

Exploratory and Confirmatory Data Analysis

Data Analysis for Psychology in R 2

dapR2 Team

Department of Psychology
The University of Edinburgh

AY 2021-2022

Learning Objectives

- What is the difference between an exploratory and confirmatory analysis?
- Some exploratory situations
 - I have a hypothesis but I'm not quite sure how to test it with the variables I have
 - I think some variables could be relevant to a DV but I'm not sure which ones
- When am I in a position to conduct a confirmatory analysis?
- Working through an example

The Issues:

- We're often interested in the relationship between variables but don't have clear predictions about how they're related
 - For example, I might be interested in why some tweets go viral and others don't
- The number of possible predictors related to this question is huge and it's not obvious which ones will be most important
 - Includes a photo? Humor? Many, many possible predictors

The Issues:

- Sometimes I might have a hypothesis, but it could be tested in multiple ways and I'm not sure how best to test it.
 - For example, I think a tweet including a photo will be retweeted more. What kind of photo though? Any photo? Happy photos?

Exploratory Analyses

- The context I've describe above is a case of exploratory analysis.
- Exploratory analyses can take many forms, but they share in common the fact that you, the researcher, don't have extremely specific predictions about the relationship between your independent variable and your dependent variable
- The exploratory phase of data analysis is a great way to learn a lot about your data, but you also need to be on-guard that you don't think you've detected signal when you've actually detected noise.
 - More on this in a bit

Exploratory analyses done wrong

- "You cannot find your starting hypothesis in your final results. It makes the stats go all wonky." - Ben Goldacre
- If you treat an exploratory result as if you had that hypothesis from the start, then it can cause problems. You will trick yourself.

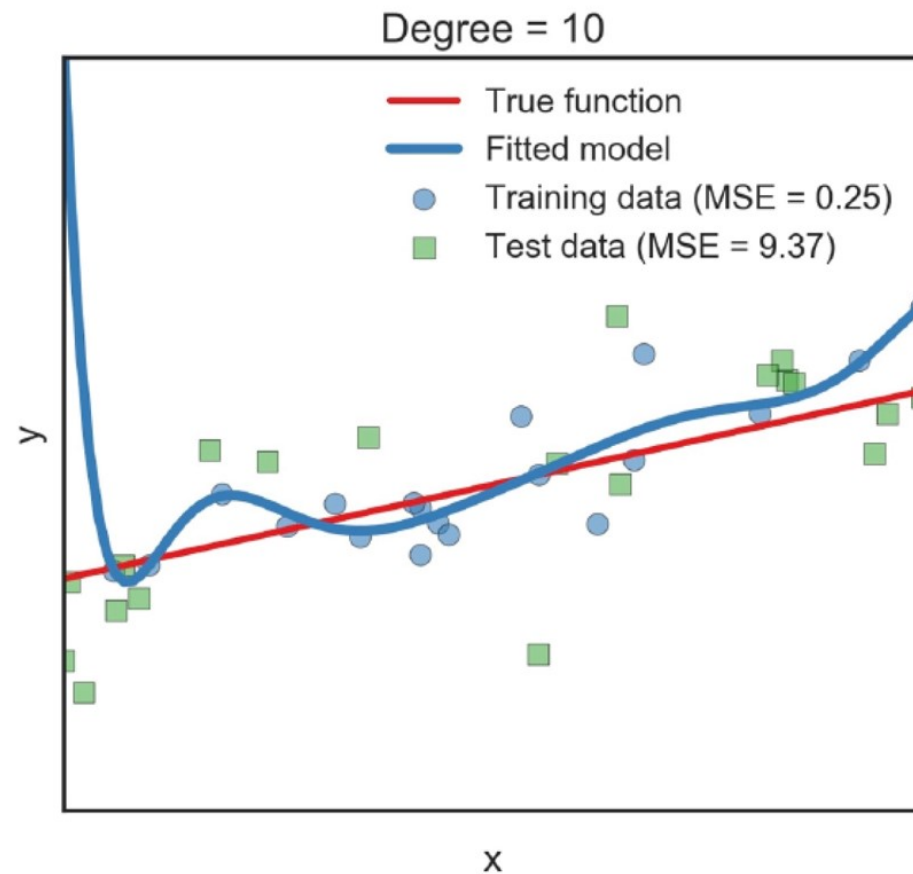
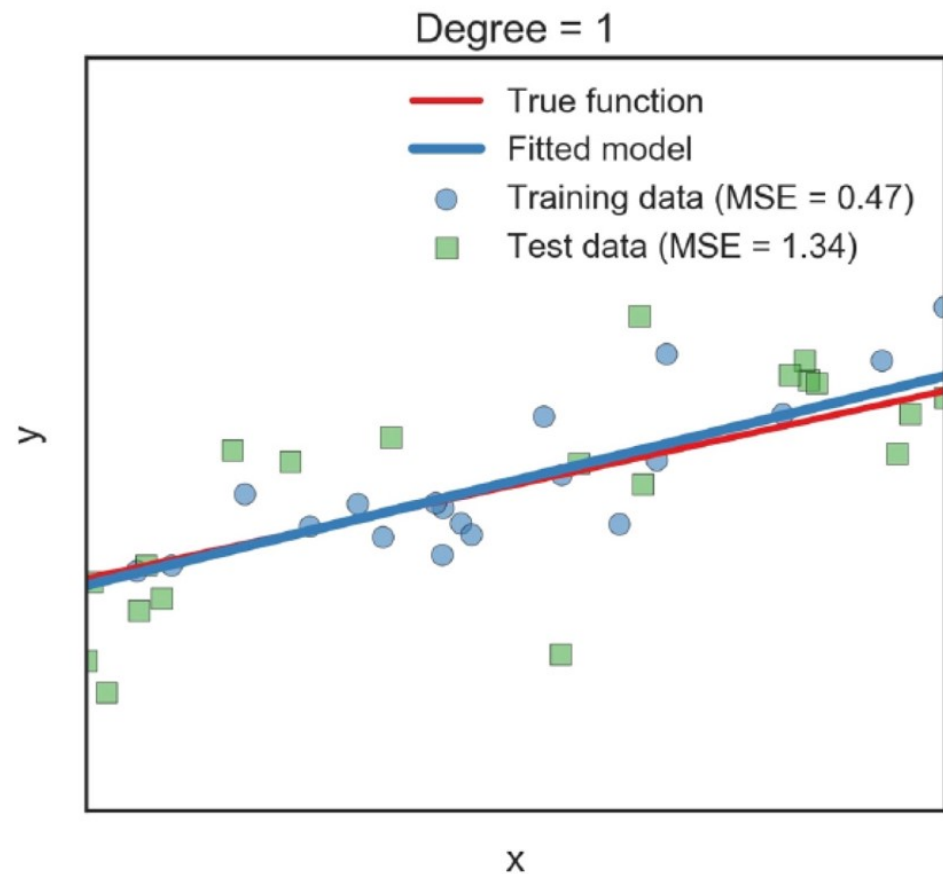
Exploratory analyses done wrong

- Exploratory analyses done poorly: Measure many variables (gender, personality characteristics, age, etc.) and only report those that yield a statistically significant result (stargazing)
 - Include in your paper only those experiments that produced the desired outcome,
 - Treat experiments or initial analyses that didn't turn out favorably as "pilots"
- Why is this problematic?
- In a word, this will lead to model overfitting.

Overfitting

- Overfitting is the tendency for statistical models to mistakenly fit sample-specific noise as if it were signal
 - In a sample of $N = 50$ with 20 uncorrelated predictors, each correlated 0.1 with the DV, the observed (and overfitted) R^2 value will, on average, be 0.45
 - Gives the impression that one could predict values of the DV rather successfully.
 - True R^2 in this situation is only 0.07. Even worse, the average out-of-sample test value of R^2 is only 0.02!

