

R Data Analysis Course – Week 7

Group Comparisons

Statistical Analysis

Ronnie Steinitz, Ph.D.
Dian Fossey Gorilla Fund

Photo: R. Steinitz | Bwindi, Uganda, 2022

Resources to reinforce your understanding

The Google Drive folder includes:

- Every week: a **complete .Rmd file** with detailed in-line explanations and example code from the lecture
- General: an updated “**Functions Learned So Far**” reference sheet

Every week: extra support with your code, project setup, or any course concepts

Group A



Group B



Group Comparison

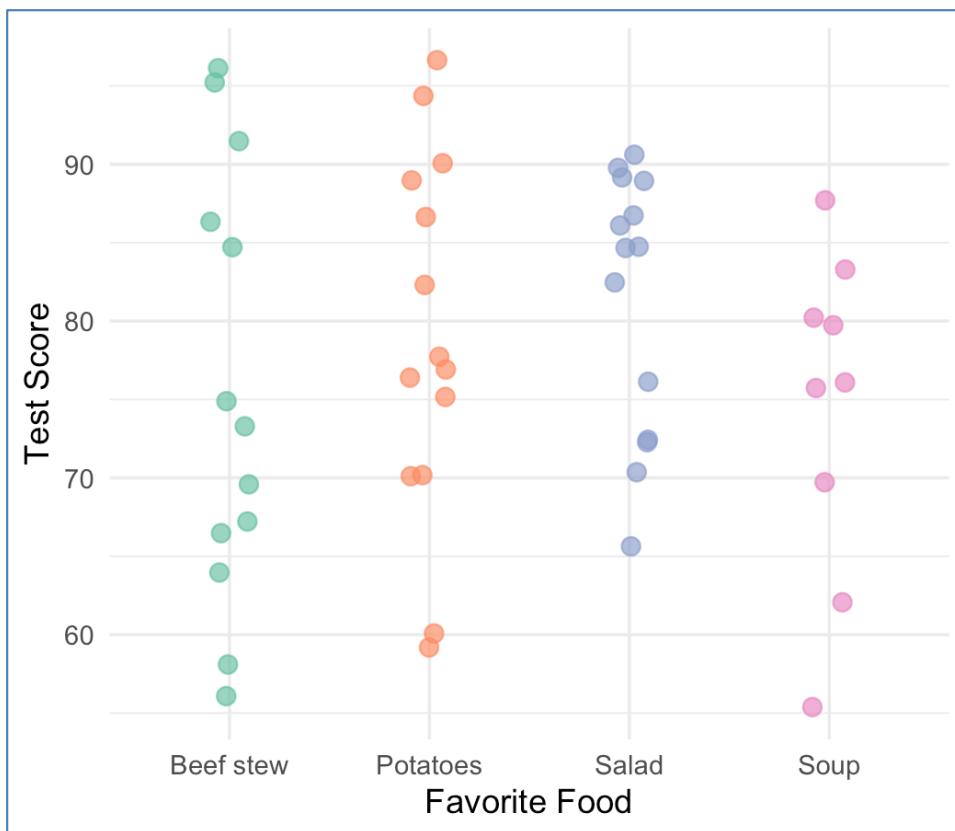
- Are differences in mean value due to real group differences or random variation?
- The larger difference between groups
... the more likely the groups truly differ

Group Comparison

id	favorite_food	studied	test_score
1	Potatoes	Yes	86.6
2	Potatoes	Yes	77.7
3	Potatoes	Yes	96.6
4	Potatoes	No	70.1
5	Salad	Yes	84.7
6	Beef stew	No	67.2
7	Salad	No	70.4
8	Salad	Yes	84.7
9	Potatoes	Yes	60.1
10	Beef stew	Yes	86.3
11	Soup	No	62.1
12	Beef stew	No	66.5
13	Soup	No	69.7
14	Beef stew	Yes	95.2
15	Potatoes	No	59.2
16	Potatoes	Yes	94.4
17	Salad	Yes	86.7
18	Beef stew	No	73.3
19	Salad	No	72.4
20	Salad	Yes	89.8
21	Soup	Yes	75.7
22	Soup	Yes	83.3
23	Potatoes	Yes	89
24	Potatoes	Yes	76.4
25	Soup	Yes	79.7

