

Stat 13, Intro. to Statistical Methods for the Life and Health Sciences.

1. ANOVA and F test continued.
2. Review list.
3. Examples.

Read ch9.

<http://www.stat.ucla.edu/~frederic/13/W23> .

The final is Fri in class and will be on ch 1-7, 10, and at most 1 question on ch9. Bring a PENCIL or pen and CALCULATOR and any books or notes you want. No computers.

If you cannot take it because of Covid or other health reasons, then email me to arrange a time to take an oral, in-person, 10-15 minute final exam in my office.

You can alternatively get an incomplete in the course and take the Spring stat 13 final, but don't come if you have covid.

ANOVA Output

- This is the kind of output you would see in most statistics packages when doing ANOVA.
- The variability between the groups is measured by the mean square treatment (40.02).
- The variability within the groups is measured by the mean square error (3.16).
- The F statistic is $40.02/3.16 = 12.67$.

Source	df	SS	MS	F	p-value
Treatment	2	80.04	40.02	12.67	0.0000
Error	54	170.53	3.16		
Total	56	250.56			

Conclusion

- Since we have a small p-value we have strong evidence against the null and can conclude at least one of the long-run mean recall scores is significantly different from the others.

Review list.

1. Meaning of SD.
2. Parameters and statistics.
3. Z statistic for proportions.
4. Simulation and meaning of pvalues.
5. SE for proportions.
6. What influences pvalues.
7. CLT and validity conditions for tests.
8. 1-sided and 2-sided tests.
9. Reject the null vs. accept the alternative.
10. Sampling and bias.
11. Significance level.
12. Type I, type II errors, and power.
13. CIs for a proportion.
14. CIs for a mean.
15. Margin of error.
16. Practical significance. (causation, extrapolation, curvature, heteroskedasticity).
17. Confounding.
18. Observational studies and experiments.
19. Random sampling and random assignment.
20. Two proportion CIs and testing.
21. IQR and 5 number summaries.
22. CIs for 2 means and testing.
23. Paired data.
24. Placebo effect, adherer bias, and nonresponse bias.
25. Prediction and causation.
26. Multiple testing and publication bias
27. Regression.
28. Correlation.
29. Calculate & interpret a & b.
30. Goodness of fit for regression.
31. Common regression problems
32. ANOVA and F-test.

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

1. What does the correlation of 0.82 imply?
 - a. 82% of the variation in weight is explained by height.
 - b. The typical variation in people's heights is 82% as large as the typical variation in their weights.
 - c. There is strong association between height and weight in this sample.
 - d. For every inch of increase in one's height, we would predict a 0.82 lb. increase in weight.
 - e. If a person weighs 100 pounds, then we typically would expect the person to be about 82 inches tall.

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

1. What does the correlation of 0.82 imply?

a. 82% of the variation in weight is explained by height.

b. The typical variation in people's heights is 82% as large as the typical variation in their weights.

c. There is strong association between height and weight in this sample.

d. For every inch of increase in one's height, we would predict a 0.82 lb. increase in weight.

e. If a person weighs 100 pounds, then we typically would expect the person to be about 82 inches tall.

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

2. What is the estimated slope, in lbs/inch, of the regression line for predicting weight from height?

a. 6.56. b. 7.12. c. 8.04. d. 9.92. e. 10.2. f. 11.4.

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

2. What is the estimated slope, in lbs/inch, of the regression line for predicting weight from height?

a. 6.56. b. 7.12. c. 8.04. d. 9.92. e. 10.2. f. 11.4.

$$r s_y / s_x = .82 \times 40 / 5 = 6.56.$$

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

3. How much would a prediction using this regression line typically be off by?

- a. 12.7 lbs. b. 13.5 lbs. c. 14.4lbs. d. 20.2 lbs. e. 22.9 lbs.

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

3. How much would a prediction using this regression line typically be off by?

- a. 12.7 lbs. b. 13.5 lbs. c. 14.4lbs. d. 20.2 lbs. e. **22.9 lbs.**

$$\sqrt{(1-r^2)} s_y = \sqrt{(1-.82^2)} \times 40 = 22.9.$$

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

4. If we were to randomly take one adult from this sample, how much would his/her height typically differ from 65 by?

a. 0.05 in. b. 0.1 in. c. 0.5 in. d. 1.0 in. e. 2.5 in. f. 5.0 in.

example problems.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

4. If we were to randomly take one adult from this sample, how much would his/her height typically differ from 65 by?

a. 0.05 in. b. 0.1 in. c. 0.5 in. d. 1.0 in. e. 2.5 in. **f. 5.0 in.**

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the median height is 64.5 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

5. Why shouldn't one trust this regression line to predict the weight of someone who is 25 inches tall?
- a. The sample size is insufficiently large.
 - b. The sample SD of weight is too small.
 - c. The value of 25 inches is too far outside the range of most observations.
 - d. The correlation of the ANOVA is a t-test confidence interval with statistical significance.
 - e. The data come from an observational study, so there may be confounding factors.
 - f. The height values are heavily right skewed, so the prediction errors are large.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the median height is 64.5 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs.

