General Linear Model

So far

- Reviewed last semester and previewed this semester
 - Theme 1: The Backstreet Boys
 - Theme 2: What is 0?
- Put some checks and balances in place to minimize errors as much as possible

Thinking in terms of models

- All the tests run tus far can be thought of as a model for how you think "the world works"
- ullet Our DV (here forth Y) is what we are trying to understand
- We hypothesize it has some relationship with your IV(s) (here forth Xs), with what is left over described as error (\$E\$)
- \bullet Y = X + E

See this in our R code

Independent samples t-test

```
t.1 <- t.test(y ~ x, data = d)
# y is cont and x is a categoriocal/nominal (dichotomous) fac</pre>
```

One-way ANOVA

```
a.1 <- aov(y ~ x, data=d)
# y is cont and x is a categoriocal/nominal factor</pre>
```

General linear model (GLM)

- This model (equation) can be very simple as in a treatment/control experiment
- It can be very complex in terms of trying to understand something like academic achievement
- The majority of our models fall under the umbrella of a general(ized) linear model
- Models imply our theory about how the data are generated (ie how the world works)

$$Y_i = b_0 + b_1 X_i + e_i$$

- ullet Y / DV / Outcome / Response / Criterion
- ullet X / IV / Predictor / Explanatory variable
- Regression coefficient (weight) / b / b* / β
- Intercept b_0 / β_0
- Error / Residuals e
- ullet Predictions \hat{Y}
- $ullet Y_i \sim Normal(\mu, \sigma)$
- \bullet The DV, Y for each person i is distributed normally, with a mean of μ and a standard deviation of σ

Regression models

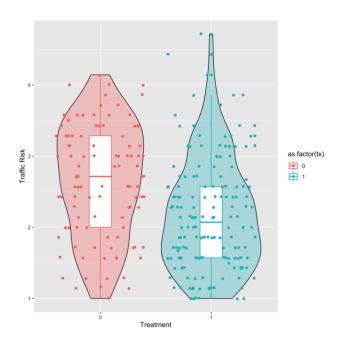
- These models are a way to convey the relationship between two (or more) variables. They translate our hypotheses into math.
- We can use these models to get information we may be interested in (e.g. means, SEs) and test hypotheses about the relationship among variables
- "All models are wrong but some are useful (and some are better than others)" - George Box

Example data can be found on GitHub, or here:

exampleData.csv

example.data

```
## # A tibble: 270 × 3
                tx traffic.risk
##
          id
      <dbl> <dbl>
                           <dbl>
##
##
           1
                            1.86
##
                            1
                            3.29
##
##
##
                            2.43
                            3.29
##
##
                            1.17
##
                            2.43
##
                            3
##
          10
                             1.71
   10
   # i 260 more rows
```



```
##
## Two Sample t-test
##
## data: traffic.risk by tx
## t = 4.9394, df = 268, p-value = 0.000001381
## alternative hypothesis: true difference in means between group 0 a
## 95 percent confidence interval:
## 0.2893360 0.6728755
## sample estimates:
## mean in group 0 mean in group 1
## 2.650641 2.169535
```

```
a.1 <- aov(traffic.risk ~ tx, data = example.data)
summary(a.1)</pre>
```

```
mod.1 <- lm(traffic.risk ~ tx, data = example.data)
summary(mod.1)</pre>
```

```
##
## Call:
## lm(formula = traffic.risk ~ tx, data = example.data)
##
## Residuals:
      Min
              10 Median
##
                             30
                                   Max
## -1.65064 -0.59811 -0.02668 0.54475 2.54475
##
## Coefficients:
            Estimate Std. Error t value
##
                                            Pr(>|t|)
## tx -0.48111 0.09740 -4.939
                                          0.00000138 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7789 on 268 degrees of freedom
## Multiple R-squared: 0.08344, Adjusted R-squared: 0.08002
## F-statistic: 24.4 on 1 and 268 DF, p-value: 0.000001381
```

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

Example summary

Same *p*-values for each test; same SS; same test!

- *t*-test gives you a *t* & df (output may also give you group mean and SD)
- ANOVA gives you an SSs, dfs, MS and Fs
- Linear model (regression) gives you an equation (and SSs and Fs)
- Which one is more useful?

