

# Chapter 4: Basic Analyses

---

Tyson S. Barrett

Summer 2017

Utah State University

Introduction

T-tests

ANOVA

Linear Regression

Reporting Results

Conclusions

# Introduction

---

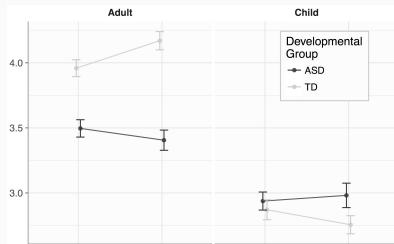
# Basic Analyses

**Basic Analyses:** The analyses taught in the first stats course

These include:

1. T-tests
2. ANOVA
3. Linear Regression

These allow us to assess relationships like that in the figure.



Maybe surprising: \ These all are doing essentially the same thing!

First, **T-TESTS!**

# T-tests

---

1. Simple
2. Independent Samples
3. Paired Samples

# Three Types

Each will be demonstrated using:

```
df <- data.frame("A"=sample(c(0,1), 100, replace = TRUE),  
                 "B"=rnorm(100),  
                 "C"=rnorm(100))
```

df

	A	B	C
1	0	-0.49302332	-1.929185797
2	0	-0.74567281	-0.501050403
3	0	0.76942036	0.972653402
4	0	0.26824581	0.325231701
5	1	0.22727798	-0.380874790
6	1	-1.30398561	0.874102358
7	0	-0.56605613	1.977208866
8	0	-1.12392072	0.146750347
9	1	-0.45675323	-0.334580377

Comparing a mean of a variable with  $\mu$ .

```
t.test(df$B, mu = 0)
```

One Sample t-test

```
data: df$B
t = -1.2338, df = 99, p-value = 0.2202
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.3102355  0.0723451
sample estimates:
mean of x
-0.1189452
```



# Independent Samples

Comparing the means of two groups (dfA is the grouping variable).

```
t.test(df$B ~ df$A)
```

Welch Two Sample t-test

```
data: df$B by df$A
```

```
t = -0.91829, df = 87.734, p-value = 0.361
```

```
alternative hypothesis: true difference in means is not equal to
```

```
95 percent confidence interval:
```

```
-0.5715761  0.2103016
```

```
sample estimates:
```

```
mean in group 0 mean in group 1
```

```
-0.19842557      -0.01778835
```

# Paired Samples

Comparing repeated measures (e.g., Pretest vs. Posttest).

```
t.test(df$B, df$C, paired = TRUE)
```

Paired t-test

data: df\$B and df\$C

t = -1.7104, df = 99, p-value = 0.09033

alternative hypothesis: true difference in means is not equal to

95 percent confidence interval:

-0.56702904 0.04202515

sample estimates:

mean of the differences

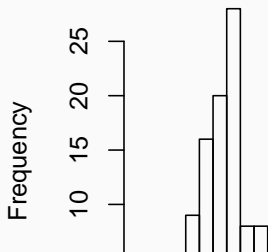
-0.2625019

# Testing Assumptions of T-Tests

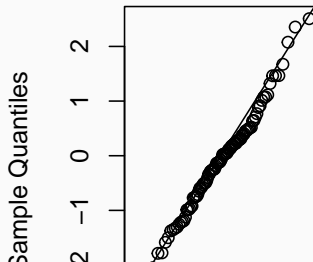
T-tests require that the data be normally distributed with approximately the same variance.

```
## Normality
par(mfrow = c(1,2))
hist(df$B)
qqnorm(df$B)
abline(a=0, b=1)
```

**Histogram of df\$B**



**Normal Q-Q Plot**



# ANOVA

---

# Analysis of Variance

The Analysis of Variance (ANOVA) is highly related to t-tests but can handle 2+ groups.

