

Assignment 2

Conner McCloney

1

a)

This dataset is a collection of books that can be found on Goodreads, which is a company considers itself ‘the world’s largest site for readers and book recommendations’, that allows users to browse its database of books, rate them, interact with other users on their platform, and other related activities. This dataset records the title of the book, its author(s), ISBN numbers, number of pages, language it’s written in, average rating, number of ratings, and number of text reviews given.

The research question I used last week was ‘Can we predict the average rating of a book from Goodreads?’ In other words, I’m curious if any of these variables, besides the title of the book, can be used as predictors to accurately model the average rating of a book as a response variable.

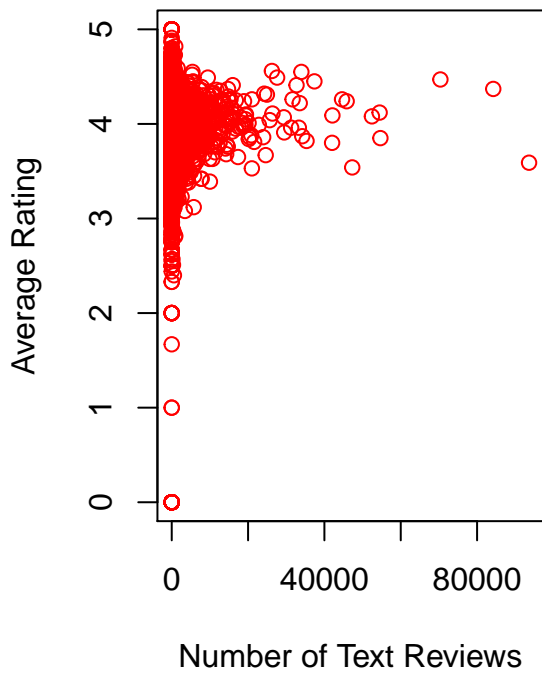
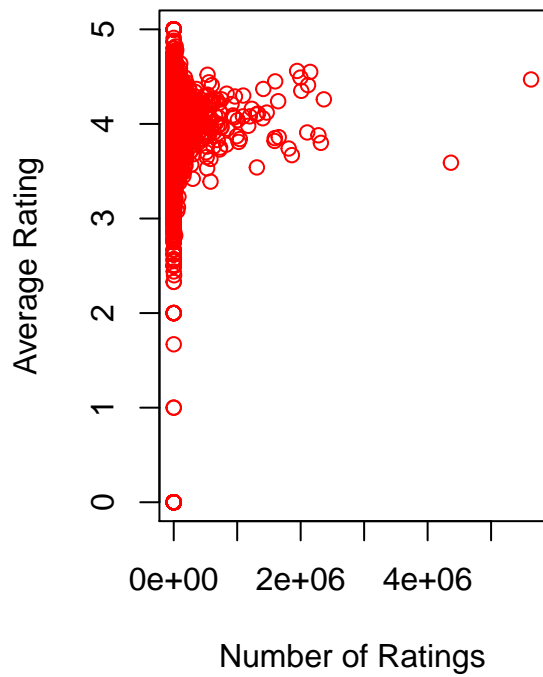
The plots below model a few of these variables against the average rating, and all of them seem to have the same distributional shape, centering around an average rating of 4. If the value of the predictor is low, there seems to be no relationship between it and the average rating, but as you increase the value, the average rating will become more and more likely to be 4. I find it surprising that this relationship also exists between average rating and the number of pages in a book as well, which is something I didn’t expect to see, and that according to the population of Goodreads users, it seems the more attention a book receives, as measured by the number of reviews and ratings given, the more the average rating will converge to a ~4/5 rating, which is interesting.

Source: [<https://www.kaggle.com/jealousleopard/goodreadsbooks>]

```
books_data <- read.csv("books.csv")
par(mfrow=c(1,2))
plot(as.numeric(as.character(average_rating))~as.numeric(as.character(ratings_count)),data=books_data, y

## Warning in eval(predvars, data, env): NAs introduced by coercion
plot(as.numeric(as.character(average_rating))~as.numeric(as.character(text_reviews_count)),data=books_d

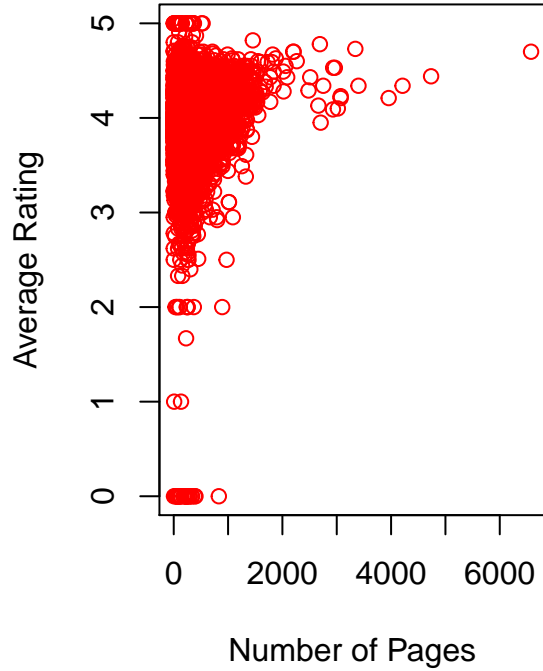
## Warning in eval(predvars, data, env): NAs introduced by coercion
```



```
plot(as.numeric(as.character(average_rating))~as.numeric(as.character(X..num_pages)),data=books_data, y=
```

```
## Warning in eval(predvars, data, env): NAs introduced by coercion
```

```
## Warning in eval(predvars, data, env): NAs introduced by coercion
```