ANOVA as regression

$$Y_i = b_0 + b_1 X_i + e_i$$
 $T. \, risk_i = b_0 + b_1 T X_i + e_i$

- ullet Each individual has a unique Y value an X value and a residual/error term
- The model only has a single b_0 and b_1 term. These are the regression parameters. b_0 is the intercept and b_1 quantifies the relationship between your model of the world and the DV.

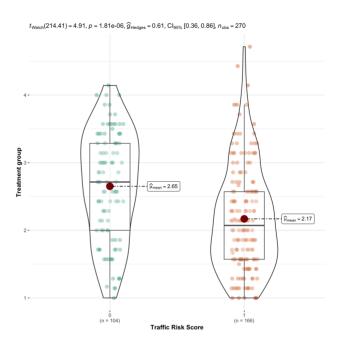
What do the estimates tell us?

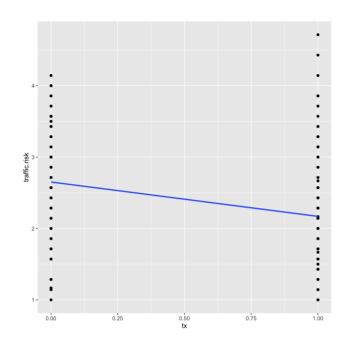
```
## # A tibble: 2 × 5
    term estimate std.error statistic p.value
##
## <chr>
              <dbl>
## 1 (Intercept) 2.65 0.0764 34.7 3.80e-101
             -0.481 0.0974 -4.94 1.38e- 6
## 2 tx
example.data %>%
  group_by(tx) %>%
  summarise(mean(traffic.risk))
## # A tibble: 2 × 2
      tx `mean(traffic.risk)`
##
## <dbl>
                     <dbl>
## 1
                      2.65
## 2
                      2.17
       1
```

How to interpret regression estimates

- Intercept is the mean of group of variable tx that is coded 0
- Regression coefficient is the slope or rise over run, scaled as a 1 unit on the x axis
- "For a one unit change in X, there is a b1 predicted change in Y."
- Regression coefficient is the difference in means between the groups, given that we coded our groups as 0 and 1.

```
library(ggstatsplot)
ggstatsplot::ggbetweenstats(
  data = example.data,
  x = tx,
  y = traffic.risk,
  xlab = "Traffic Risk Score
  ylab = "Treatment group",
  bf.message = FALSE,
  messages = FALSE
)
```





This is what it looks like if we wanted to put a "regression line" to the data. Note that the same interpretation for a regression line holds: for a 1 unit change in X (tx) there is a predicted b change in Y (traffic risk)

Interpretations

- ullet Intercept (b_0) signifies the level of Y when your model IVs (Xs) are zero
- ullet Regression (b1) signifies the difference for a one unit change in your X
- As with last semester you have estimates (like \bar{x}) and standard errors, which you can then ask whether they are likely assuming a null or create a CI

Predicted scores

 Based on the output how do I calculate means for each group?

ANOVA as regression

- "For a one unit change in X, there is a b_1 predicted change in Y" -- always true
- Nominal/categorical variables do not have any inherent numbers associated with them so we need to assign them numbers
- What numbers you assign will impact the equation/estimates/hypothesis you can test
- Make them useful! O and 1 are useful and are the default in R

ANOVA as regression

$$T. risk_i = b_0 + b_1 TX_i + e_i$$

```
## # A tibble: 270 × 3
     id tx traffic.risk
 <dbl> <dbl>
          <dbl>
         1 1.86
## 1
## 2 2 1
## 3 3 1 3.29
## 4 4 1
## 5 5
         1 2.43
## 6 6 1 3.29
1.17
## 8 8
         0 2.43
##
## 10
     10
               1.71
  ## # i 260 more rows
```

```
library(dplyr)
example.data$tx.r <- as.factor(example.data$tx)
example.data$tx.r <- recode_factor(example.data$tx.r, "0" = "</pre>
```