1. Provides the same p-value as t-tests
2.  $t^2 = F$

For example:

```
fit_ano = aov(df$B ~ df$A)
summary(fit_ano)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
df\$A	1	0.80	0.8040	0.864	0.355
Residuals	98	91.21	0.9307		

```
t.test(df$B ~ df$A)$p.value
```

```
[1] 0.3609868
```

# Analysis of Variance

```
fit_ano = aov(df$B ~ df$A)
summary(fit_ano)
t.test(df$B ~ df$A)$p.value
```

Notice in the code:

- We assigned the `aov()` the name `fit_ano` (which we could have called anything)
- We used the `summary()` function to see the F and p values.
- We pulled the p-value right out of the `t.test()` function.

1. One-Way
2. Two-Way (Factorial)
3. Repeated Measures
4. A combination of Factorial and Repeated Measures

# Types

We will use the following data set for the examples:

```
library(tidyverse)
df <- data.frame("A"=sample(c(0,1), 100, replace = TRUE) %>% fac
               "B"=rnorm(100),
               "C"=rnorm(100),
               "D"=sample(c(1:4), 100, replace = TRUE) %>% fac
df
```

	A	B	C	D
1	0	-0.41875713	0.303400934	3
2	0	-0.59433161	-0.975528435	3
3	1	-1.40803921	0.064343332	3
4	1	-1.10251217	-0.223957259	4
5	0	0.11895720	-0.528506918	3
6	1	-1.92523501	-0.908708582	1
7	0	0.16546402	0.254348633	1
8	1	0.91605559	1.704428403	2



# One-Way

A One-Way ANOVA can be run using `aov()`.

```
fit1 = aov(B ~ D, data = df)
summary(fit1)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
D	3	6.45	2.1489	2.326	0.0796 .
Residuals	96	88.68	0.9238		

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

## Two-Way

A Two-Way ANOVA uses essentially the exact same code with a minor change—including the other variable in an interaction.

```
fit2 = aov(B ~ D * A, data = df)
summary(fit2)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
D	3	6.45	2.1489	2.473	0.0666 .
A	1	0.47	0.4738	0.545	0.4622
D:A	3	8.26	2.7518	3.166	0.0281 *
Residuals	92	79.95	0.8691		

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

The D:A line highlights the interaction term whereas the others show the main effects.

## Repeated Measures

To show this, we will add a fake ID variable to our already fake data set df.

```
df$ID = 1:100
```

And change our data to long (Can you remember how to do it?)

```
library(tidyverse)
df_long = gather(df, "var", "value", 2:3)
df_long
```

	A	D	ID	var	value
1	0	3	1	B	-0.418757129
2	0	3	2	B	-0.594331614
3	1	3	3	B	-1.408039212
4	1	4	4	B	-1.102512174
5	0	3	5	B	0.118957199
6	1	1	6	B	-1.925235006
7	0	1	7	B	0.165464016

The repeated measures, besides using a long-form of the data, is very similar in code. In addition to our usual formula (e.g., `something ~ other + stuff`), we have the `Error()` function. This function tells R how the repeated measures are clustered. In general, you'll provide the subject ID. The next slide highlights this.

# Repeated Measures

```
fit3 = aov(value ~ var + Error(ID), data = df_long)
summary(fit3)
```

Error: ID

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Residuals	1	0.03218	0.03218		

Error: Within

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
var	1	2.34	2.3419	2.405	0.123
Residuals	197	191.81	0.9737		

Here, value was the value of the repeated measures where var is the time. That means our outcome is testing if there were any differences from pre-test to post-test across all the groups.

To take the repeated measures a step further, we can do a Three-Way Repeated Measures ANOVA.