Group Comparison

Raw Values

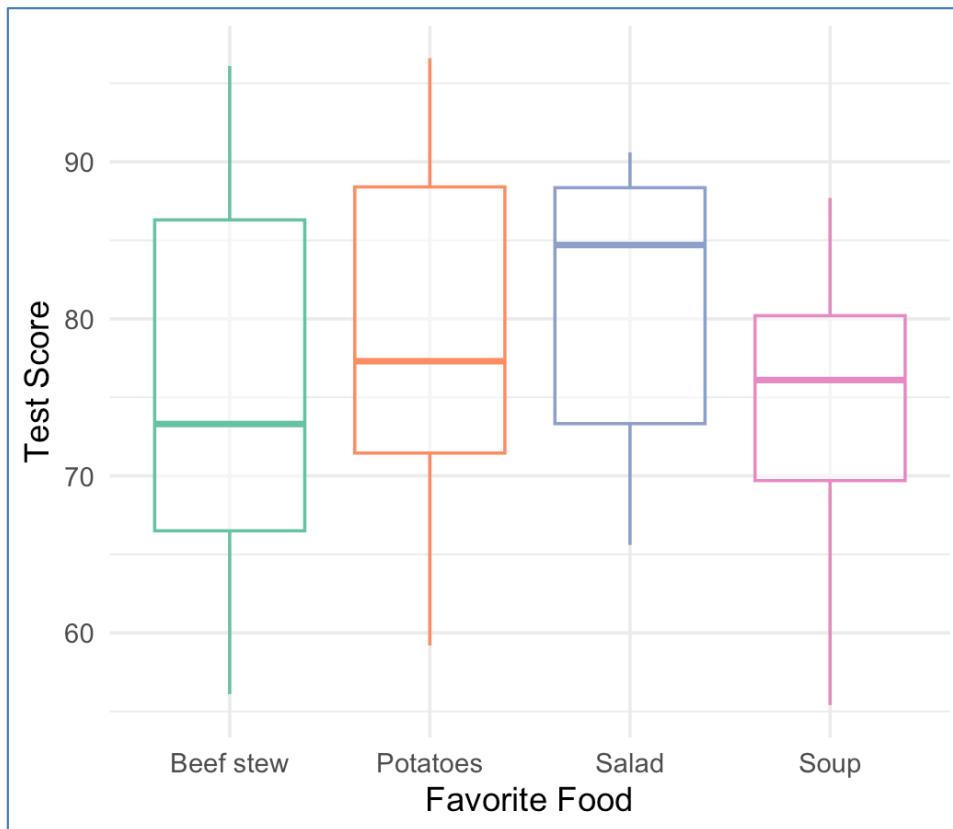


Test Score ~ Favorite Food

id	favorite_food	studied	test_score
1	Potatoes	Yes	86.6
2	Potatoes	Yes	77.7
3	Potatoes	Yes	96.6
4	Potatoes	No	70.1
5	Salad	Yes	84.7
6	Beef stew	No	67.2
7	Salad	No	70.4
8	Salad	Yes	84.7
9	Potatoes	Yes	60.1
10	Beef stew	Yes	86.3
11	Soup	No	62.1
12	Beef stew	No	66.5
13	Soup	No	69.7
14	Beef stew	Yes	95.2
15	Potatoes	No	59.2
16	Potatoes	Yes	94.4
17	Salad	Yes	86.7
18	Beef stew	No	73.3
19	Salad	No	72.4
20	Salad	Yes	89.8
21	Soup	Yes	75.7
22	Soup	Yes	83.3
23	Potatoes	Yes	89
24	Potatoes	Yes	76.4
25	Soup	Yes	79.7

Group Comparison

Boxplot – summary stats

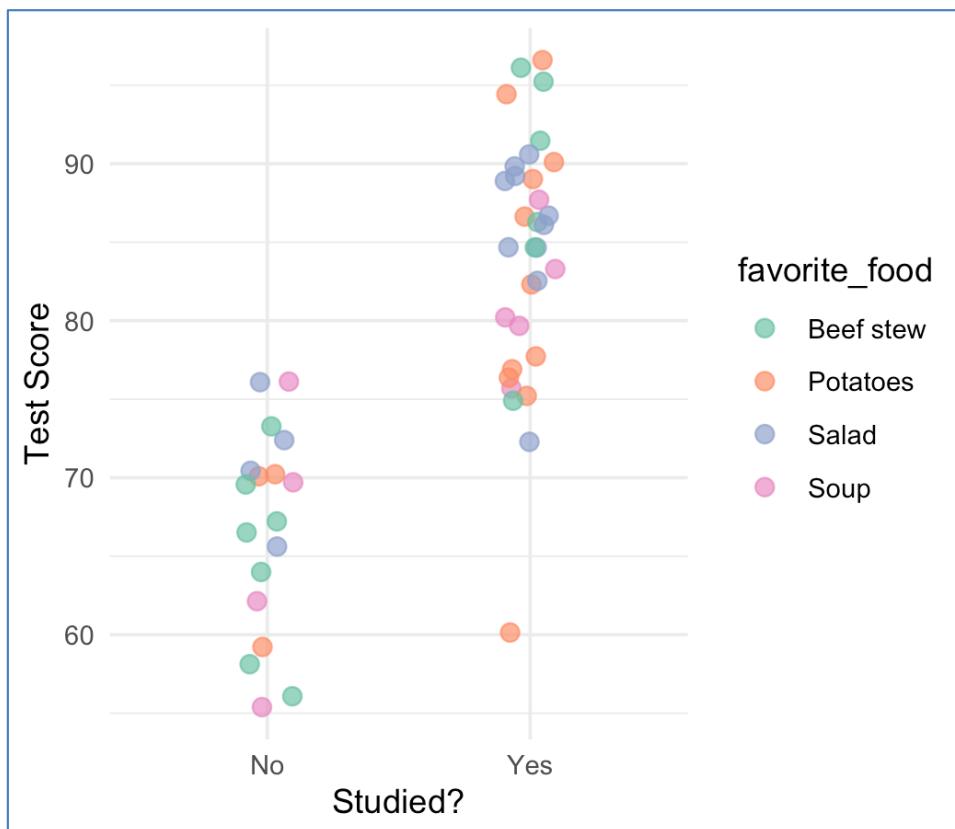


Test Score ~ Favorite Food

id	favorite_food	studied	test_score
1	Potatoes	Yes	86.6
2	Potatoes	Yes	77.7
3	Potatoes	Yes	96.6
4	Potatoes	No	70.1
5	Salad	Yes	84.7
6	Beef stew	No	67.2
7	Salad	No	70.4
8	Salad	Yes	84.7
9	Potatoes	Yes	60.1
10	Beef stew	Yes	86.3
11	Soup	No	62.1
12	Beef stew	No	66.5
13	Soup	No	69.7
14	Beef stew	Yes	95.2
15	Potatoes	No	59.2
16	Potatoes	Yes	94.4
17	Salad	Yes	86.7
18	Beef stew	No	73.3
19	Salad	No	72.4
20	Salad	Yes	89.8
21	Soup	Yes	75.7
22	Soup	Yes	83.3
23	Potatoes	Yes	89
24	Potatoes	Yes	76.4
25	Soup	Yes	79.7

