

Analysis of Variance

BIOE 498/598

2/19/2020

Does our model do anything?

Let's return to our data of my son throwing the stuffed monkey.

```
attach(read.csv("AndersThrow.csv"))
model <- lm(distance ~ 0 + hand + hat + boots)
summary(model)

##
## Call:
## lm(formula = distance ~ 0 + hand + hat + boots)
##
## Residuals:
##      1      2      3      4      5      6      7      8
## -0.375 -1.125  0.625  0.875  1.125  0.375 -1.375 -0.125
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## handleleft      5.375      0.857   6.272 0.003298 **
## handright       7.625      0.857   8.898 0.000882 ***
## hatyes         -1.500      0.857  -1.750 0.154947
## bootsyes        1.000      0.857   1.167 0.308065
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.212 on 4 degrees of freedom
## Multiple R-squared:  0.9825, Adjusted R-squared:  0.9649
## F-statistic: 56.02 on 4 and 4 DF, p-value: 0.0009119
```

The sum of squares

Our analysis is based on the *sum of squares*, or SS . In particular, the total SS is the combination of the SS explained by our model and the SS that is residual (or unexplained).

$$SS_{\text{total}} = SS_{\text{explained}} + SS_{\text{residual}}$$

The sum of squares

Our analysis is based on the *sum of squares*, or SS. In particular, the total SS is the combination of the SS explained by our model and the SS that is residual (or unexplained).

$$SS_{\text{total}} = SS_{\text{explained}} + SS_{\text{residual}}$$

How do we calculate each one for a model $\mathbf{y} = \mathbf{X}\beta$?

The sum of squares

Our analysis is based on the *sum of squares*, or SS. In particular, the total SS is the combination of the SS explained by our model and the SS that is residual (or unexplained).

$$SS_{\text{total}} = SS_{\text{explained}} + SS_{\text{residual}}$$

How do we calculate each one for a model $\mathbf{y} = \mathbf{X}\beta$?

$$SS_{\text{total}} = \sum_i (y_i - \text{mean}(\mathbf{y}))^2$$

The sum of squares

Our analysis is based on the *sum of squares*, or SS. In particular, the total SS is the combination of the SS explained by our model and the SS that is residual (or unexplained).

$$SS_{\text{total}} = SS_{\text{explained}} + SS_{\text{residual}}$$

How do we calculate each one for a model $\mathbf{y} = \mathbf{X}\beta$?

$$SS_{\text{total}} = \sum_i (y_i - \text{mean}(\mathbf{y}))^2$$

$$SS_{\text{residual}} = \sum_i (y_i - \text{predicted}(y_i))^2$$

The sum of squares

Our analysis is based on the *sum of squares*, or SS. In particular, the total SS is the combination of the SS explained by our model and the SS that is residual (or unexplained).

$$SS_{\text{total}} = SS_{\text{explained}} + SS_{\text{residual}}$$

How do we calculate each one for a model $\mathbf{y} = \mathbf{X}\beta$?

$$SS_{\text{total}} = \sum_i (y_i - \text{mean}(\mathbf{y}))^2$$

$$SS_{\text{residual}} = \sum_i (y_i - \text{predicted}(y_i))^2$$

$$SS_{\text{explained}} = SS_{\text{total}} - SS_{\text{residual}}$$

For our throwing data

```
ss <- function(x) sum(x^2)
sst <- ss(distance - 0)
ssr <- ss(residuals(model))
ssx <- sst - ssr
c(sst, ssr, ssx)
```

```
## [1] 335.000    5.875 329.125
```


Degrees of freedom

The amount of variation we see depends on the number of independent parameters in the model. These are the *degrees of freedom*.

- ▶ For $SS_{\text{explained}}$, $DF = \# \text{ of parameters}$
- ▶ For SS_{residual} , $DF = (\# \text{ data points}) - (\# \text{ parameters})$

The F -statistic

The value of our model is explained by the ratio between the explained variance and the residual (unexplained) variance *after adjusting for the DF*.

$$F = \frac{SS_{\text{explained}}/DF(SS_{\text{explained}})}{SS_{\text{residual}}/DF(SS_{\text{residual}})}$$

The F -statistic

The value of our model is explained by the ratio between the explained variance and the residual (unexplained) variance *after adjusting for the DF*.

$$F = \frac{SS_{\text{explained}}/DF(SS_{\text{explained}})}{SS_{\text{residual}}/DF(SS_{\text{residual}})}$$

For our throwing example

$$F = \frac{329.125/4}{5.875/(8-4)} = 56.02$$

How big should the F -statistic be?

That depends on the number of degrees of freedom. The F -statistic follows the F -distribution. We can use this distribution to convert the F -statistic into a p -value.

```
summary(model)
```

```
##
## Call:
## lm(formula = distance ~ 0 + hand + hat + boots)
##
## Residuals:
##      1      2      3      4      5      6      7      8
## -0.375 -1.125  0.625  0.875  1.125  0.375 -1.375 -0.125
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## handleft      5.375      0.857   6.272 0.003298 **
## handright     7.625      0.857   8.898 0.000882 ***
## hatyes       -1.500      0.857  -1.750 0.154947
## bootsyes      1.000      0.857   1.167 0.308065
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.212 on 4 degrees of freedom
## Multiple R-squared:  0.9825, Adjusted R-squared:  0.9649
## F-statistic: 56.02 on 4 and 4 DF,  p-value: 0.0009119
```

Testing single factors

We previously compared the entire model against the residuals to see if the model added value. We can apply the same procedure to a single variable.

This is called the *analysis of variance*, or ANOVA.

ANOVA on handedness

Let's find the explained variance for a model with only handedness:

```
model_hand <- lm(distance ~ 0 + hand)
sst - ss(residuals(model_hand))
```

```
## [1] 322.625
```

Now let's compare this to the residuals of the entire model:

ANOVA on handedness

Let's find the explained variance for a model with only handedness:

```
model_hand <- lm(distance ~ 0 + hand)
sst - ss(residuals(model_hand))
```

```
## [1] 322.625
```

Now let's compare this to the residuals of the entire model:

$$F = \frac{322.625/2}{5.875/(8-4)} = 109.83$$

ANOVA on a linear model

We can repeat this procedure for every variable, or we can use R's built-in ANOVA command.

```
anova(model)
```

```
## Analysis of Variance Table
```

```
##
```

```
## Response: distance
```

##		Df	Sum Sq	Mean Sq	F value	Pr(>F)	
##	hand	2	322.62	161.312	109.8298	0.0003198	***
##	hat	1	4.50	4.500	3.0638	0.1549474	
##	boots	1	2.00	2.000	1.3617	0.3080650	
##	Residuals	4	5.87	1.469			
##	---						
##	Signif. codes:	0	'***'	0.001	'**'	0.01	'*' 0.05 '.' 0.1

Conclusions

- ▶ p -values on effect sizes tell us if the effect size is nonzero.
- ▶ A significant effect size does not mean the effect matters.
- ▶ ANOVA can tell us which variables (not parameters!) explain a significant fraction of the variance in our data.
- ▶ Significance is relative to the unexplained variance in the model.