```
fit4 = aov(value ~ var * D * A + Error(ID), data = df_long)
summary(fit4)
```

The output is on the next slide...

# Combination

Error: ID

	Df	Sum Sq	Mean Sq
D	1	0.03218	0.03218

Error: Within

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
var	1	2.34	2.3419	2.461	0.1185
D	3	6.65	2.2179	2.330	0.0758 .
A	1	0.20	0.1953	0.205	0.6511
var:D	3	1.20	0.3993	0.420	0.7392
var:A	1	0.29	0.2892	0.304	0.5821
D:A	3	3.68	1.2257	1.288	0.2800
var:D:A	3	5.63	1.8753	1.970	0.1200
Residuals	183	174.17	0.9518		

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

# Checking Assumptions

Of course, as with any statistical analysis, there are assumptions.

Many of these we can test.

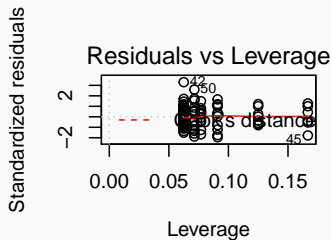
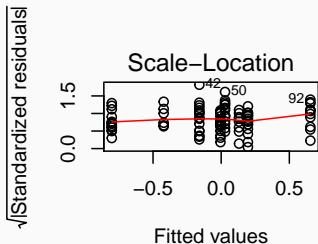
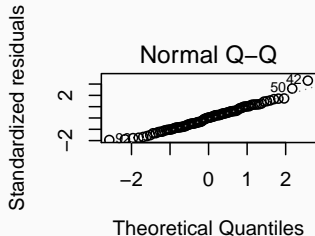
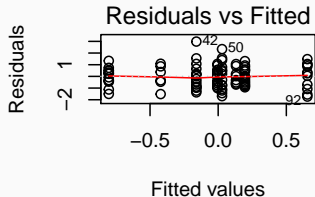
Using our `fitX` objects from our ANOVAs above, we can look at our assumptions:

```
par(mfrow = c(2,2)) ## puts the four plots on a 2 x 2 grid
plot(fit2)
```

Again, the output is on the next slide...



# Checking Assumptions



# Checking Assumptions

They don't fit great on the slides but trust me that normality looks good.  
The assumption of homogeneity of variance looks good as well.

But, if you wanted to test it, you could.

```
library(car)
leveneTest(fit2)
```

```
Levene's Test for Homogeneity of Variance (center = median)
      Df F value Pr(>F)
group  7  1.7072 0.1169
      92
```

Large p-value here is a good thing: `emo::ji("smile")` <sup>1</sup>

<sup>1</sup>This shows a smiley in 'R', just not on these slides—from the 'emo' package on GitHub.

# Linear Regression

---

Once again, linear regression is essentially the more flexible twin of ANOVA and t-tests.<sup>2</sup>

It can:

1. Handle continuous and categorical predictors (i.e., independent variables)
2. Less stringent assumption of equality of variances
3. Is what many other methods are built on (Chapter 5 and 6 will talk about some of these)

---

<sup>2</sup>It mainly only differs from ANOVA in the way it takes a dummy code rather than an effect code of the categorical variables.

# Linear Regression

We will use `lm()` (Linear Model) to fit these models.

```
fit5 = lm(B ~ A, data = df)
summary(fit5)
```

Call:

```
lm(formula = B ~ A, data = df)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.2052	-0.6925	-0.1018	0.6285	2.9226

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.03416	0.14493	0.236	0.814
A1	-0.13420	0.19722	-0.680	0.498

Residual standard error: 0.9829 on 98 degrees of freedom

# Linear Regression

We can add an interaction with the \*.

```
fit6 = lm(B ~ A*D, data = df)
summary(fit6)
```

Call:

```
lm(formula = B ~ A * D, data = df)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.73077	-0.69977	0.02393	0.55829	2.98275

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.19623	0.23306	0.842	0.40198
A1	-0.99870	0.34809	-2.869	0.00511 **
D2	-0.06529	0.40367	-0.162	0.87187
D3	-0.20149	0.32960	-0.611	0.54250

## Other Specifications

We can also make adjustments to the variables within the model.

First, we can transform the variables (e.g., log transformation).

```
fit7 = lm(log(B) ~ A*D, data = df)
summary(fit7)
```

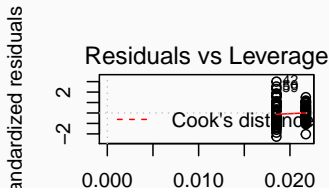
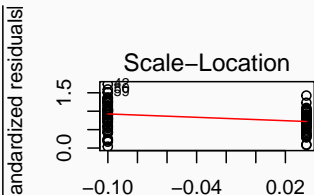
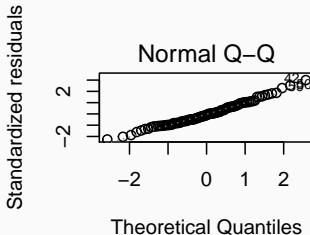
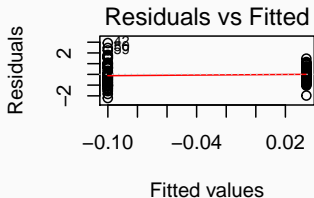
We can change the reference level of a variable, too.