R Squared is very optimistic



Overfitting continued

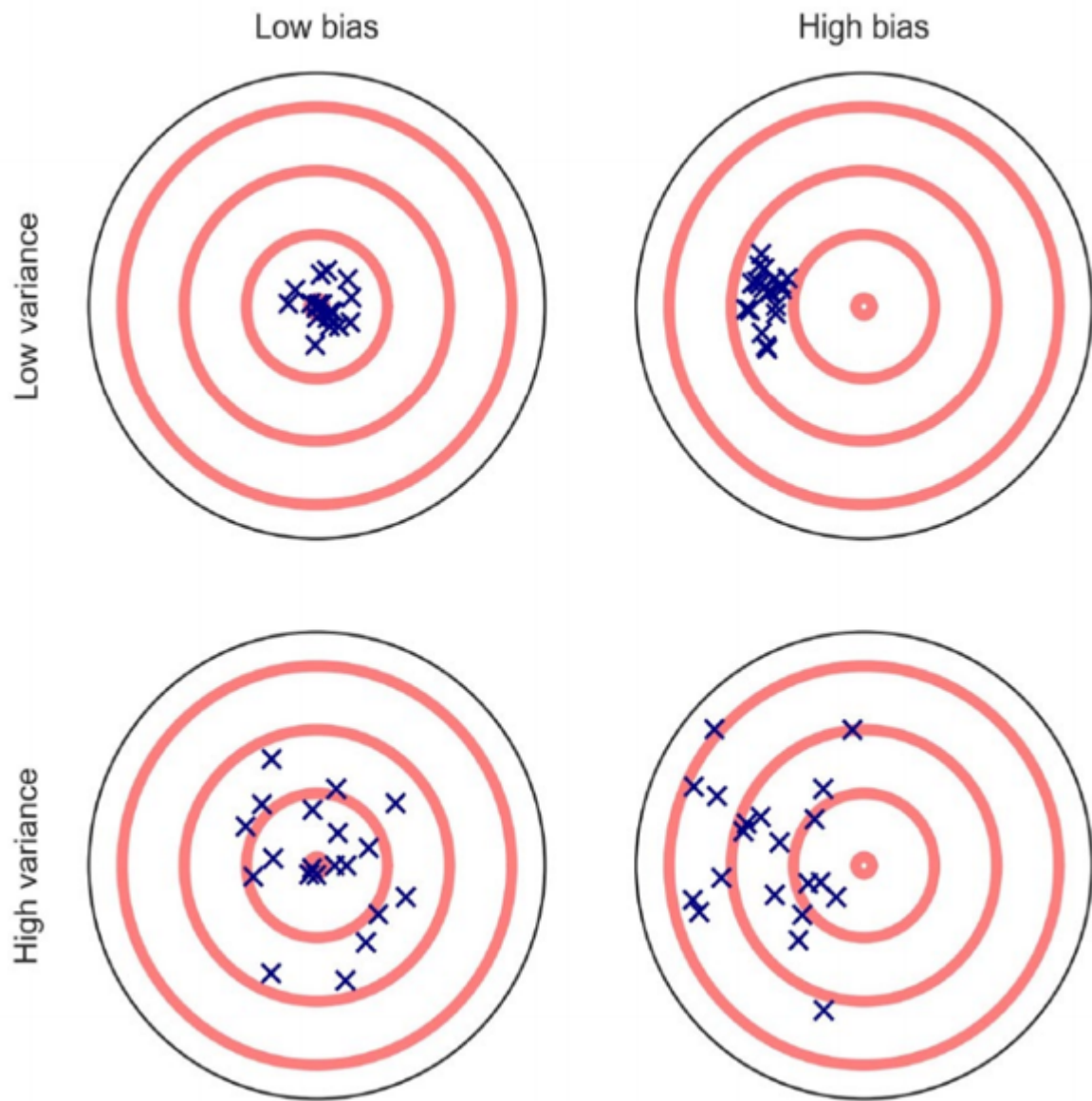
- Don't trust estimates of model performance if those estimates are obtained by "testing" the model on the same data on which it was originally trained
- We need a method for doing exploratory analyses without tricking ourselves.

An aside: The link between p-hacking and overfitting

- p-hacking is a special case of overfitting. Specifically, it is procedural overfitting (Yarkoni and Westfall, 2018). It takes place prior to (or in parallel with) model estimation
 - For example, during data cleaning, model selection, or choosing which analyses to report

We often want to explore though. How do we do that in a principled way?

- First, need to distinguish bias and variance
 - Bias: the tendency for a model to consistently produce answers that are wrong in a particular direction (e.g., estimates that are consistently too high).
 - Variance: the extent to which a model's fitted parameters will tend to deviate from their central tendency across different datasets.



Bias-Variance Tradeoff

- Liberal, flexible data analysis is a low-bias but high-variance approach
 - Almost any pattern in data can potentially be detected, at the cost of a high rate of spurious identifications
 - This is exploratory data analysis
- An approach that favors strict adherence to a fixed set of procedures as a high bias, low-variance approach
 - Only a limited range of patterns can be identified, but the risk of pattern hallucination is low
 - This is confirmatory data analysis

What to do? Consider lots of possibilities but focus on minimizing prediction error (no stargazing!)

- What's required to do exploratory data analysis that gives you *information* on which you can do confirmatory research?
 - Datasets large enough to support training models
 - Accurately estimate prediction error to assess performance and improve model
 - Exert control over the bias-variance tradeoff when appropriate

Cross-validation

- All of these are directly related to cross-validation and replication
 - To assess our models, we need to quantify out-of-sample prediction error
 - Cross-validation: various techniques involved in training and testing a model on different samples of data

Cross-validation

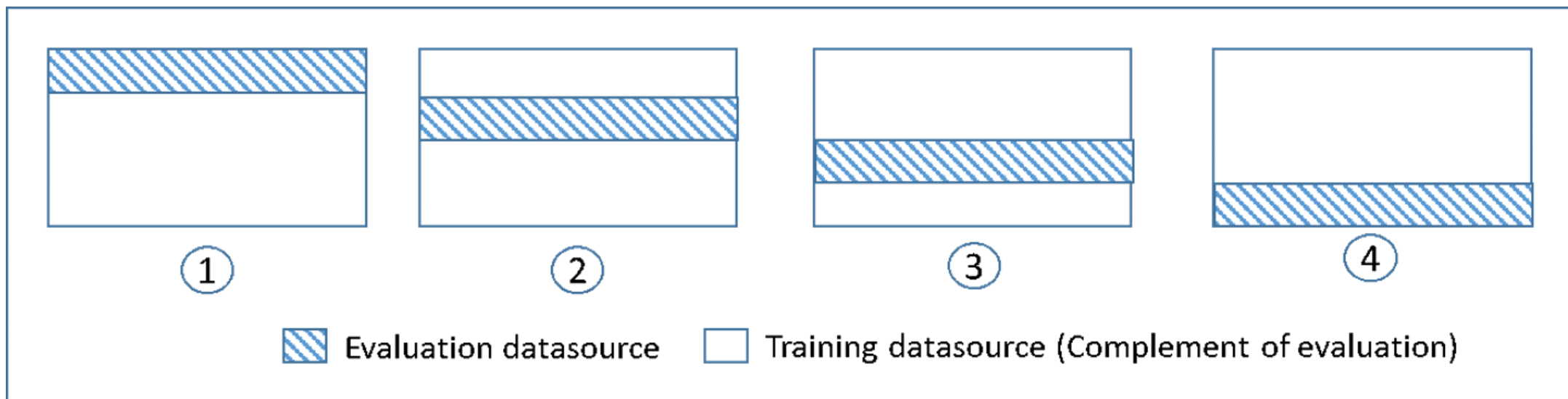
- Canonical cross-validation
 - The classical replication study, where a model is trained on one dataset and then tested on a completely independent dataset. Most typical of experimental research. Less common in correlational research.

Cross-validation

- Sometimes you can't collect more data though
 - One giant study you want to analyze was run once
 - There is a limited population
 - Limited funds to collect more data

Recycle your dataset.

- Don't assign each observation exclusively to either the training or the test datasets - do both!
 - Known as K-Folding where K is the number of folds
 - In one "fold" (essentially a subset of your data), one half of the data is used for training and the other half is for testing
 - In a second fold, the datasets are reversed, and the training set and test sets exchange roles.
 - Typical number of folds is 10



Time for a break

Welcome Back!