T-test

Raw Values



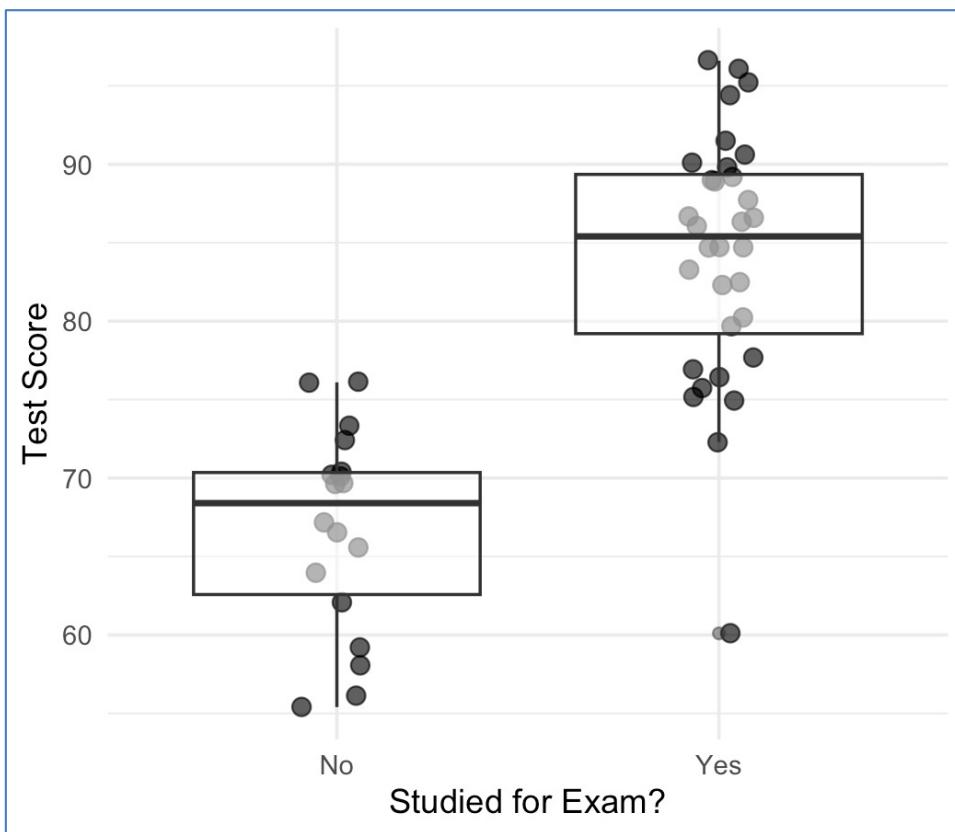
Test Score ~ Studied?

What about favorite food?

id	favorite_food	studied	test_score
1	Potatoes	Yes	86.6
2	Potatoes	Yes	77.7
3	Potatoes	Yes	96.6
4	Potatoes	No	70.1
5	Salad	Yes	84.7
6	Beef stew	No	67.2
7	Salad	No	70.4
8	Salad	Yes	84.7
9	Potatoes	Yes	60.1
10	Beef stew	Yes	86.3
11	Soup	No	62.1
12	Beef stew	No	66.5
13	Soup	No	69.7
14	Beef stew	Yes	95.2
15	Potatoes	No	59.2
16	Potatoes	Yes	94.4
17	Salad	Yes	86.7
18	Beef stew	No	73.3
19	Salad	No	72.4
20	Salad	Yes	89.8
21	Soup	Yes	75.7
22	Soup	Yes	83.3
23	Potatoes	Yes	89
24	Potatoes	Yes	76.4
25	Soup	Yes	79.7

T-test

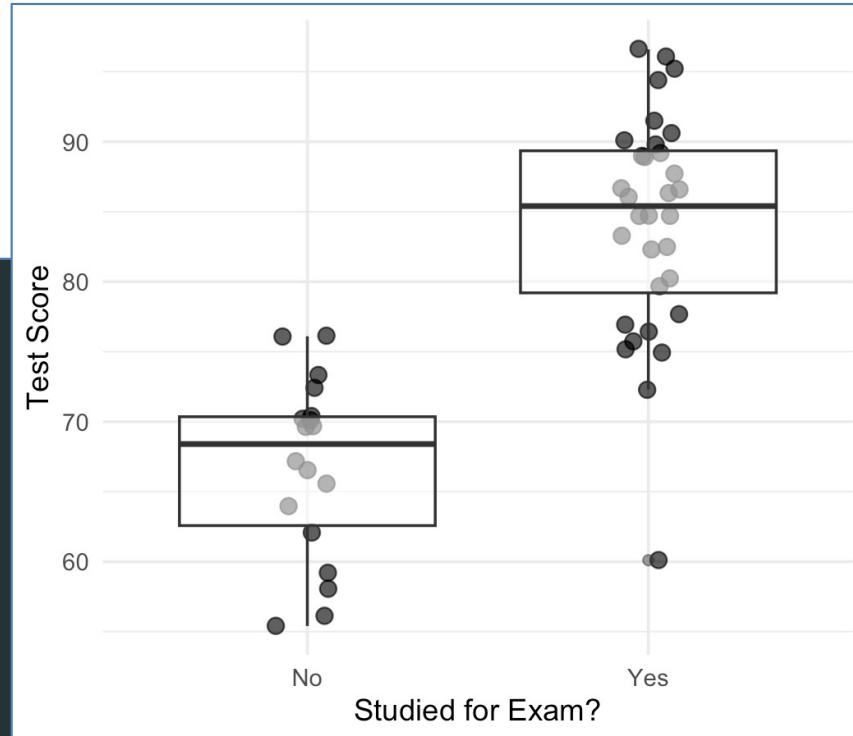
Raw values and stats



Test Score ~ Studied?