```
fit8 = lm(B ~ relevel(D,ref = "4"), data = df)
summary(fit8)
```

# Checking Assumptions

Assumption checking is similar to that of the linear model.

```
par(mfrow = c(2,2))  
plot(fit5)
```





# Reporting Results

---

## Making This into a Table

Often we want to present this information in a table. This can be done is several ways:

1. Pulling information out of the model objects directly
2. Using a package like `stargazer` to do that work for you
3. Manually by hand

We can certainly do number 3 but why? So we'll look at both 1 and 2.

# Pull information out of the model objects

The model objects contain loads of information that we can pull out:

1. Coefficients
2. Standard Errors and P-values
3. Confidence Intervals
4. Fit Statistics
5. Predicted Values
6. and more! <sup>3</sup>

---

<sup>3</sup>For a low cost of \$49.99! Kidding. . .

# Pull information out of the model objects

To see what the model object holds:

```
names(fit5)
```

```
[1] "coefficients" "residuals"    "effects"      "rank"
[5] "fitted.values" "assign"        "qr"           "df.residual"
[9] "contrasts"     "xlevels"      "call"         "terms"
[13] "model"
```

```
names(summary(fit5))
```

```
[1] "call"          "terms"         "residuals"     "coefficients"
[5] "aliased"       "sigma"         "df"            "r.squared"
[9] "adj.r.squared" "fstatistic"    "cov.unscaled"
```

## Pull information out of the model objects

Using that information we can grab:

```
summary(fit5)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.0341622	0.144925	0.2357233	0.8141393
A1	-0.1342035	0.197218	-0.6804832	0.4978031

or

```
summary(fit5)$fstatistic
```

value	numdf	dendf
0.4630573	1.0000000	98.0000000

## Pull information out of the model objects

Put it in a table:

```
rbind(data.frame(summary(fit5)$coefficients, "Type"="Simple Regr  
      data.frame(summary(fit6)$coefficients, "Type"="Interaction
```

	Estimate	Std..Error	t.value	Pr...t..
(Intercept)	0.03416220	0.1449250	0.2357233	0.814139298
A1	-0.13420350	0.1972180	-0.6804832	0.497803125
(Intercept)1	0.19623455	0.2330595	0.8419935	0.401975067
A11	-0.99869651	0.3480920	-2.8690588	0.005106257
D2	-0.06528975	0.4036708	-0.1617401	0.871865234
D3	-0.20148573	0.3295959	-0.6113115	0.542500969
D4	-0.61820637	0.4462749	-1.3852591	0.169322861
A1:D2	1.52192513	0.5557046	2.7387304	0.007406153
A1:D3	1.03230395	0.4874023	2.1179709	0.036873921
A1:D4	1.26047333	0.5659765	2.2270773	0.028381762

Type

(Intercept) Simple Regression

## Pull information out of the model objects

On the previous slide we:

1. Created two `data.frame` with the coefficients and a variable called "Type"
2. Glued them together by row with `rbind()`

This is a simple way of putting a table together that you can later export.

# Use a package like stargazer to do that work for you

A simpler but less flexible way is using a package like stargazer.

```
library(stargazer)
stargazer(fit5, fit6, type = "text")
```

Dependent variable:		
	B	
	(1)	(2)
A1	-0.134 (0.197)	-0.999*** (0.348)
D2		-0.065 (0.404)
D3		-0.201 (0.330)
D4		-0.618 (0.446)
A1:D2		1.522*** (0.556)
A1:D3		1.032** (0.487)
A1:D4		1.260** (0.566)
Constant	0.034	0.196



## Use a package like `stargazer` to do that work for you

This particular package can take several model objects and produce a nice table. It is hard to see but it includes the number of observations, fit statistics, the coefficients, and f-statistics.

Other packages exist that do similar things (e.g., `texreg`).

```
library(texreg)
screenreg(list(fit5, fit6))
```

# Conclusions

---

1. Performing linear models is straightforward in 'R'
2. With a few lines of code, we can fit a model and check model assumptions
3. We can easily turn our model information into an informative table