Confirmatory research

- The confirmatory phase of research is characterized by the fact that you specify prior to data collection the exact statistical analyses you intend to run, and your expectations about the relationships between the variables
 - For example, "x1 will positively predict dv1 with an effect size of approximately Cohen's $d = .2$ I will test this by fitting the model $y \sim x1 + x2$ ")



Read data into R

```
Exploredf <- read_csv("df1.csv")
```

```
## New names:
## * `` -> ...1
## * ...1 -> ...2

## Rows: 5000 Columns: 7
## -- Column specification -----
## Delimiter: ","
## chr (3): screen_name, media_type, text
## dbl (4): ...1, ...2, favorite_count, retweet_count
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
head(Exploredf)
```

```
## # A tibble: 6 x 7
##   ...1  ...2 screen_name favorite_count retweet_count media_type text
##   <dbl> <dbl> <chr>           <dbl>           <dbl> <chr>   <chr>
## 1      1      1 78462 DogsTrust           41             11 Photo  "With our #Ca~
```

```
Exploredf1 <- Exploredf%>%  
  mutate(Postnumber = 1:n())%>%  
  select(-c(...1))  
  
sentiment <- Exploredf1 %>%  
  unnest_tokens(output = "word", input = "text")
```



```
sentiment_dictionary1 <- get_sentiments("bing")  
head(sentiment_dictionary1)
```

```
## # A tibble: 6 x 2  
##   word      sentiment  
##   <chr>    <chr>  
## 1 2-faces  negative  
## 2 abnormal negative  
## 3 abolish negative  
## 4 abominable negative  
## 5 abominably negative  
## 6 abominate negative
```

```
sentiment_dictionary2 <- get_sentiments("afinn")  
head(sentiment_dictionary2)
```

```
## # A tibble: 6 x 2  
##   word      value  
##   <chr>    <dbl>  
## 1 abandon      -2  
## 2 abandoned    -2  
## 3 abandons     -2  
## 4 abducted     -2  
## 5 abduction    -2  
## 6 abductions   -2
```

```
sentiment_dictionary3 <- get_sentiments("nrc")  
head(sentiment_dictionary3)
```

```
## # A tibble: 6 x 2  
##   word      sentiment  
##   <chr>    <chr>  
## 1 abacus   trust  
## 2 abandon  fear  
## 3 abandon  negative  
## 4 abandon  sadness  
## 5 abandoned anger  
## 6 abandoned fear
```

```
sentiment1df <- merge(sentiment, sentiment_dictionary1, by = "word")
head(sentiment1df)
```

```
##           word    ...2  screen_name favorite_count retweet_count media_type
## 1 abominably  72496      peta           78           34      Nophoto
## 2   absence  85636      WWF          220           70      Nophoto
## 3 abundance 122417 AWF_Official          92           24        Photo
## 4 abundance  73701      peta            0            0      Nophoto
## 5 abundance  93624  Defenders          53           23      Nophoto
## 6 abundance  1507      oceana         126           23        Photo
## Postnumber sentiment
## 1         3665  negative
## 2         2455  negative
## 3         2377  positive
## 4         2664  positive
## 5         1056  positive
## 6          407  positive
```

```
library(summarytools)
```

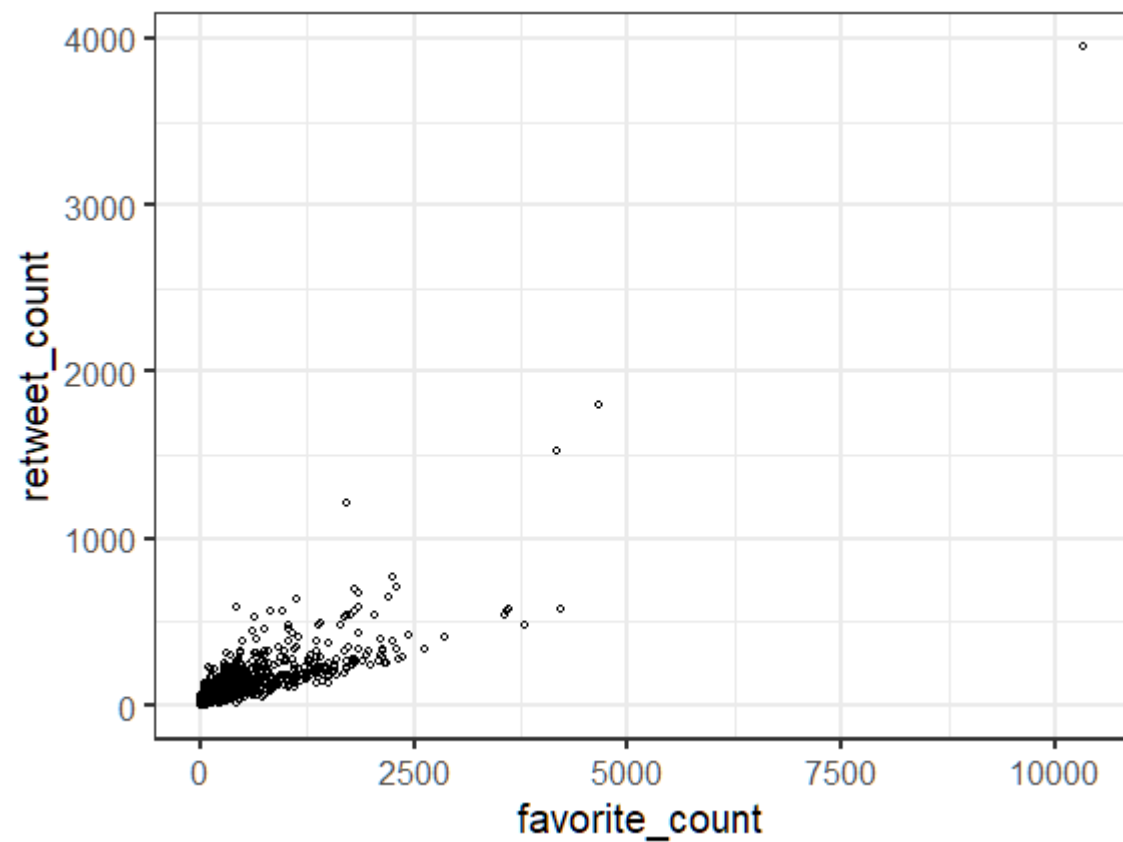
```
##  
## Attaching package: 'summarytools'  
  
## The following object is masked from 'package:tibble':  
##  
##      view
```

```
view(dfSummary(sentiment1df))
```

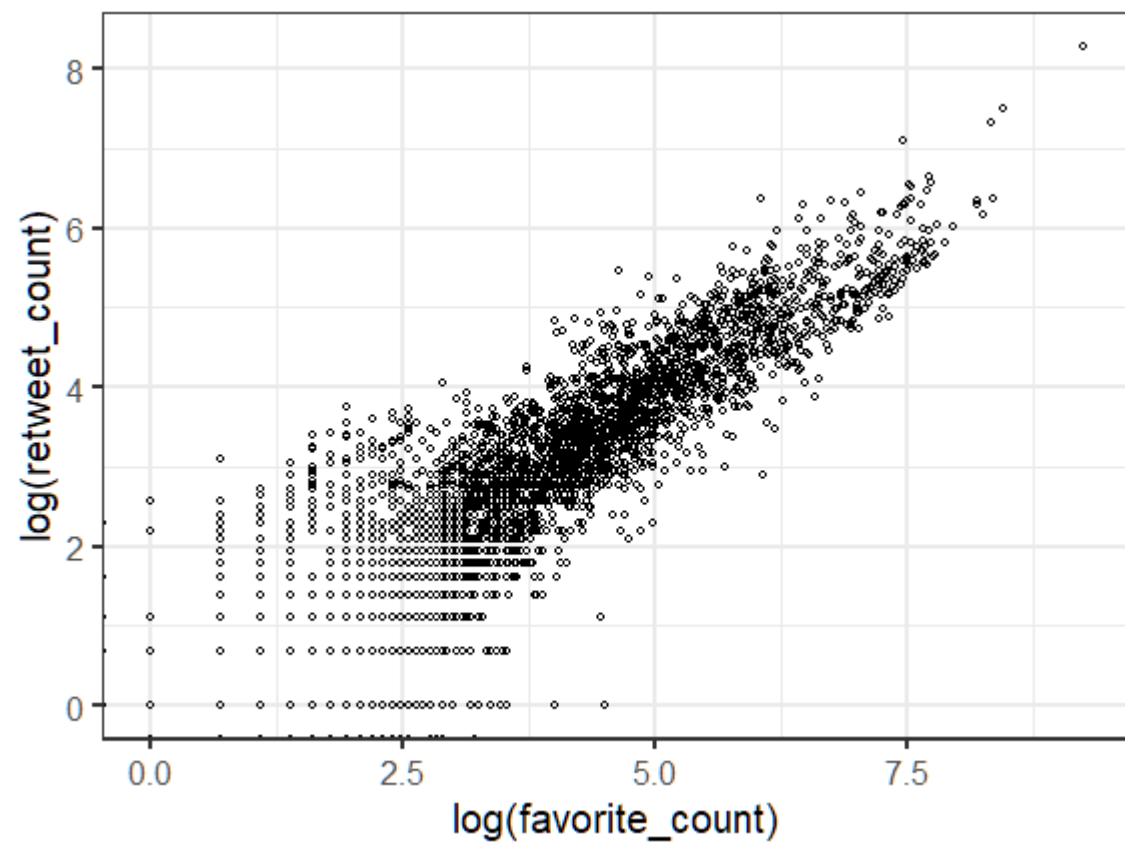
```
## Switching method to 'browser'
```

```
## Output file written: C:\Users\zachs\AppData\Local\Temp\RtmpIX2Iqy\file6c085dd6138b.html
```