b)

This dataset is a sample of pulsar candidates measured during the High Time Resolution Universe Survey in 2011. A pulsar is a rare type of rotating Neutron star that emits electromagnetic radiation, and can only be measured when this beam is aimed directly at Earth. These pulsars were searched for using large radio telescopes looking for periodic radio signals that a pulsar would produce. Many measurements taken that could be a pulsar, known as a candidate, are recorded, however in practice most observations recorded are caused by radio frequency interference (RFI) and noise. This dataset records the mean, standard deviation, skewness, and excess kurtosis of the integrated profile, which is an array of variables that describe the signal recorded, and of the DM-SNR curve. The DM-SNR curve, which stands for Dispersion Measure - Signal-to-Noise Ratio, describes the relationship between the two for the observed signal, where a curve whose SNR peaks at a DM of zero is likely RFI, or if it is a legitimate signal it should peak at a DM greater than zero. Lastly, there is the `target_class` variable, which is a binary classification variable of whether the given signal was truly a pulsar or not.

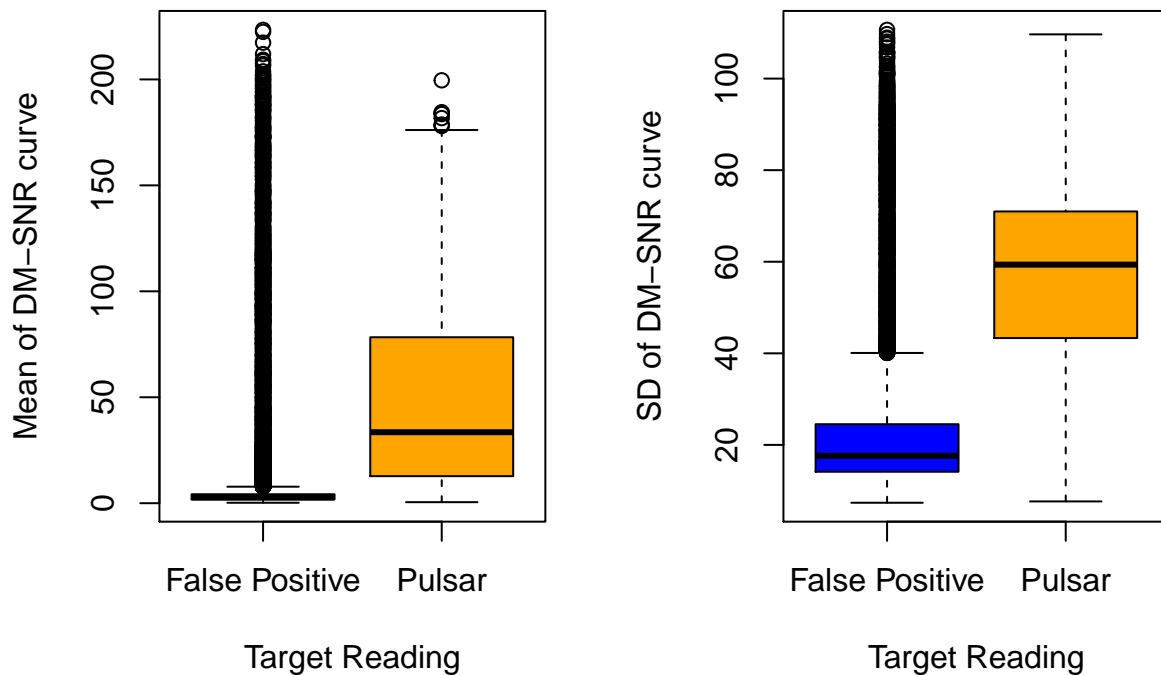
The research question I used last week was ‘Is there a relationship between the dispersion measure signal-to-noise ratio (DM-SNR) curve and the classification of an observed measurement as a pulsar or a false positive?’, which the plots below were used to explore, although extending this question to include the integrated profile would be equally interesting.