Create a new variable that is not numeric

example.data

```
## # A tibble: 270 × 4
               tx traffic risk tx.r
##
         id
      <dbl> <dbl>
                         <dbl> <fct>
##
##
   1
          1
                          1.86 treatment
                                treatment
##
          3
##
   3
                          3.29 treatment
                                treatment
##
          4
          5
##
    5
                          2.43 treatment
                1
##
                          3.29 treatment
                          1.17 control
##
                0
##
    8
                          2.43 control
                0
                                control
##
                0
                          1.71 treatment
##
  10
         10
                1
## # i 260 more rows
```

```
mod.1 <- lm(traffic.risk ~ tx, data = example.data)</pre>
tidv(mod.1)
## # A tibble: 2 × 5
## term estimate std.error statistic p.value
## <chr> <dbl> <dbl> <dbl> <dbl>
## 1 (Intercept) 2.65 0.0764 34.7 3.80e-101
## 2 tx -0.481 0.0974 -4.94 1.38e- 6
mod.1r <- lm(traffic.risk ~ tx.r, data = example.data)</pre>
tidv(mod.1r)
## # A tibble: 2 × 5
## term estimate std.error statistic p.value
## <chr> <dbl> <dbl> <dbl> <dbl>
## 1 (Intercept) 2.65 0.0764 34.7 3.80e-101
## 2 tx.rtreatment -0.481 0.0974 -4.94 1.38e- 6
```

What if we changed 0 and 1 to other values?

- Infinite number of ways to code categorical/nominal variables, only a few meaningful ways
- The R default is called "dummy coding"
- Uses 0s and 1s to put numbers to categories. We will soon see what this looks like when you have more than 2 groups.
- Changing the numbers changes...?

Effect coding

```
example.data$tx.effect <- dplyr::recode(example.data$tx,</pre>
```

example.data

```
## # A tibble: 270 × 5
               tx traffic.risk tx.r tx.effect
##
         Ьi
      <dbl> <dbl>
                          <dbl> <fct>
                                               <dbl>
##
##
                           1.86 treatment
                                treatment
##
##
                           3.29 treatment
                                treatment
##
          5
                           2.43 treatment
##
                           3.29 treatment
##
                           1.17 control
##
##
                0
                           2.43 control
                                control
##
    9
                0
                                                   -1
         10
                           1.71 treatment
##
   10
## # i 260 more rows
```

Effect coding

```
mod.1.eff <- lm(traffic.risk ~ tx.effect, data = example.data
tidy(mod.1.eff)</pre>
```

• systematically changes both the intercept and the regression estimate

```
## # A tibble: 2 × 5
##
   term estimate std.error statistic
                                 p.value
         ## <chr>
                                   <dbl>
## 1 (Intercept) 2.41 0.0487 49.5 8.15e-137
## 2 tx.effect -0.241 0.0487 -4.94 1.38e- 6
## # A tibble: 2 × 5
##
   term estimate std.error statistic
                                 p.value
## <chr>
          < db1 >
## 1 (Intercept) 2.65 0.0764 34.7 3.80e-101
            -0.481 0.0974 -4.94 1.38e- 6
## 2 tx
```

- Intercept: value when your predictor X is zero. WHAT DOES ZERO MEAN?!
- ullet Regression coefficient: one unit increase in X is associated with a (regression estimate) predicted increase in Y

Effect coding

Consists of -1, 1s (And zeros for more than 2 groups)

- 1. The intercept is the "grand mean" or "mean of means" if unbalanced
- 2. The regression coefficient represents the group "effect":
 - the difference between the mean of means and the group labeled 1 (we will revisit this when we have more than 2 groups)

Dummy coding

- More appropriate when you are interested in comparing to a specific group rather than an "average person"
- Intercept: value of the group coded zero
- Regression coefficient: mean difference between groups

Contrast coding

- As our models get more complex our coding schemes can too
- What happens if you code the groups -.5,
 -.5, and 1?
- These make more sense when we have more groups. More groups require more independent variables

Statistical Inference

- The way the world is = our model + error
- How good is our model? Is it a good representation of reality? Does it "fit" the data well?
- Need to go beyond asking if it is significant, because what does that mean?
- We are going to make predictions and see if the predictions (based on our model) matches our data
- We can then compare one model to another to see which one matches our data better.
 Which one is a better representation of reality?

Predictions

- ullet predictions \hat{Y} are of the form of E(Y|X)
- ullet They are created by simply plugging a persons Xs into the created model
- If you have b_0 , b_1 , and Xs then you can create a prediction