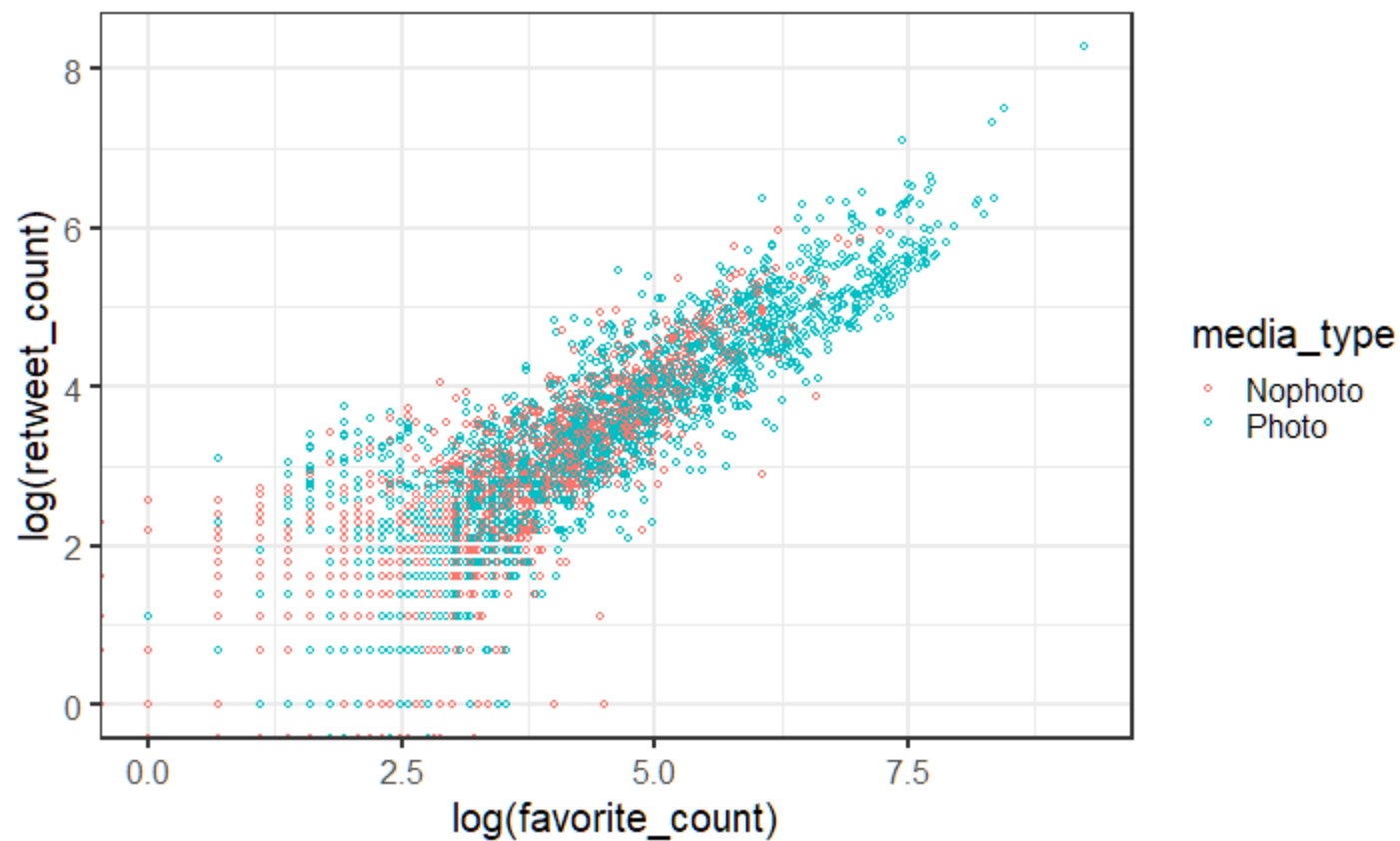
```
ggplot(Exploredf1)+  
  geom_point(aes(y=retweet_count, x = favorite_count), shape=1)+  
  theme_bw(20)
```



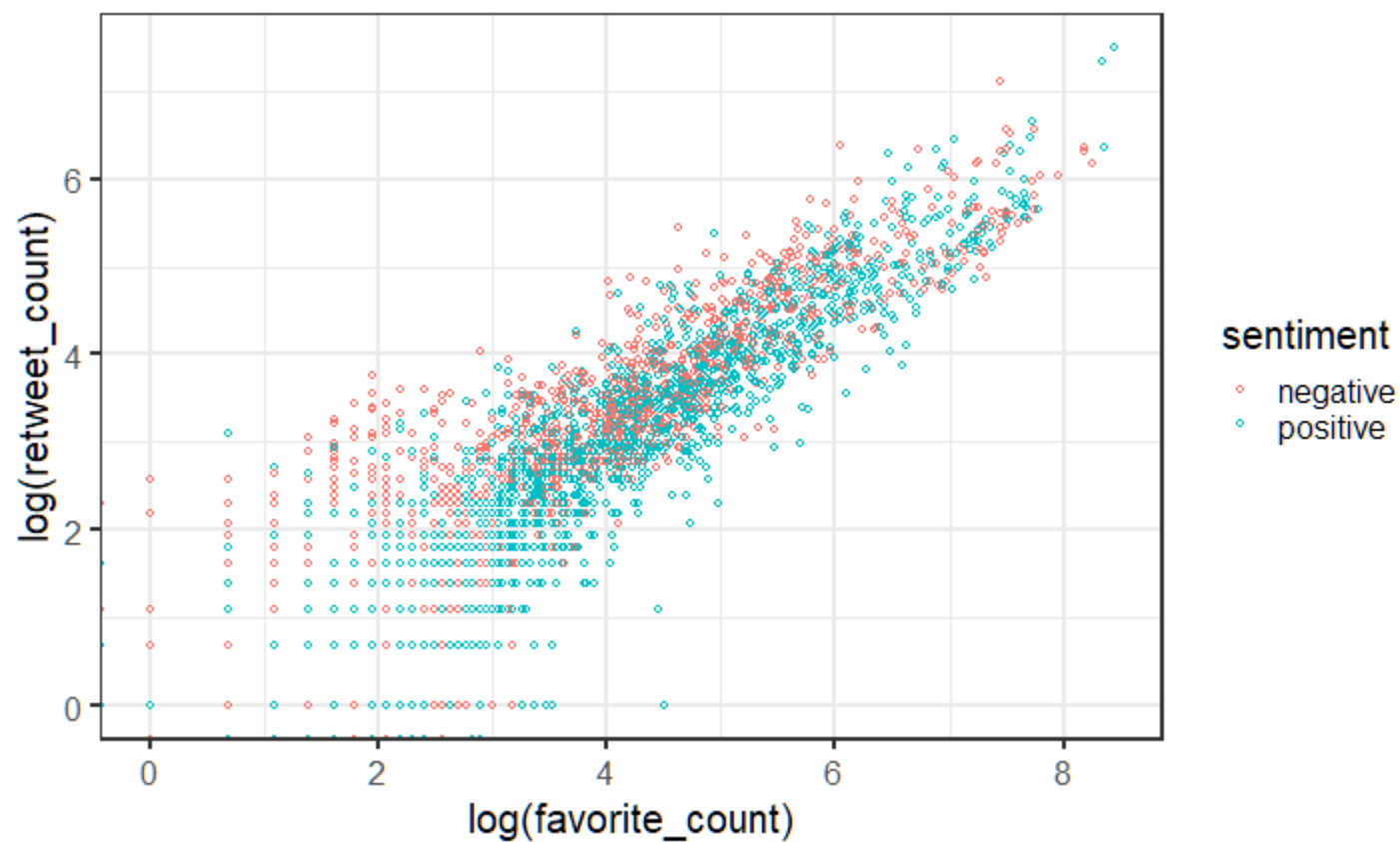
```
ggplot(Exploredf1)+  
  geom_point(aes(y=log(retweet_count), x = log(favorite_count)), shape=1)+  
  theme_bw(20)
```



```
ggplot(Exploredf1)+  
  geom_point(aes(y=log(retweet_count), x = log(favorite_count), colour=media_type), shape=1)+  
  theme_bw(20)
```




```
ggplot(sentiment1df)+  
  geom_point(aes(y=log(retweet_count), x = log(favorite_count), colour=sentiment), shape=1)+  
  theme_bw(20)
```



```
m1 <- lm(retweet_count ~ media_type, data = sentiment1df)
summary(m1)
```

```
##
## Call:
## lm(formula = retweet_count ~ media_type, data = sentiment1df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -68.95  -42.95  -18.95   4.24 1731.05
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    19.762     1.443   13.70  <2e-16 ***
## media_typePhoto  49.191     1.988   24.75  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 89.9 on 8202 degrees of freedom
## Multiple R-squared:  0.06947,    Adjusted R-squared:  0.06936
## F-statistic: 612.4 on 1 and 8202 DF,  p-value: < 2.2e-16
```

```
m2 <- lm(retweet_count ~ media_type + sentiment, data=sentiment1df)
summary(m2)
```

```
##
## Call:
## lm(formula = retweet_count ~ media_type + sentiment, data = sentiment1df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -78.32  -43.76  -13.77   4.67  1736.24
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      28.333      1.873   15.12 < 2e-16 ***
## media_typePhoto    49.984      1.985   25.18 < 2e-16 ***
## sentimentpositive -14.560      2.039   -7.14 1.01e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 89.62 on 8201 degrees of freedom
## Multiple R-squared:  0.07522,    Adjusted R-squared:  0.075
## F-statistic: 333.5 on 2 and 8201 DF,  p-value: < 2.2e-16
```