id	favorite_food	studied	test_score
1	Potatoes	Yes	86.6
2	Potatoes	Yes	77.7
3	Potatoes	Yes	96.6
4	Potatoes	No	70.1
5	Salad	Yes	84.7
6	Beef stew	No	67.2
7	Salad	No	70.4
8	Salad	Yes	84.7
9	Potatoes	Yes	60.1
10	Beef stew	Yes	86.3
11	Soup	No	62.1
12	Beef stew	No	66.5
13	Soup	No	69.7
14	Beef stew	Yes	95.2
15	Potatoes	No	59.2
16	Potatoes	Yes	94.4
17	Salad	Yes	86.7
18	Beef stew	No	73.3
19	Salad	No	72.4
20	Salad	Yes	89.8
21	Soup	Yes	75.7
22	Soup	Yes	83.3
23	Potatoes	Yes	89
24	Potatoes	Yes	76.4
25	Soup	Yes	79.7

T-test

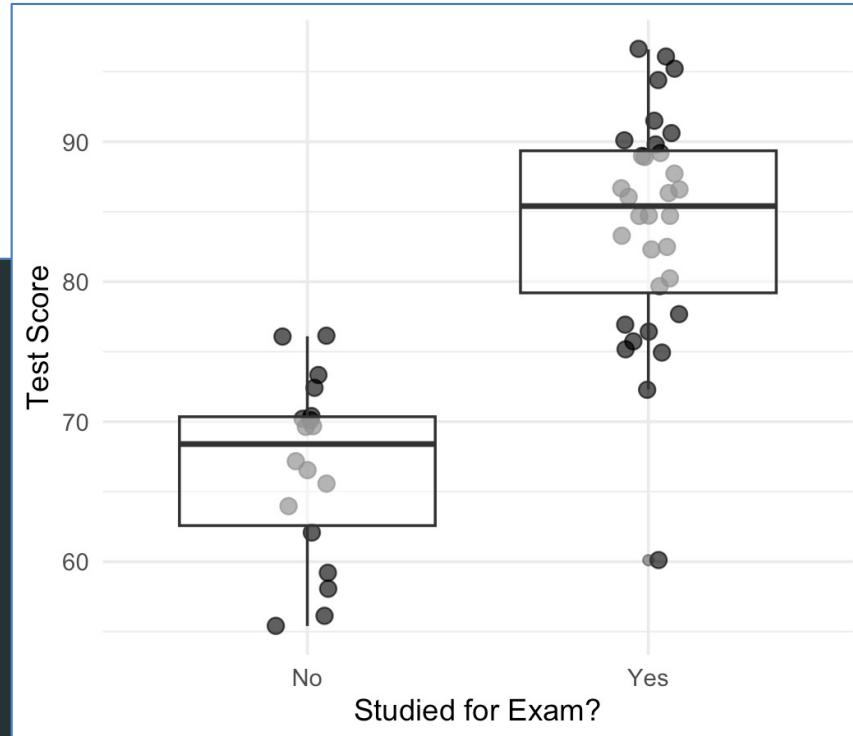


```
> t.studied <- t.test(test_score ~ studied, data = mock_data)  
> t.studied
```

Welch Two Sample t-test

```
data: test_score by studied  
t = -8.4496, df = 41.64, p-value = 0.0000000001426  
alternative hypothesis: true difference in means between group No and group Yes is not equal to 0  
95 percent confidence interval:  
-21.64329 -13.29629  
sample estimates:  
mean in group No mean in group Yes  
66.78333 84.25312
```

T-test



```
> t.studied <- t.test(test_score ~ studied, data = mock_data)
```

t = difference **df** = ~total

relative to data points in
variance of data the study

```
Welch Two Sample t-test
```

data: test score by studied
t = -8.4496, **df** = 41.64, **p-value** = 0.0000000001426

alternative hypothesis: true difference in means between group No and group Yes is not equal to 0

95 percent confidence interval:
-21.64329 -13.29629

sample estimates:

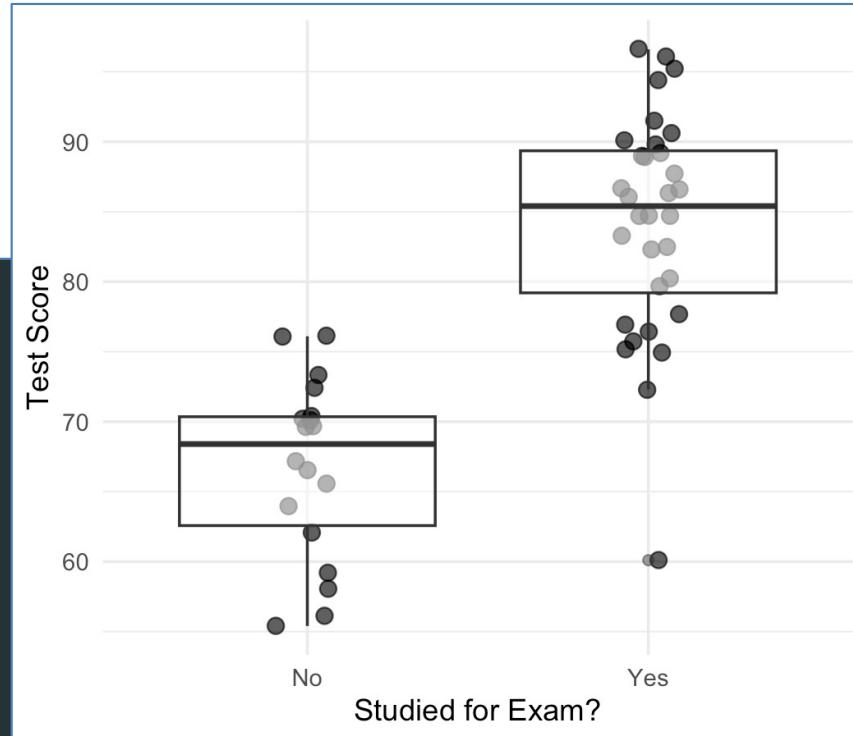
mean in group No	mean in group Yes
66.78333	84.25312

