \right.$

$F_{test} = \frac{MSTR}{MSE}$ "signal to noise ratio"

$df_1 = k-1, df_2 = N-k$

Note: Under the null hypothesis $H_0 (\mu_1 = \mu_2 = \dots \mu_k)$, $MSTR \approx MSE$.

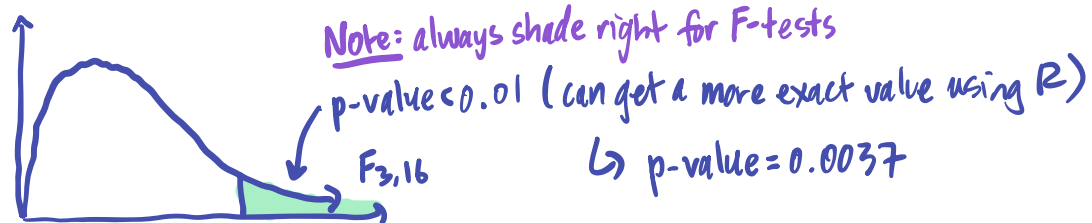
Note: Under the ANOVA framework, we are assuming the distributions of the k populations are approximately normally distributed and have equal variances.

- Calculate MSTR, MSE, F_{test} , and the p-value for the GPA data set. $k=4, N=20$

$$MSTR = \frac{SSTR}{df_{TR}} = \frac{3.338}{4-1} = 1.113 = MSTR (df_{TR}=3)$$

$$MSE = \frac{SSE}{df_E} = \frac{2.627}{20-4} = 0.164 = MSE (df_E=16)$$

$$F_{test} = \frac{MSTR}{MSE} = \frac{1.113}{0.164} = 6.79 = F_{test} (df_1=3, df_2=16)$$



Conclusion: p-value = 0.0037 $< 0.05 = \alpha \rightarrow$ Reject $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$

\rightarrow conclude that one or more of the class levels has a different true/pop. mean GPA.

ANOVA Assumptions

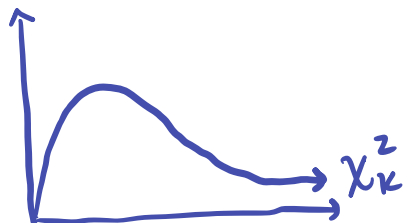
(or approx)

(pop.) ANOVA requires that the population or treatment distributions are **normal** and have the **same variance**. These assumptions can be checked in R although we assume these conditions hold in the problems we analyze.

↖ "chi-square"
 χ^2 distribution

↖ parameter of χ^2 -dist.

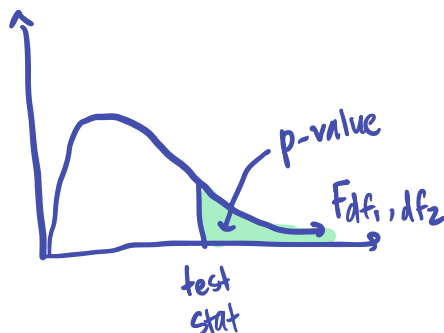
1. χ_k^2 is a continuous random variable, where **k** is the degrees of freedom.
2. A χ_k^2 random variable is the **sum of k independent squared standard normal random variables**.
3. Suppose Z_1, Z_2, \dots, Z_k are independent standard normal random variables. Then $\sum_{i=1}^k Z_i^2 \sim \chi_k^2$.
4. $\chi_k^2 \geq 0$
5. pdf of χ^2 distribution:



F distribution

1. F_{df_1, df_2} is a continuous random variable with two parameters, df_1 and df_2 , that correspond to the numerator and denominator degrees of freedom.
2. An F random variable is the ratio of two χ^2 distributions, so it will always be non-negative.
3. $F_{df_1, df_2} \sim \frac{\chi_{df_1}^2/df_1}{\chi_{df_2}^2/df_2}$.
4. For our purposes, any hypothesis test will be right-tailed.
5. pdf of F_{df_1, df_2} distribution:

(F-test)



$k = \#$ of groups/treatments
 $N =$ total sample size

One Way ANOVA F-test

Step 1: State the hypotheses.

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

H_a : at least one μ_i differs ($i=1, 2, \dots, k$)

Step 2: State the level of significance.