Time for a break

Welcome Back!

```
library(purrr)
library(modelr)

cv <- crossv_kfold(sentiment1df, k = 10)

cv
```

```
## # A tibble: 10 x 3
##   train                test                .id
##   <named list>         <named list>         <chr>
## 1 <resample [7,383 x 8]> <resample [821 x 8]> 01
## 2 <resample [7,383 x 8]> <resample [821 x 8]> 02
## 3 <resample [7,383 x 8]> <resample [821 x 8]> 03
## 4 <resample [7,383 x 8]> <resample [821 x 8]> 04
## 5 <resample [7,384 x 8]> <resample [820 x 8]> 05
## 6 <resample [7,384 x 8]> <resample [820 x 8]> 06
## 7 <resample [7,384 x 8]> <resample [820 x 8]> 07
## 8 <resample [7,384 x 8]> <resample [820 x 8]> 08
## 9 <resample [7,384 x 8]> <resample [820 x 8]> 09
## 10 <resample [7,384 x 8]> <resample [820 x 8]> 10
```

```
models0 <- map(cv$train, ~lm(retweet_count ~ 1, data = .))  
models1 <- map(cv$train, ~lm(retweet_count ~ media_type, data = .))  
models2 <- map(cv$train, ~lm(retweet_count ~ media_type + sentiment, data = .))
```

```
get_pred <- function(model, test_data){  
  data <- as.data.frame(test_data)  
  pred <- add_predictions(data, model)  
  return(pred)  
}  
  
pred0 <- map2_df(models0, cv$test, get_pred, .id = "Run")  
pred1 <- map2_df(models1, cv$test, get_pred, .id = "Run")  
pred2 <- map2_df(models2, cv$test, get_pred, .id = "Run")
```


Mean Squared Error to assess model fit

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

MSE = mean squared error

n = number of data points

Y_i = observed values

\hat{Y}_i = predicted values

- To find the MSE, take the observed value, subtract the predicted value, and square that difference. Repeat that for all observations. Then, sum all of those squared values and divide by the number of observations.

Interpreting Mean Squared Error

- An MSE of zero, meaning that the estimator $\hat{\theta}$ predicts observations of the parameter θ with perfect accuracy.
- Two or more statistical models may be compared using their MSEs as a measure of how well they explain a given set of observations.
- Mean squared error has the disadvantage of heavily weighting outliers. This property, undesirable in many applications, has led researchers to use alternatives such as the mean absolute error.

```
MSE0 <- pred0 %>% group_by(Run) %>%  
  summarise(MSE = mean((retweet_count - pred)^2))
```

```
MSE0
```

```
## # A tibble: 10 x 2
```

```
##   Run      MSE
```

```
##   <chr> <dbl>
```

```
## 1 1      9043.
```

```
## 2 10     11685.
```

```
## 3 2       7814.
```

```
## 4 3     11480.
```

```
## 5 4       9992.
```

```
## 6 5       5616.
```

```
## 7 6     10028.
```

```
## 8 7       7736.
```

```
## 9 8       8351.
```

```
## 10 9      5103.
```

```
MSE1 <- pred1 %>% group_by(Run) %>%  
  summarise(MSE = mean( (retweet_count - pred)^2))
```

```
MSE1
```

```
## # A tibble: 10 x 2
```

```
##   Run      MSE
```

```
##   <chr> <dbl>
```

```
## 1 1      8396.
```

```
## 2 10     10831.
```

```
## 3 2       7187.
```

```
## 4 3     10678.
```

```
## 5 4       9376.
```

```
## 6 5       5148.
```

```
## 7 6       9413.
```

```
## 8 7       7357.
```

```
## 9 8       7750.
```

```
## 10 9      4697.
```

```
MSE2 <- pred2 %>% group_by(Run) %>%  
  summarise(MSE = mean( (retweet_count - pred)^2))
```

```
MSE2
```

```
## # A tibble: 10 x 2  
##   Run      MSE  
##   <chr> <dbl>  
## 1 1      8324.  
## 2 10     10752.  
## 3 2       7101.  
## 4 3     10607.  
## 5 4       9421.  
## 6 5       5148.  
## 7 6       9366.  
## 8 7       7309.  
## 9 8       7655.  
## 10 9      4671.
```

```
mean(MSE0$MSE)
```

```
## [1] 8684.763
```

```
mean(MSE1$MSE)
```

```
## [1] 8083.495
```

```
mean(MSE2$MSE)
```

```
## [1] 8035.393
```

Time for a break

Welcome Back!

Load new, larger data set

```
Confirmdf <- read_csv("df2.csv")
```

```
## New names:
## * `` -> ...1
## * ...1 -> ...2

## Rows: 20000 Columns: 7
## -- Column specification -----
## Delimiter: ","
## chr (3): screen_name, media_type, text
## dbl (4): ...1, ...2, favorite_count, retweet_count
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
head(Confirmdf)
```

```
## # A tibble: 6 x 7
##   ...1    ...2 screen_name favorite_count retweet_count media_type text
##   <dbl> <dbl> <chr>             <dbl>         <dbl> <chr>    <chr>
## 1      1      1    539 oceana             438         121 Photo    GOOD NEWS~
```

Similar computations as on Exploredf

```
Confirmdf <- Confirmdf%>%  
  mutate(Postnumber = 1:n())%>%  
  select(-c(...1))  
  
confirmsentiment <- Confirmdf %>%  
  unnest_tokens(output = "word", input = "text")
```

```
sentiment_dictionary1 <- get_sentiments("bing")  
head(sentiment_dictionary1)
```

```
## # A tibble: 6 x 2  
##   word      sentiment  
##   <chr>    <chr>  
## 1 2-faces  negative  
## 2 abnormal negative  
## 3 abolish negative  
## 4 abominable negative  
## 5 abominably negative  
## 6 abominate negative
```

```
sentiment_dictionary2 <- get_sentiments("afinn")  
head(sentiment_dictionary2)
```

```
## # A tibble: 6 x 2  
##   word      value  
##   <chr>    <dbl>  
## 1 abandon      -2  
## 2 abandoned    -2  
## 3 abandons     -2  
## 4 abducted     -2  
## 5 abduction    -2  
## 6 abductions   -2
```