p = likelihood that
difference between
groups is random

confidence interval = probably range of true difference

Mean of each group

T-test

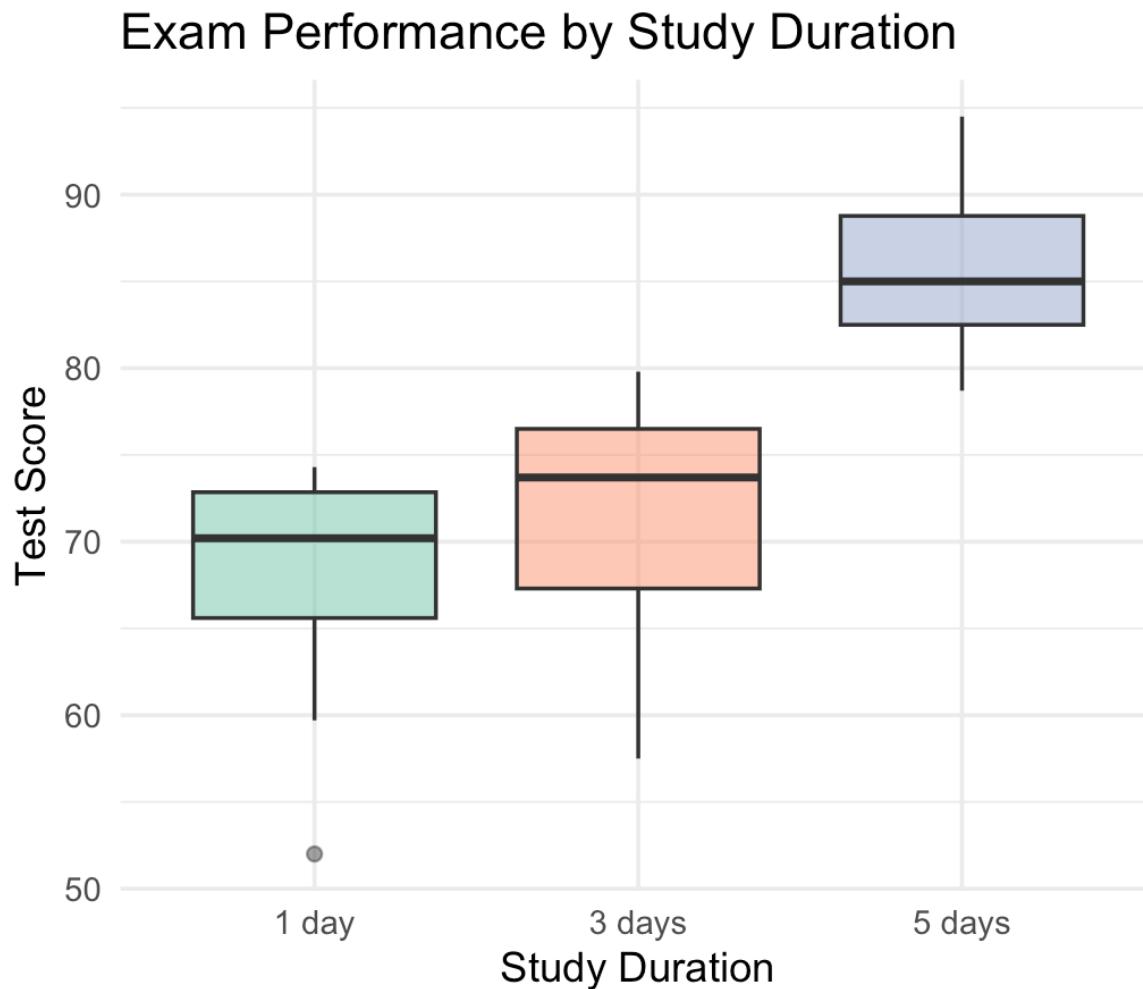


Students who studied scored significantly higher on the test than those who did not ($t(41.6) = -8.45, p < 0.001$).