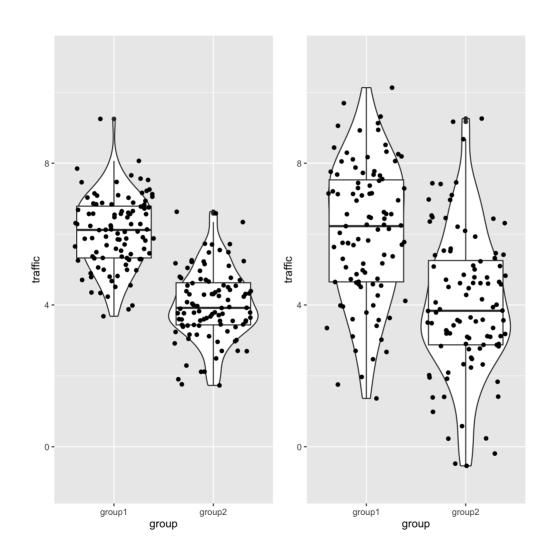
$$\hat{Y}_i$$
 = 2.65064 + -0.48111* X_i

 You have already done this with dummy codes above

Predictions

- ullet We want our predictions to be close to our actual data for each person (Y_i)
- The difference between the actual data and our our prediction ($Y_i \hat{Y}_i = e$) is the residual, how far we are "off". This tells us how good our fit is.
- You can have the same estimates for two models but completely different fit.
- Previously you have evaluated fit by looking at Eta Squared, SS Error and visualizing observations around group means

Which one has better fit?



Easy to examine fit with **lm** objects

 These are automatically created anytime you run a lm in R

```
mod.1 <- lm(traffic.risk ~ tx, data = example.data)</pre>
```

```
coefficients(mod.1) # coefficients
residuals(mod.1) # residuals
fitted.values(mod.1) # fitted values ie predicted
summary(mod.1)$r.squared # R-sq for the model
summary(mod.1)$sigma # sd of residuals
```

coefficients(mod.1)

```
## (Intercept) tx
## 2.6506410 -0.4811057
```

fitted.values(mod.1)

	7	6	5	4	3	2	1	##
2.6506	2.650641	2.169535	2.169535	2.169535	2.169535	2.169535	2.169535	##
	18	17	16	15	14	13	12	##
2.6506	2.650641	2.650641	2.650641	2.650641	2.650641	2.650641	2.650641	##
;	29	28	27	26	25	24	23	##
2.6506	2.650641	2.650641	2.650641	2.650641	2.650641	2.169535	2.169535	##
)	40	39	38	37	36	35	34	##
2.6506	2.169535	2.169535	2.169535	2.169535	2.650641	2.650641	2.650641	##
	51	50	49	48	47	46	45	##
2.1695	2.169535	2.650641	2.169535	2.650641	2.650641	2.650641	2.169535	##
. •	62	61	60	59	58	57	56	##
2.1695	2.169535	2.650641	2.169535	2.169535	2.169535	2.169535	2.650641	##
-	73	72	71	70	69	68	67	##
2.1695	2.169535	2.169535	2.169535	2.169535	2.169535	2.650641	2.650641	##
	84	83	82	81	80	79	78	##
2.6506	2.650641	2.650641	2.650641	2.169535	2.650641	2.169535	2.169535	##
9	95	94	93	92	91	90	89	##
2.1695	2.169535	2.169535	2.169535	2.650641	2.169535	2.169535	2.169535	##
10	106	105	104	103	102	101	100	##
2.6506	2.650641	2.650641	2.169535	2.169535	2.169535	2.169535	2.169535	##
1.	117	116	115	114	113	112	111	##
2.6506	2.650641	2.650641	2.650641	2.650641	2.650641	2.650641	2.650641	##

122 123 124 125 126 127 128 39/59¹²
2.650641 2.650641 2.169535 2.169535 2.169535 2.169535 2.169535 2.65064

residuals(mod.1)

##	1	2	3	4	5	
##	-0.312392427	-1.169535284	1.116179002	-0.169535284	0.259036145	1.11
##	9	10	11	12	13	
##	0.349358974	-0.455249570	1.687607573	-0.079212455	0.206501831	0.92
##	17	18	19	20	21	
##	-0.936355312	-1.364926740	-0.364926740	-0.312392427	-0.598106713	-0.02
##	25	26	27	28	29	
##	1.063644688	-0.650641026	-0.079212455	0.206501831	0.777930403	-1.07
##	33	34	35	36	37	
##	0.777930403	-0.936355312			0.116179002	-0.59
##	41	42	43	44	45	
##	0.635073260	-0.598106713	-0.507783883	0.635073260	0.116179002	-0.79
##	49	50	51	52	53	
##	0.830464716			-0.598106713	-0.312392427	-0.88
##	57	58	59	60	61	
##	-1.169535284	0.259036145	-0.455249570	-0.169535284	-0.079212455	0.54
##	65	66	67	68	69	
##	0.830464716	2.259036145	-0.650641026	-1.650641026	-0.169535284	1.68
##	73	74	75	76	77	
##	0.973321859				0.849358974	-0.59
##	81	82	83	84	85	
##	1.401893287	-1.364926740			0.492216117	1.68
##	89	90	91	92	930	0 / 59
##	-0.312392427	-0.312392427	1.973321859	0.920787545	-1.169535284	