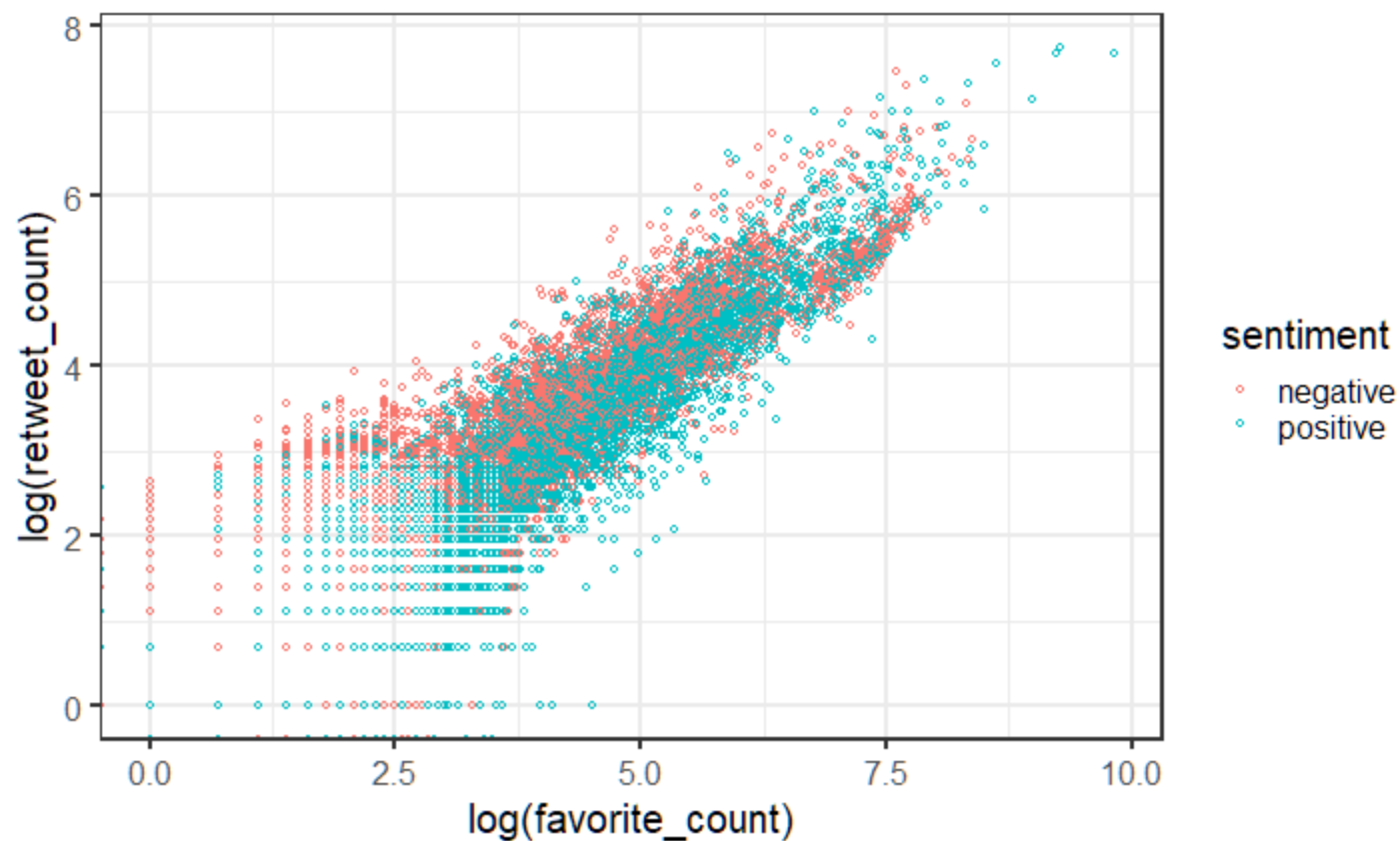
```
sentiment_dictionary3 <- get_sentiments("nrc")  
head(sentiment_dictionary3)
```

```
## # A tibble: 6 x 2  
##   word      sentiment  
##   <chr>    <chr>  
## 1 abacus    trust  
## 2 abandon   fear  
## 3 abandon   negative  
## 4 abandon   sadness  
## 5 abandoned anger  
## 6 abandoned fear
```

```
confirmsentiment1df <- merge(confirmsentiment, sentiment_dictionary1, by = "word")
head(confirmsentiment1df)
```

```
##      word    ...2    screen_name favorite_count retweet_count media_type
## 1 abnormal  70074    HSiGlobal         420          168    Nophoto
## 2 abnormal  80025    MoveTheWorld       315          331     Photo
## 3 abnormal  32344    savingoceans         6           2     Photo
## 4 abolish 112289    Network4Animals       23          21     Photo
## 5 abolish 101041    FarmSanctuary        48          10     Photo
## 6 abound  48231     Greenpeace         46          19    Nophoto
## Postnumber sentiment
## 1      12913  negative
## 2      11822  negative
## 3       6203  negative
## 4       9532  negative
## 5       3559  negative
## 6      16354  positive
```

```
ggplot(confirmsentiment1df)+  
  geom_point(aes(y=log(retweet_count), x = log(favorite_count), colour=sentiment), shape=1)+  
  theme_bw(20)
```



Issues though...

```
predictedvalues <- predict(m1, data=sentiment1df, interval="prediction")
```

```
## Warning in predict.lm(m1, data = sentiment1df, interval = "prediction"): predictions on current data re
```

```
View(predictedvalues)
```

```
v1 <- c(19.76, -156, 196)
v2 <- c(68.95279, -107, 245)
v3 <- c("Nophoto", "Photo")

smallldf <- rbind(v1,v2)

smallldf <- as.data.frame(smallldf)%>%
  rename("meanretweet" = "V1",
         "lower" = "V2",
         "upper" = "V3")

smallldf <- cbind(smallldf, v3)

smallldf <- smallldf %>%
  rename("media_type" = "v3")

sampleConfirmdf <- confirmsentimentldf %>%
  sample_n(400)
```

What's wrong with this picture?

```
ggplot()+  
  geom_jitter(data=sampleConfirmdf, aes(x = media_type, y = retweet_count), width = .05, height=.01)  
  geom_point(data=smalldf, aes(x= media_type, y= meanretweet), shape = 1, colour = "blue")+  
  geom_errorbar(data=smalldf, aes(x= media_type, ymin = lower,ymax=upper), width=.1, colour = "blue")  
  scale_y_continuous(limits = c(-200,500))+  
  theme_bw(20)
```

Warning: Removed 5 rows containing missing values (geom_point).

Impossible values in our error bars

- Need a GLM that takes into account the data are counts!

```
m1poss <- glm(retweet_count ~ media_type, data = sentiment1df, family = "poisson")
summary(m1poss)
```

```
##
## Call:
## glm(formula = retweet_count ~ media_type, family = "poisson",
##      data = sentiment1df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -11.743   -6.287   -4.713    0.718   91.003
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    2.983737   0.003610   826.5  <2e-16 ***
## media_typePhoto 1.249685   0.004048   308.7  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
```

```
m2poss <- glm(retweet_count ~ media_type + sentiment , data = sentiment1df, family = "poisson")
summary(m2poss)
```

```
##
## Call:
## glm(formula = retweet_count ~ media_type + sentiment, family = "poisson",
##      data = sentiment1df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -12.911    -6.852    -4.662     0.719    93.321
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.155850   0.003983   792.2  <2e-16 ***
## media_typePhoto  1.267198   0.004053   312.7  <2e-16 ***
## sentimentpositive -0.312762   0.003292   -95.0  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 770433  on 8203  degrees of freedom
## Residual deviance: 644619  on 8201  degrees of freedom
## AIC: 677966
##
## Number of Fisher Scoring iterations: 6
```

We've conducted a bunch of exploratory analyses!

- We seem to find that both media type and the specific sentiment of the post predict retweeting!
- Let's turn to conducting a purely confirmatory analysis. How should we do that?

First, let's consider our prior findings and then write down our predictions.