WELCH TWO-SAMPLE T-TEST

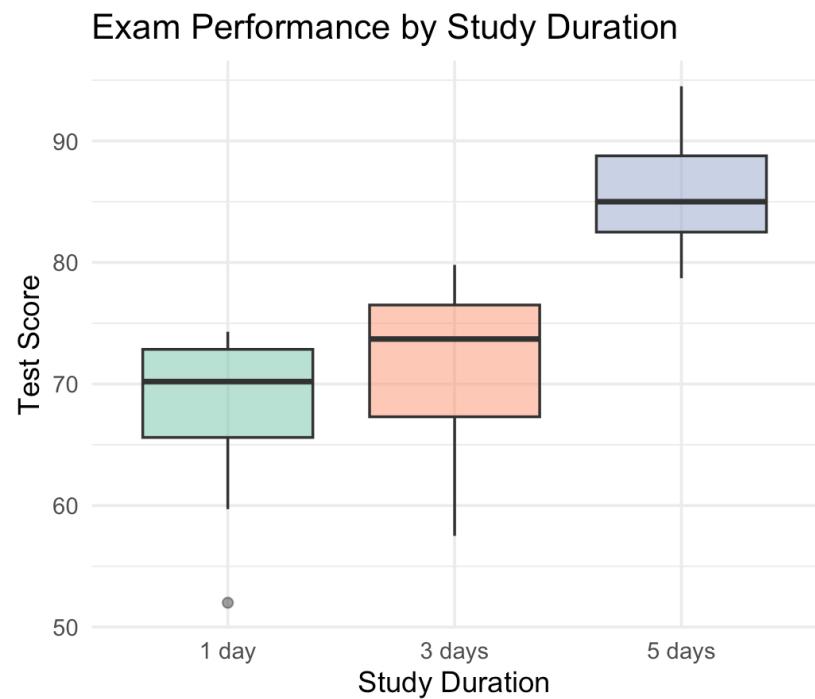
```
data: test_score by studied
t = -8.4496, df = 41.64, p-value = 0.0000000001426
alternative hypothesis: true difference in means between group No and group Yes is not equal to 0
95 percent confidence interval:
-21.64329 -13.29629
sample estimates:
mean in group No mean in group Yes
66.78333      84.25312
```

More than 2 groups: ANOVA = Analysis of Variance



More than 2 groups: ANOVA = Analysis of Variance

```
> anova_studied <- aov(test_score ~ studied_qty, data = mock_study)
> summary(anova_studied)
   Df Sum Sq Mean Sq F value    Pr(>F)
studied_qty  2  3066   1533.0   44.99 0.000000000000121 ***
Residuals  47  1601     34.1
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



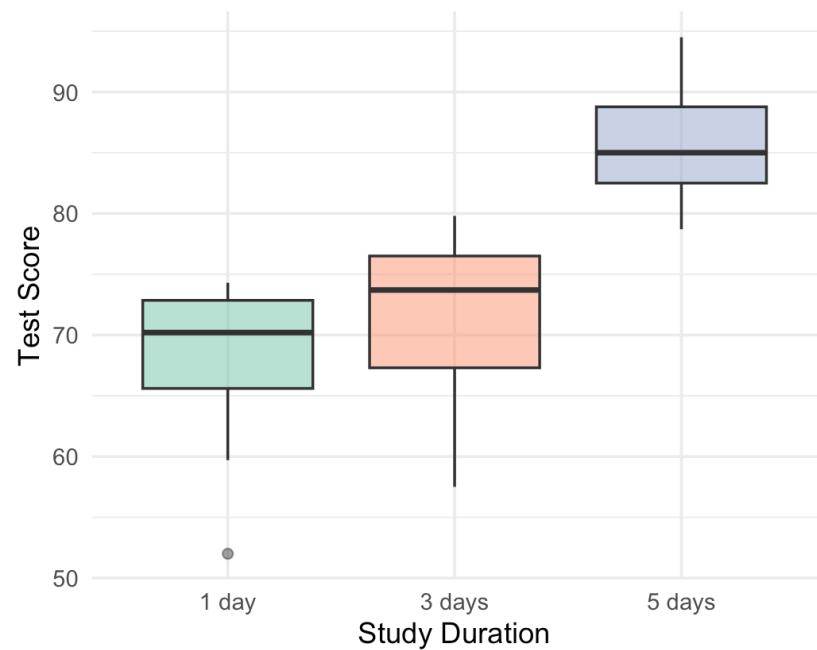
More than 2 groups: ANOVA = Analysis of Variance

```
> anova_studied <- aov(test_score ~ studied_qty, data = mock_study)
> summary(anova_studied)
  Df Sum Sq Mean Sq F value    Pr(>F)
studied_qty  2  3066   1533.0 44.99 0.0000000000121 ***
Residuals   47  1601     34.1
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

F = ratio of variation p-value

*between groups vs. within
groups*

Exam Performance by Study Duration

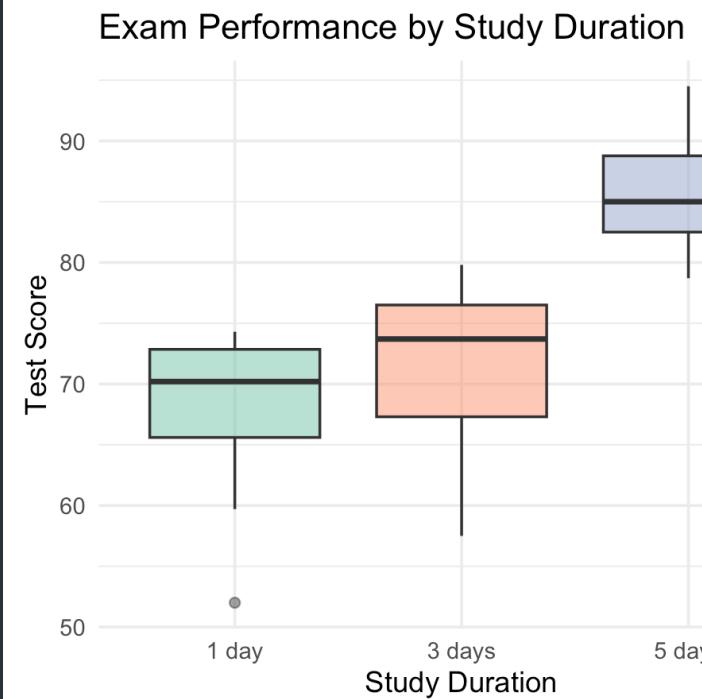


More than 2 groups: ANOVA = Analysis of Variance

```
> anova_studied <- aov(test_score ~ studied_qty, data = mock_study)
> summary(anova_studied)
   Df Sum Sq Mean Sq F value    Pr(>F)
studied_qty  2  3066   1533.0   44.99 0.0000000000121 ***
Residuals  47  1601     34.1
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> TukeyHSD(anova_studied)
  Tukey multiple comparisons of means
    95% family-wise confidence level

Fit: aov(formula = test_score ~ studied_qty, data = mock_study)

$studied_qty
      diff      lwr      upr      p adj
3 days-1 day  3.801215 -1.283410  8.885839 0.1776357
5 days-1 day 17.558480 12.912041 22.204918 0.0000000
5 days-3 days 13.757265  8.615582 18.898948 0.0000002
```