Pop quiz

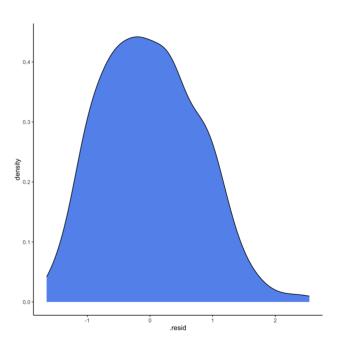
$$\hat{Y_i} = b_0 + b_1 X_i$$
 $Y_i = b_0 + b_1 X_i + e_i$
 $Y_i - \hat{Y_i} = e$

- Can you plug in numbers and calculate subject 3's predicted and residual scores without explicitly asking for Im object residuals and fitted values? (using the same model from slide 37 on)
- Post answers in Slack -- see if you and your peers get the same results!

example.data

```
## # A tibble: 270 × 5
##
         id
                tx traffic.risk tx.r tx.effect
      <dbl> <dbl>
                          <dbl> <fct>
                                                <dbl>
##
##
    1
          1
                 1
                           1.86 treatment
    2
                                 treatment
##
                 1
                           1
##
                           3.29 treatment
##
          4
                                 treatment
          5
##
                           2.43 treatment
##
          6
                           3.29 treatment
##
                 0
                           1.17 control
                                                   -1
##
                 0
                           2.43 control
                                                   -1
##
          9
                 0
                           3
                                 control
                                                   -1
##
   10
         10
                           1.71 treatment
                                                    1
##
   # i 260 more rows
```

Residuals



lm objects

- Im objects consist
 of the information
 embedded in your
 linear model
- the broom package makes model objects into dataframes

```
library(broom)
fit.1.tidy <- tidy(mod.1)
fit.1.tidy</pre>
```

 Augment function from the broom package amends the original dataset with Im object content. The new variable names of have a "." in front of the name to distinguish

fit.1.data <- augment(mod.1)</pre>

3.29

2.43

3.29

3

4

5

6

```
head(fit.1.data)
## # A tibble: 6 × 8
                   tx .fitted .resid
    traffic.risk
                                        .hat .sigma .cooksd .std.resid
##
           <dbl> <dbl>
                        <dbl> <dbl>
                                       <dbl> <dbl>
                                                      <dbl>
                                                                 <dbl>
##
## 1
            1.86
                     1 2.17 -0.312 0.00602 0.780 0.000490
                                                                -0.402
                      2.17 -1.17 0.00602 0.777 0.00687
## 2
            1
                                                                -1.51
```

0.00602 0.777 0.00626

0.00602 0.777 0.00626

2.17 -0.170 0.00602 0.780 0.000144

2.17 0.259 0.00602 0.780 0.000337

2.17 1.12

2.17 1.12

1.44 -0.218

0.334

1.44

- tidy = model components like β
- glance is similar but for model fit measures; η^2 or R^2
- augment = adds to existing dataset

```
tidy(mod.1)
## # A tibble: 2 × 5
          estimate std.error stat
##
    term
    <chr>
                   <dbl>
                          <dbl>
##
## 1 (Intercept) 2.65
                           0.0764
## 2 tx
                  -0.481 0.0974
glance(mod.1)
## # A tibble: 1 × 12
    r.squared adj.r.squared sigma stati
##
        <dbl>
                      <dbl> <dbl>
##
## 1
       0.0834
                    0.0800 0.779
```

Pop quiz #2

- For a four group Oneway ANOVA how many different predicted values will we have?
 Residuals?
- Post your answer in Slack and see how you compare to your peers!

Statistical Inference

- In making predictions, we have to compare our prediction to some alternative prediction to see if we are doing well or not.
- What is our best guess (ie prediction) if we didn't collect any data?

$$\hat{Y}=ar{Y}$$

• Regression can be thought of as: is E(Y|X) better than E(Y)?