5. Why shouldn't one trust this regression line to predict the weight of someone who is 25 inches tall?
- a. The sample size is insufficiently large.
 - b. The sample SD of weight is too small.
 - c. The value of 25 inches is too far outside the range of most observations.**
 - d. The correlation of the ANOVA is a t-test confidence interval with statistical significance.
 - e. The data come from an observational study, so there may be confounding factors.
 - f. The height values are heavily right skewed, so the prediction errors are large.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the median height is 64.5 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs. The estimated slope in predicting weight from height is 6.56 lbs/in.

6. How should one interpret the estimated slope of 6.56?

- a. Each extra inch you grow causes you to increase your weight by 6.56 lbs on average.
- b. Each extra lb. you weigh causes you to grow 6.56 inches.
- c. The amount of weight Americans average is 6.56 standard errors above the mean.
- d. The Z-score corresponding to the correlation between height and weight is 6.56.
- e. For each extra inch taller you are, your predicted weight increases by 6.56 lbs.
- f. The proportion of variance in weight explained by the regression equation is 6.56%.

Suppose that among a sample of 100 adults in a given town, the correlation between height (inches) and weight (lbs.) is 0.82, and the mean height is 65 inches, the median height is 64.5 inches, the sd of height is 5 inches, the mean weight is 160 lbs., and the sd of weight is 40 lbs. The estimated slope in predicting weight from height is 6.56 lbs/in.

6. How should one interpret the estimated slope of 6.56?

- a. Each extra inch you grow causes you to increase your weight by 6.56 lbs on average.
- b. Each extra lb. you weigh causes you to grow 6.56 inches.
- c. The amount of weight Americans average is 6.56 standard errors above the mean.
- d. The Z-score corresponding to the correlation between height and weight is 6.56.
- e. For each extra inch taller you are, your predicted weight increases by 6.56 lbs.**
- f. The proportion of variance in weight explained by the regression equation is 6.56%.

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

7. What is the pooled sample percentage with acne, in both groups combined?
a. 16.1%. b. 17.2%. c. 18.4%. d. 19.7%. e. 21.2%. f. 23.0%.

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

7. What is the pooled sample percentage with acne, in both groups combined?
a. 16.1%. b. 17.2%. c. 18.4%. **d. 19.7%.** e. 21.2%. f. 23.0%.

$$118/600 = 19.7\%.$$

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

8. Using this pooled sample percentage, under the null hypothesis that the two groups have the same acne rate, what is the standard error for the difference between the two percentages?

- a. 1.42%. b. 1.88%. c. 2.02%. d. 2.99%. e. 3.08%. f. 3.44%.

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

8. Using this pooled sample percentage, under the null hypothesis that the two groups have the same acne rate, what is the standard error for the difference between the two percentages?

a. 1.42%. b. 1.88%. c. 2.02%. d. 2.99%. e. 3.08%. **f. 3.44%.**

$$\sqrt{(.197 * (1-.197)/400 + .197*(1-.197)/200)} = .0344.$$

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

9. What is the Z statistic for the difference between the two group percentages?

- a. 1.02. b. 1.23. c. 1.55. d. 1.88. e. 2.03. f. 2.43.

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

9. What is the Z statistic for the difference between the two group percentages?

- a. 1.02. b. 1.23. c. 1.55. d. 1.88. **e. 2.03.** f. 2.43.

$$(88/400 - 30/200) \div 0.0344 = 2.03.$$

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

10. Using the unpooled standard error for the difference between the two percentages, find a 95% confidence interval for the percentage with acne among those age 21-30 minus the percentage of acne among those age 31-40.

a. 7% +/- 5.02%. b. 7% +/- 5.54%. c. 7% +/- 5.92%. d. 7% +/- 6.03%. e. 7% +/- 6.41%.