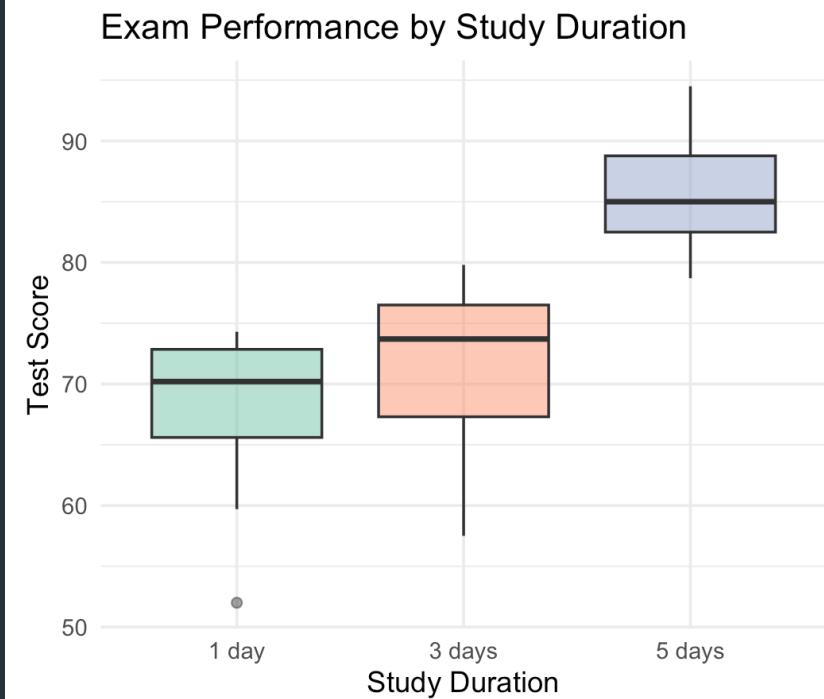
More than 2 groups: ANOVA = Analysis of Variance

```
> anova_studied <- aov(test_score ~ studied_qty, data = mock_study)
> summary(anova_studied)
   Df Sum Sq Mean Sq F value    Pr(>F)
studied_qty  2  3066   1533.0   44.99 0.0000000000121 ***
Residuals  47  1601     34.1
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> TukeyHSD(anova_studied)
  Tukey multiple comparisons of means
    95% family-wise confidence level

Fit: aov(formula = test_score ~ studied_qty, data = mock_study)

$studied_qty
      diff      lwr      upr   p adj
3 days-1 day 3.801215 -1.283410 8.885839 0.1776357
5 days-1 day 17.558480 12.912041 22.204918 0.0000000
5 days-3 days 13.757265  8.615582 18.898948 0.0000002
```

*Difference
between
mean* *Upper and lower
boundaries of
confidence interval*



Study duration had a significant effect on test scores ($F(2, 47) = 44.99$, $p < 0.001$). Post-hoc comparisons showed that students who studied for 5 days scored significantly higher than those who studied for 1 or 3 days, while the 1- and 3-day groups did not differ.

```
> anova_studied <- aov(test_score ~ studied_qty, data = mock_study)
> summary(anova_studied)
   Df Sum Sq Mean Sq F value    Pr(>F)
studied_qty  2  3066   1533.0   44.99 0.0000000000121 ***
Residuals  47  1601     34.1
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
> TukeyHSD(anova_studied)
  Tukey multiple comparisons of means
  95% family-wise confidence level
```

```
Fit: aov(formula = test_score ~ studied_qty, data = mock_study)
```

\$studied_qty	diff	lwr	upr	p adj
3 days-1 day	3.801215	-1.283410	8.885839	0.1776357
5 days-1 day	17.558480	12.912041	22.204918	0.0000000
5 days-3 days	13.757265	8.615582	18.898948	0.0000002

Difference between mean **Upper and lower boundaries of confidence interval**

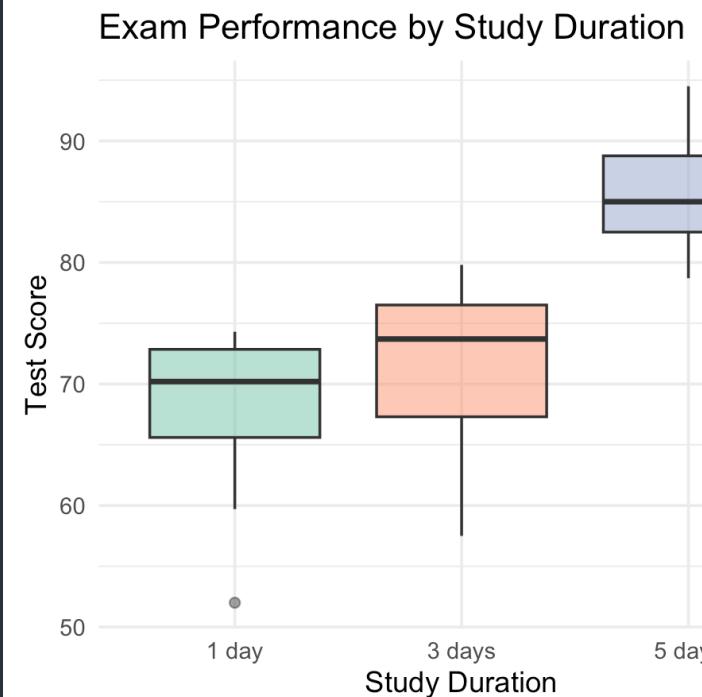


Table S4.1. Age-sex Class Differences in Biomarker Values. Results from a series of ANOVA tests followed by post-hoc Welch's t-tests to identify pair-wise differences for each biomarker between age-sex classes (AF = adult female; AM = adult male; SF = subadult females; SM = subadult male) of red-tailed monkeys. SF were excluded from $\delta^{15}\text{N}$ tests because no isotope data exists for this age-sex class.