Statistical Inference

- ullet To the extent that we can generate different predicted values of Y based on the values of the predictors, our model is doing well
- ullet The closer our model is to the "actual" data generating model, our guesses (\hat{Y}) will be closer to our actual data (Y)

 We formally test how well we are doing with our guesses by partitioning variation

$$Y = \hat{Y} + e$$
 $Y = \hat{Y} + (Y - \hat{Y})$
 $Y - \bar{Y} = (\hat{Y} - \bar{Y}) + (Y - \hat{Y})$
 $(Y - \bar{Y})^2 = [(\hat{Y} - \bar{Y}) + (Y - \hat{Y})]^2$
 $\sum (Y - \bar{Y})^2 = \sum (\hat{Y} - \bar{Y})^2 + \sum (Y - \hat{Y})^2$

Partitioning the variation in Y

$$\sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum (Y_i - \hat{Y}_i)^2$$

- SS total = SS between + SS within
- SS total = SS regression + SS residual (or error)
- Completely the same as last semester because ANOVA IS REGRESSION

What can we do with this?

- ullet Last semester you did omnibus F tests
- What hypothesis does the omnibus F test test, generally?

$$s_y^2 = s_{regression}^2 + s_{residual}^2$$

$$1 = rac{s^2_{regression}}{s^2_y} + rac{s^2_{residual}}{s^2_y}$$

Coefficient of Determination

$$rac{s_{regression}^2}{s_y^2} = rac{SS_{regression}}{SS_Y} = R^2$$

- Percent (of total) variance explained by your model...which currently are groups
- Another way of asking how much variance group status explains

R^2 and Eta squared

```
summary(mod.1)$r.squared

## [1] 0.08344007

library(lsr)
etaSquared(mod.1)

## eta.sq eta.sq.part
## tx 0.08344007 0.08344007
```

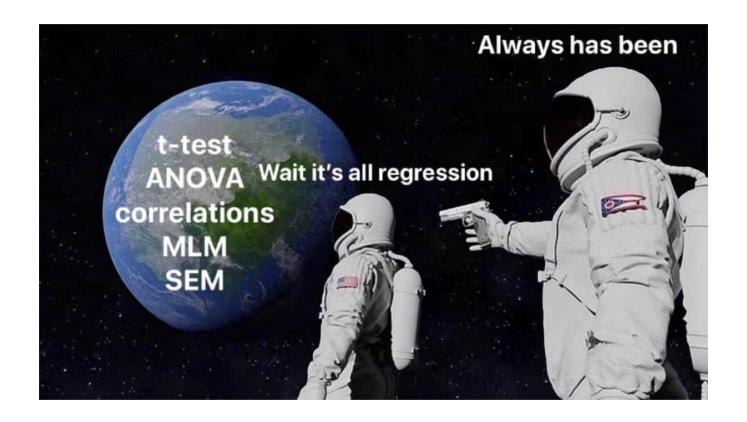
${\mathbb R}^2$ for different coding schemes

```
glance(mod.1)
## # A tibble: 1 × 12
## r.squared adj.r.squared sigma statistic p.value df logLik
     ##
## 1 0.0834 0.0800 0.779 24.4 0.00000138 1 -315.
glance(mod.1.eff)
## # A tibble: 1 × 12
   r.squared adj.r.squared sigma statistic p.value df logLik
##
     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
##
## 1 0.0834 0.0800 0.779 24.4 0.00000138 1 -315.
```

Note the R^2 p-value

```
##
## Call:
## lm(formula = traffic.risk ~ tx, data = example.data)
##
## Residuals:
      Min 1Q Median 3Q
##
                                  Max
## -1.65064 -0.59811 -0.02668 0.54475 2.54475
##
## Coefficients:
           Estimate Std. Error t value
                                           Pr(>|t|)
##
## tx -0.48111 0.09740 -4.939 0.00000138 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7789 on 268 degrees of freedom
## Multiple R-squared: 0.08344, Adjusted R-squared: 0.08002
## F-statistic: 24.4 on 1 and 268 DF, p-value: 0.000001381
```

Summary



Summary

- We are using linear models to do the exact same tests as *t*-tests and ANOVAs
- It is the exact same because t-tests and ANOVAs are part of the general linear model
- Using linear models gives us the same information, and more!
- Provides a more systematic way at 1) building and testing your theoretical model and 2) comparing between alternative theoretical models
- You can get 1) estimates and 2) fit statistics from the model. Both are important.

Next time

Revisiting correlations, and turning those into regression