The plots below model each of the characteristics of the DM-SNR curve against the classification of whether a reading is a pulsar or a false positive. It seems fairly clear from the first two plots that the mean and SD of the curve tend to be larger if the signal is a legitimate pulsar signal, although the false positive variation, especially for the mean of the curve, are much smaller, resulting in a high number of outliers according to these plots. For the remaining two plots, the reverse seems to be true, as the excess kurtosis and skewness of

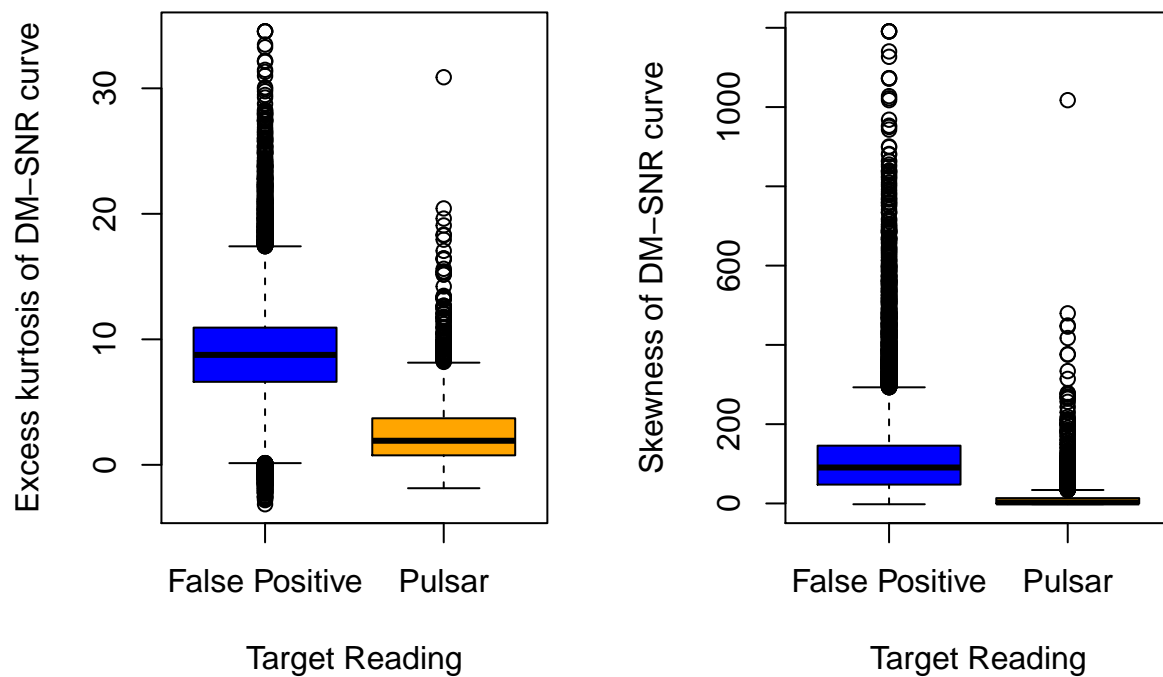
the curve tend to be larger if the signal is a false positive. The variation in the kurtosis plot seems to be roughly equal for the two readings, although it is very small for true pulsars in the skewness plot.

Source: [https://www.kaggle.com/pavanraj159/predicting-a-pulsar-star]

```
pulsar_data <- read.csv("pulsar_stars.csv")
pulsar_data$target_class[pulsar_data$target_class==1] <- "Pulsar"
pulsar_data$target_class[pulsar_data$target_class==0] <- "False Positive"
par(mfrow=c(1,2))
boxplot(Mean.of.the.DM.SNR.curve~target_class,data=pulsar_data,ylab="Mean of DM-SNR curve", xlab="Target Reading")
boxplot(Standard.deviation.of.the.DM.SNR.curve~target_class,data=pulsar_data, ylab="SD of DM-SNR curve")
```



```
par(mfrow=c(1,2))
boxplot(Excess.kurtosis.of.the.DM.SNR.curve~target_class,data=pulsar_data, ylab="Excess kurtosis of DM-SNR curve", xlab="Target Reading")
boxplot(Skewness.of.the.DM.SNR.curve~target_class,data=pulsar_data, ylab="Skewness of DM-SNR curve", xlab="Target Reading")
```



3

The KNN classifier examines the K-nearest neighbors of a given point, and classifies the response value of it based on the mode of those K-nearest neighbors. The KNN regression model examines the K-nearest neighbors of a given point x_0 , and generates an estimate for $f(x_0)$ based on the mean of the response values of those points. This will generate $\hat{f}(X)$, in an attempt to find the true relationship between the response and the predictor variables.

4

```
library(MASS)
library(ISLR)
```