- Tweets with photos were retweeted more not just on training data but also on test data
- Sentiment of a tweet, specifically negative sentiments, seemed to predict more retweeting as well (also on test data)

Issues:

- We used one metric of sentiment. Ideally, our findings should hold for other metrics of sentiment. *We should predict they will.*
- We've focused on retweeting entirely but we also see that retweets and favorites are extremely strongly correlated. *We should predict all of the same predictions will hold for favoriting just like retweeting. Or we need a good reason to distinguish them*
- We initially fit linear models but it's pretty clear those models are problematic. Need to fit a poisson model.
- Now let's write down our models that correspond to these hypotheses.

Initial measure of sentiment with poisson model

```
confirm.m2.poss.rt <- glm(retweet_count ~ media_type + sentiment , data = confirmsentiment1df, family = "poisson")
summary(confirm.m2.poss.rt)
```

```
##
## Call:
## glm(formula = retweet_count ~ media_type + sentiment, family = "poisson",
##      data = confirmsentiment1df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -12.182   -6.568   -4.642    0.860   111.759
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.071279   0.002021  1519.3  <2e-16 ***
## media_typePhoto  1.235566   0.002041   605.5  <2e-16 ***
## sentimentpositive -0.220538   0.001694  -130.2  <2e-16 ***
```

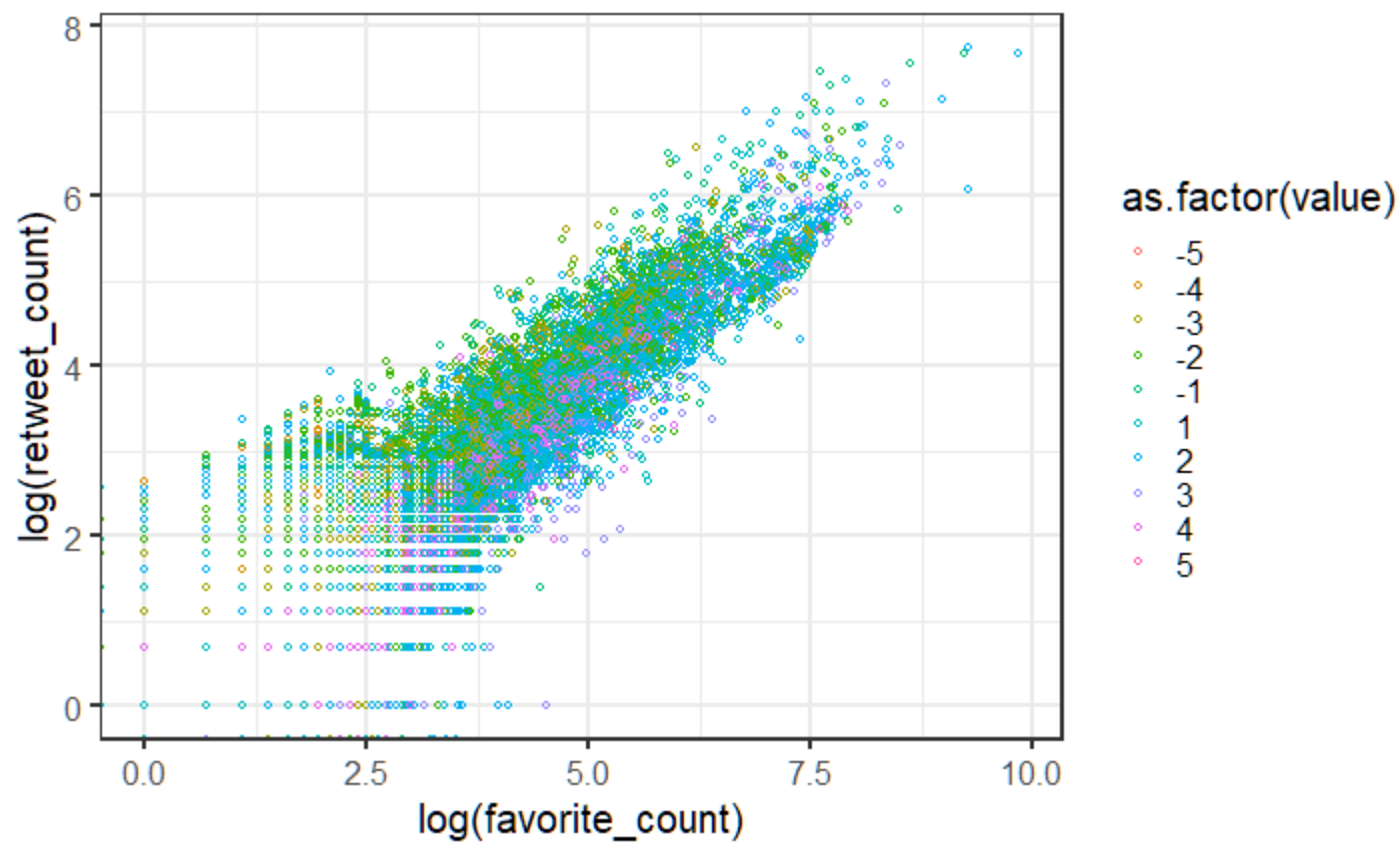
```
confirm.m2.poss.fav <- glm(favorite_count ~ media_type + sentiment , data = confirmsentiment1df, family = "poisson")
summary(confirm.m2.poss.fav)
```

```
##
## Call:
## glm(formula = favorite_count ~ media_type + sentiment, family = "poisson",
##      data = confirmsentiment1df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -22.13  -14.33   -8.14   -0.25   352.86
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.7995430   0.0013039  2914.06  <2e-16 ***
## media_typePhoto    1.7010002   0.0012809  1327.97  <2e-16 ***
## sentimentpositive -0.0108747   0.0009329   -11.66  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 13291568  on 33248  degrees of freedom
## Residual deviance: 10800231  on 33246  degrees of freedom
## AIC: 10965750
```

```
confirmsentiment2df <- merge(confirmsentiment, sentiment_dictionary2, by = "word")
head(confirmsentiment2df)
```

```
##      word    ...2    screen_name favorite_count retweet_count media_type
## 1 abandon  68029    BornFreeFDN             0             0      Nophoto
## 2 abandon  92123      Defenders             72            34      Nophoto
## 3 abandon 118030    Animals1st             49            26      Nophoto
## 4 abandon  43511             350             14             7      Nophoto
## 5 abandon 119983  SheldrickTrust          2425           356        Photo
## 6 abandon  13551      whalesorg            151           106        Photo
##   Postnumber value
## 1         1245   -2
## 2        12522   -2
## 3         2593   -2
## 4         5336   -2
## 5          128   -2
## 6        12715   -2
```