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

10. Using the unpooled standard error for the difference between the two percentages, find a 95% confidence interval for the percentage with acne among those age 21-30 minus the percentage of acne among those age 31-40.

a. 7% +/- 5.02%. b. 7% +/- 5.54%. c. 7% +/- 5.92%. d. 7% +/- 6.03%. **e. 7% +/- 6.41%.**

$$\begin{aligned} & 88/400 - 30/200 \pm 1.96 * \sqrt{(88/400 * (1-88/400) / 400 + 30/200 * (1-30/200) / 200)} \\ & = 7\% \pm 6.41\%. \end{aligned}$$

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

11. What can we conclude from the 95% CI in the problem above?
- a. There is no statistically significant difference in the percentage with acne among those age 21-30 and those age 31-40.
 - b. There is statistically significant correlation between the sample size and the effect size for the confounding factors of a randomized controlled experiment.
 - c. Of those with acne, there is no statistically significant difference between the percentage who are age 21-30 and the percentage who are 31-40.
 - d. A statistically significantly higher percentage of people age 21-30 have acne than those age 31-40.
 - e. The sample sizes are too small for a conclusion to be valid here.

Suppose a researcher studying acne, or oily skin, takes a simple random sample of 400 Americans age 21-30 and calls them group A, and a simple random sample of 200 Americans age 31-40 and calls them group B. In group A, 88 people have acne, and in group B, 30 people have acne.

11. What can we conclude from the 95% CI in the problem above?

- a. There is no statistically significant difference in the percentage with acne among those age 21-30 and those age 31-40.
- b. There is statistically significant correlation between the sample size and the effect size for the confounding factors of a randomized controlled experiment.
- c. Of those with acne, there is no statistically significant difference between the percentage who are age 21-30 and the percentage who are 31-40.
- d. A statistically significantly higher percentage of people age 21-30 have acne than those age 31-40.**
- e. The sample sizes are too small for a conclusion to be valid here.

For the following two problems, suppose a researcher studies a treatment for a certain condition. The researcher divides 230 subjects randomly into three groups. Group A receives a placebo. Group B receives the treatment at a low dose, and group C receives the treatment at a high dose. The table below shows the output from an ANOVA F-test on the mean levels of the condition in the 3 groups.

Analysis of Variance Table

Response: outcome

	Df	Sum Sq	MeanSq	F value	Pr(>F)
group	2	0.467	0.233	0.221	0.803
Residuals	227	28.5	1.06		

12. In the table, what number is a measure of the variability between groups?

a. 2. b. 0.233 c. 0.221 d. 0.803 e. 28.5 f. 1.06.

For the following two problems, suppose a researcher studies a treatment for a certain condition. The researcher divides 230 subjects randomly into three groups. Group A receives a placebo. Group B receives the treatment at a low dose, and group C receives the treatment at a high dose. The table below shows the output from an ANOVA F-test on the mean levels of the condition in the 3 groups.

Analysis of Variance Table

Response: outcome

	Df	Sum Sq	MeanSq	F value	Pr(>F)
group	2	0.467	0.233	0.221	0.803
Residuals	227	28.5	1.06		

12. In the table, what number is a measure of the variability between groups?

a. 2. **b. 0.233** c. 0.221 d. 0.803 e. 28.5 f. 1.06.

For the following two problems, suppose a researcher studies a treatment for a certain condition. The researcher divides 230 subjects randomly into three groups. Group A receives a placebo. Group B receives the treatment at a low dose, and group C receives the treatment at a high dose. The table below shows the output from an ANOVA F-test on the mean levels of the condition in the 3 groups.

Analysis of Variance Table

Response: outcome

	Df	Sum Sq	MeanSq	F value	Pr(>F)
group	2	0.467	0.233	0.221	0.803
Residuals	227	28.5	1.06		

13. What can we conclude from the result of the F test?

- a. The treatment has a statistically significant effect, though we cannot be sure which dose is attributable to this significant effect or in which direction the effect goes.
- b. The treatment seems to have significantly greater effect at large doses than at small doses.
- c. The treatment does not seem to have a statistically significant effect.
- d. We fail to reject the null hypothesis that the treatment has a significant effect on the condition.
- e. The correlation between dose and outcome appears to be statistically significant.

For the following two problems, suppose a researcher studies a treatment for a certain condition. The researcher divides 230 subjects randomly into three groups. Group A receives a placebo. Group B receives the treatment at a low dose, and group C receives the treatment at a high dose. The table below shows the output from an ANOVA F-test on the mean levels of the condition in the 3 groups.

Analysis of Variance Table

Response: outcome

	Df	Sum Sq	MeanSq	F value	Pr(>F)
group	2	0.467	0.233	0.221	0.803
Residuals	227	28.5	1.06		

13. What can we conclude from the result of the F test?

- a. The treatment has a statistically significant effect, though we cannot be sure which dose is attributable to this significant effect or in which direction the effect goes.
- b. The treatment seems to have significantly greater effect at large doses than at small doses.
- c. The treatment does not seem to have a statistically significant effect.**
- d. We fail to reject the null hypothesis that the treatment has a significant effect on the condition.
- e. The correlation between dose and outcome appears to be statistically significant.