Biomarker	Classes	N₁	N₂	t	df	p-value	adj. p-value	ns
CR:SG	AF - AM	1146	122	-2.53	152	0.01	0.06	ns
	AF - SF	1146	20	-2.94	19.7	0.01	0.05	
	AF - SM	1146	54	-1.27	59.9	0.21	0.42	
	AM - SF	122	20	-1.81	24.8	0.08	0.25	
	AM - SM	122	54	0.50	109	0.62	0.62	
	SF - SM	20	54	2.00	30.4	0.05	0.22	
UCP	AF - AM	1676	182	-0.13	211	0.90	1.00	ns
	AF - SF	1676	43	6.65	82.5	< 0.001	< 0.001	
	AF - SM	1676	114	0.53	128	0.60	1.00	
	AM - SF	182	43	3.31	222	0.001	0.01	
	AM - SM	182	114	0.50	265	0.62	1.00	
	SF - SM	43	114	-2.44	143	0.02	0.06	
T3	AF - AM	235	12	0.86	11.8	0.41	1.00	ns
	AF - SF	235	8	-1.45	7.47	0.19	0.94	
	AF - SM	235	16	-0.78	15.5	0.45	1.00	
	AM - SF	12	8	-1.67	16.7	0.11	0.68	
	AM - SM	12	16	-1.14	24.8	0.26	1.00	
	SF - SM	8	16	0.21	22	0.84	1.00	

Table S4.1. Age-sex Class Differences in Biomarker Values. Results from a series of ANOVA tests followed by post-hoc Welch's t-tests to identify pair-wise differences for each biomarker between age-sex classes (AF = adult female; AM = adult male; SF = subadult females; SM = subadult male) of red-tailed monkeys. SF were excluded from $\delta^{15}\text{N}$ tests because no isotope data exists for this age-sex class.

Biomarker	Classes	N₁	N₂	t	df	p-value	p-value	adj.
CR:SG	AF - AM	1146	122	-2.53	152	0.01	0.06	ns
	AF - SF	1146	20	-2.94	19.7	0.01	0.05	*
	AF - SM	1146	54	-1.27	59.9	0.21	0.42	ns
	AM - SF	122	20	-1.81	24.8	0.08	0.25	ns
	AM - SM	122	54	0.50	109	0.62	0.62	ns
	SF - SM	20	54	2.00	30.4	0.05	0.22	ns

Adult females had significantly lower relative muscle mass compared to subadult females (CR:SG: $t(19.7) = -2.94$, $p.\text{adj} < 0.05$; Table S4.1).

To evaluate the differences

in $\delta^{15}\text{N}$ isotope values across different reproductive statuses (cycling, noncycling, and pregnant) in red-tailed monkeys, a one-way ANOVA was conducted, followed by pairwise comparisons using Welch's t-tests (Derrick et al., 2016).

Adult females had significantly lower relative muscle mass compared to subadult females (CR:SG: $t(19.7) = -2.94$, $p.\text{adj} < 0.05$; Table S4.1). No significant differences were observed for other comparisons across age-sex classes (all $p > 0.05$). Subadult females had lower UCP values compared to adult females ($t(82.5) = 6.65$, $p.\text{adj} < 0.001$), and compared to adult males ($t(222) = 3.31$, $p.\text{adj} < 0.01$), but other comparisons of UCP did not show significant differences (all $p > 0.06$). All other pairwise comparisons of hormones across classes were non-significant ($p > 0.18$ for T3, $p > 0.10$ for $\delta^{15}\text{N}$, $p > 0.14$ for ketones).

Skills Learning

code-along lecture



Mount Sabyinyo 2025

Skills Application

Laboratory exercise

R Data Anlaysis Course: Final Project Instructions

Ronnie Bailey-Steinitz

2025-11-03

Week 8: Independent Data Exploration and Analysis

Self-Guided Exercise

This week, you will apply all the skills you've learned so far to analyze new data.

Your goal is to create a *clear, reproducible* R Markdown report that walks through your full analysis: **from data import to final visualization.**

You'll share one plot from your analysis in next week's *Presentation Day*.

1. Create a New R Markdown File

- Open RStudio and create a new .Rmd file called Week8_YourName_DataAnalysis.Rmd.
- Save it inside your course project folder.
- Add a title, author, and date in the header.
- Use headings (#, ##) to organize your markdown file into sections.

2. Load Required Libraries

- Determine which libraries you need and load them.

3. Import Your Dataset

Week 7 - Skills Application

Download Final Project Instructions from [Google Drive folder](#)

If you are working on your own dataset:

- Keep working on your *Lastname_Firstname_Data.Rmd*

If you are working on sample dataset:

- create a Markdown file just for this week.
- Name it: *Lastname_Firstname_final.Rmd* in `/**Week 7**`

In your new Markdown file: Use code blocks (Ctrl+Alt+i) to load packages, import dataset into the environment, and save data as an object (e.g., `*data <-* `). Then follow *instructions*.