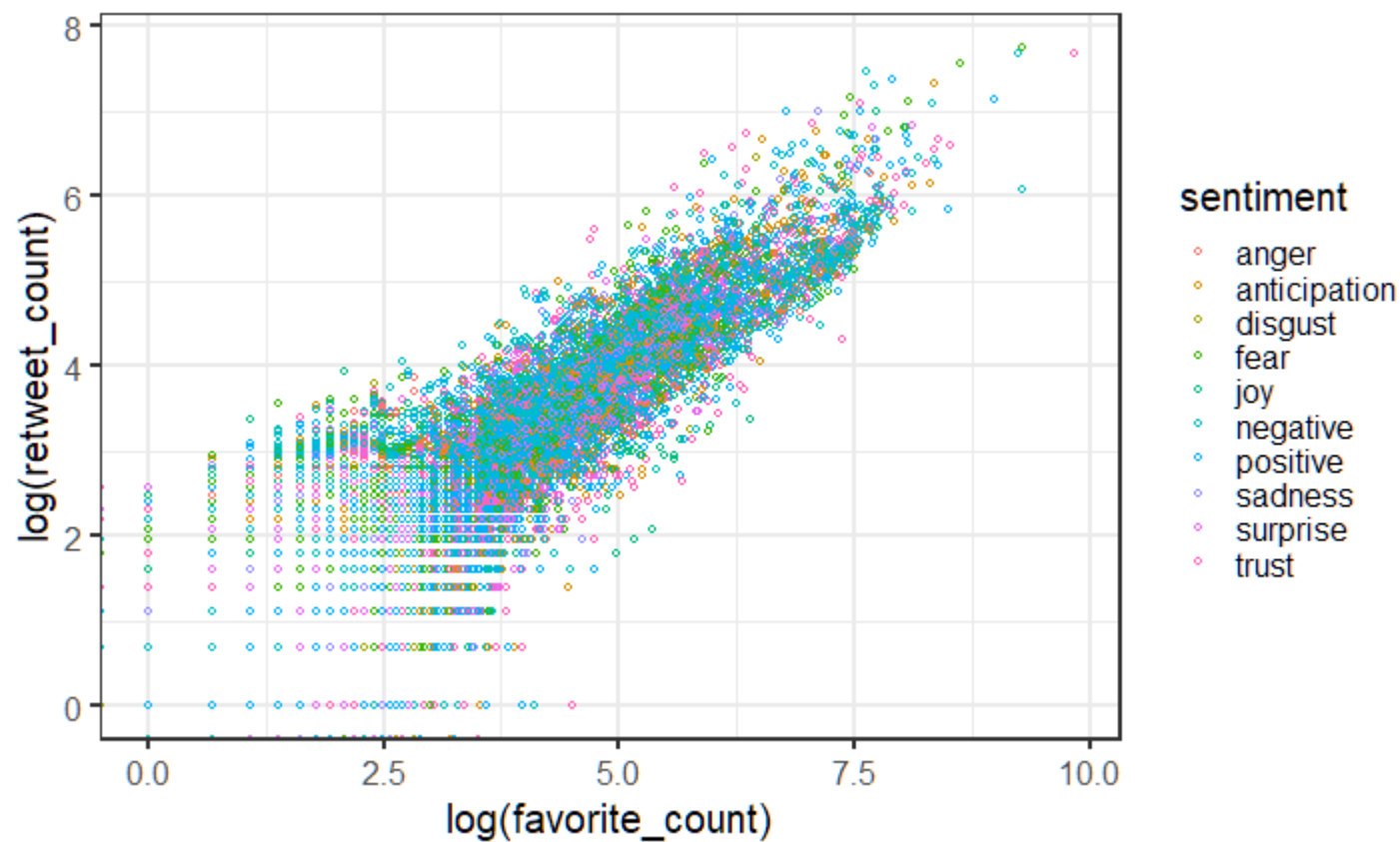
```
ggplot(confirmsentiment2df)+  
  geom_point(aes(y=log(retweet_count), x = log(favorite_count), colour=as.factor(value)), shape=1)+  
  theme_bw(20)
```



```
confirmsentiment3df <- merge(confirmsentiment, sentiment_dictionary3, by = "word")  
  
head(confirmsentiment3df)
```

```
##      word  ...2 screen_name favorite_count retweet_count media_type Postnumber  
## 1 abandon 68029 BornFreeFDN             0              0    Nophoto        1245  
## 2 abandon 68029 BornFreeFDN             0              0    Nophoto        1245  
## 3 abandon 68029 BornFreeFDN             0              0    Nophoto        1245  
## 4 abandon 43511           350            14             7    Nophoto        5336  
## 5 abandon 43511           350            14             7    Nophoto        5336  
## 6 abandon 43511           350            14             7    Nophoto        5336  
##      sentiment  
## 1  negative  
## 2      fear  
## 3  sadness  
## 4  negative  
## 5      fear  
## 6  sadness
```

```
ggplot(confirmsentiment3df)+  
  geom_point(aes(y=log(retweet_count), x = log(favorite_count), colour=sentiment), shape=1)+  
  theme_bw(20)
```



```
confirm.m2.poss.rt.s2 <- glm(retweet_count ~ media_type + value , data = confirmsentiment2df, family = "poisson")
summary(confirm.m2.poss.rt.s2)
```

```
##
## Call:
## glm(formula = retweet_count ~ media_type + value, family = "poisson",
##      data = confirmsentiment2df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -12.243    -6.418    -4.718     0.536   112.026
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    2.9763279   0.0016787   1773.0  <2e-16 ***
## media_typePhoto  1.1942001   0.0018970    629.5  <2e-16 ***
## value          -0.0487216   0.0004004   -121.7  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 3458813  on 38550  degrees of freedom
## Residual deviance: 2976550  on 38548  degrees of freedom
## AIC: 3127612
```

```
confirm.m2.poss.fav.s2 <- glm(favorite_count ~ media_type + value , data = confirmsentiment2df, fam
summary(confirm.m2.poss.fav.s2)
```

```
##
## Call:
## glm(formula = favorite_count ~ media_type + value, family = "poisson",
##      data = confirmsentiment2df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -22.11  -13.43   -8.41   -0.62  355.71
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.7768794  0.0011059  3415.2  <2e-16 ***
## media_typePhoto 1.5934463  0.0011920  1336.8  <2e-16 ***
## value          0.0320925  0.0002331   137.7  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 14333425  on 38550  degrees of freedom
## Residual deviance: 11873546  on 38548  degrees of freedom
## AIC: 12059771
```



```
confirm.m2.poss.rt.s3 <- glm(retweet_count ~ media_type + sentiment , data = confirmsentiment3df, fa  
summary(confirm.m2.poss.rt.s3)
```

```
##  
## Call:  
## glm(formula = retweet_count ~ media_type + sentiment, family = "poisson",  
##      data = confirmsentiment3df)  
##  
## Deviance Residuals:  
##      Min        1Q    Median        3Q        Max   
## -12.012    -6.846    -4.387     0.909    112.548   
##  
## Coefficients:  
##              Estimate Std. Error z value Pr(>|z|)      
## (Intercept)      3.142474   0.002010 1563.574 < 2e-16 ***  
## media_typePhoto      1.082290   0.001029 1052.293 < 2e-16 ***  
## sentimentanticipation -0.132585   0.002271  -58.384 < 2e-16 ***  
## sentimentdisgust      0.030135   0.002806   10.740 < 2e-16 ***  
## sentimentfear         0.040933   0.002338   17.511 < 2e-16 ***  
## sentimentjoy        -0.097065   0.002298  -42.245 < 2e-16 ***  
## sentimentnegative     0.042245   0.002184   19.340 < 2e-16 ***  
## sentimentpositive    -0.177523   0.002057  -86.312 < 2e-16 ***  
## sentimentsadness      0.053854   0.002543   21.179 < 2e-16 ***  
## sentimentsurprise    -0.018466   0.002607   -7.083 1.41e-12 ***  
## sentimenttrust      -0.150338   0.002176  -69.088 < 2e-16 ***
```

```
confirm.m2.poss.fav.s3 <- glm(favorite_count ~ media_type + sentiment , data = confirmsentiment3df,  
summary(confirm.m2.poss.fav.s3)
```

```
##  
## Call:  
## glm(formula = favorite_count ~ media_type + sentiment, family = "poisson",  
##      data = confirmsentiment3df)  
##  
## Deviance Residuals:  
##      Min        1Q    Median        3Q        Max   
## -22.85   -14.60    -8.28    -0.09   357.77   
##  
## Coefficients:  
##              Estimate Std. Error z value Pr(>|z|)      
## (Intercept)    3.8165626  0.0012317 3098.69  <2e-16 ***  
## media_typePhoto  1.5277198  0.0006477 2358.61  <2e-16 ***  
## sentimentanticipation 0.0720990  0.0013125  54.93  <2e-16 ***  
## sentimentdisgust   -0.0024745  0.0016946  -1.46    0.144   
## sentimentfear      0.1019829  0.0013825  73.77  <2e-16 ***  
## sentimentjoy       0.2206427  0.0013033 169.30  <2e-16 ***  
## sentimentnegative  0.1388681  0.0012892 107.71  <2e-16 ***  
## sentimentpositive  0.0502798  0.0012052  41.72  <2e-16 ***  
## sentimentsadness   0.1608966  0.0014839 108.43  <2e-16 ***  
## sentimentsurprise  0.2185901  0.0014712 148.58  <2e-16 ***  
## sentimenttrust     0.0597045  0.0012656  47.17  <2e-16 ***
```


Summary

- We began with very open ended questions. We first discussed weaknesses of standard exploring modeling.
- Then we learned about folding data to make sure our predictions were more likely to hold.
- We used some of these strategies on a Twitter dataset
- Then we attempted to confirm our initial hypotheses we formed during an exploratory modeling phase.
- Many but not all of our hypotheses, or claims entailed by our hypothesis, were confirmed. But we also considered the fact that the meaningfulness of these effects may be questionable.