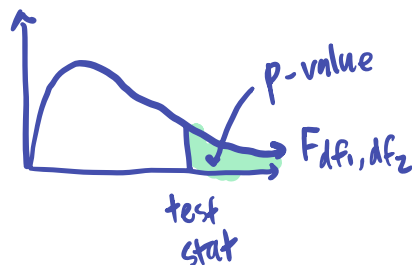
$$\alpha = \underline{\hspace{2cm}} \quad (\text{typically, } \alpha = 0.05)$$

Step 3: Calculate the test statistic.

$$F_{\text{test}} = \frac{MSTR}{MSE}, \quad \begin{matrix} df_1 = k-1 \\ df_2 = N-k \end{matrix}$$

Step 4: Calculate the p-value and plot.

↳ options: F-table, R, calculator



Step 5: Make a statistical decision.

Is $p\text{-value} \leq \alpha$? If yes, reject H_0 .
If no, FTR H_0 .

Step 6: Interpret your decision in the context of the problem.

$k = \# \text{ of groups/treatments/populations}$

Using R to complete ANOVA

Since many of these calculations can be tedious to calculate, a statistical program like R can be used to complete the analysis.

ANOVA table:

Source of Variation	df	SS	MS	F	p-value
Treatment	$k-1$	SSTR	MSTR	F_{test}	p-value
Error (Residual)	$N-k$	SSE	MSE	$\frac{MSTR}{MSE}$	
Total	$N-1$	SST			

find in R or F-table

R output:

partial ANOVA table

```
> summary(fit)
      Df Sum Sq Mean Sq F value Pr(>F)
year (group/term) 3  3.338  1.1125   6.781 0.00367 **
Residuals      16  2.625  0.1641
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Annotations: SSTR, MSE, MSTR, SSE, Df, Sum Sq, Mean Sq, F value, Pr(>F), p-value

R code:

GPA data set

```
fresh.gpa = c(2.05, 2.20, 2.00, 2.50, 2.75)
soph.gpa = c(3.00, 2.80, 3.80, 3.60, 3.30)
junior.gpa = c(3.30, 3.20, 2.30, 2.90, 3.30)
senior.gpa = c(3.50, 3.10, 2.90, 3.00, 4.00)
```

Make one list of all GPAs

```
gpa = c(fresh.gpa, soph.gpa, junior.gpa, senior.gpa)
```

Create labels

```
year = c(rep("fresh",5), rep("soph",5), rep("junior",5), rep("senior",5))
```

Make a boxplot

```
boxplot(gpa~year)
```

calculate ANOVA

```
fit = aov(gpa~year)
summary(fit)
```

ex) response: GPA
predictor: year

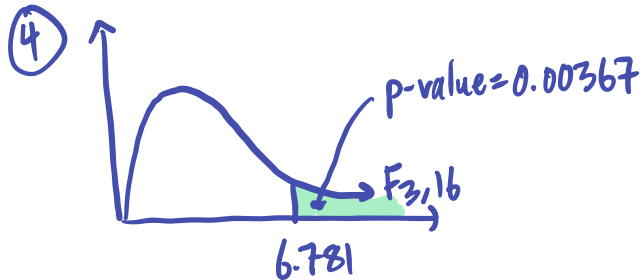
(ANOVA F-test)

- Using the R output from the previous page, complete a six-step hypothesis test to whether the mean GPAs at the four college levels are the same at $\alpha = 0.05$.

① $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ ② $\alpha = 0.05$

H_a : at least one μ_i differs

③ $F_{\text{test}} = 6.781$ ($df_1 = 3, df_2 = 16$)



⑤ p-value $\leq \alpha$? Yes, reject H_0 .

⑥ At $\alpha = 0.05$, there is sufficient evidence to conclude that one or more of the college levels have a different true mean GPA.

10.2: Multiple Comparison in ANOVA

Suppose that we conclude that the means of the k populations are not all equal. Which populations means do differ?

Motivation

Q: Can we just do two-sample t-tests for all pairs of populations?

A: Yes and no. We need to make an adjustment for simultaneously completing a group of tests.

- In the GPA example, we had four populations. Suppose the null hypothesis is true ($\mu_1 = \mu_2 = \mu_3 = \mu_4$) and $\alpha = 0.05$.

→ If we want to do all pairwise comparisons, there are $\binom{4}{2} = 6$ pairwise comparisons (e.g., compare 2nd and 4th year GPAs)

Q: What is the probability of zero Type I errors if the tests are independent?

A: $(0.95)^6 = 0.735$

→ since we are testing multiple hypotheses simultaneously, we need to make an adjustment to control the family-wise error rate.

Tukey's Procedure

Tukey's Procedure (Tukey's Honest Significant Difference or Tukey HSD) is a multiple comparisons procedure to account for completing a number of simultaneous tests.

Confidence intervals for the difference in means ($\mu_i - \mu_j$) can be computed using Tukey's Procedure. These calculations utilize another continuous distribution, the studentized range distribution (Q).

While these calculations can be completed by hand, we'll focus on calculating them in R.

Example: Returning to the GPA data set, determine which population means differ at $\alpha = 0.05$ using a Tukey adjustment. Some R code and output is provided below.

R code:

```
TukeyHSD(fit)
plot(TukeyHSD(fit))
```

R output:

```
> TukeyHSD(fit)
```

```
Tukey multiple comparisons of means
95% family-wise confidence level
```

```
Fit: aov(formula = gpa ~ year)
```

\$year	diff	lwr	upr	p adj
junior-fresh	0.7	-0.03291823	1.4329182	0.0638299
senior-fresh	1.0	0.26708177	1.7329182	0.0062342
soph-fresh	1.0	0.26708177	1.7329182	0.0062342
senior-junior	0.3	-0.43291823	1.0329182	0.6527735
soph-junior	0.3	-0.43291823	1.0329182	0.6527735
soph-senior	0.0	-0.73291823	0.7329182	1.0000000

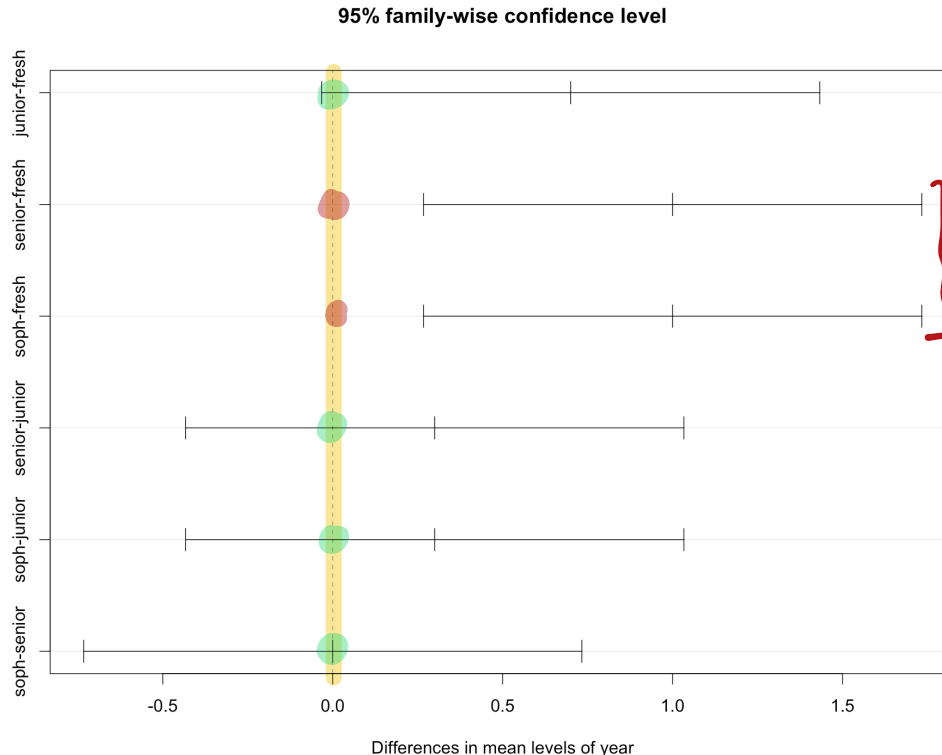
p-values (adjusted for multiple comparisons)

testing $H_0: \mu_1 = \mu_3$
(p-value = 0.064 > 0.05 = α ,
FTR H_0 : this difference is not
stat. sig.)

these diffs are
stat. sig.

Conclusion: 1st and 2nd year differ
1st and 4th year differ

R plot:



these CI do not
contain zero
↓
these diffs are
stat. sig.