code cademy

OKCupid – Date –A- Scientist

Cletus Norton
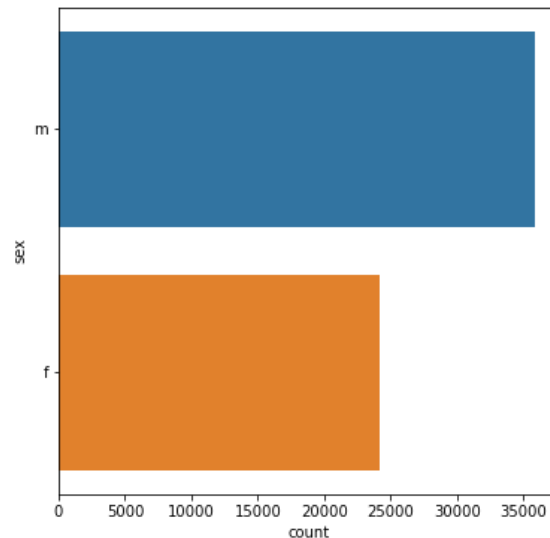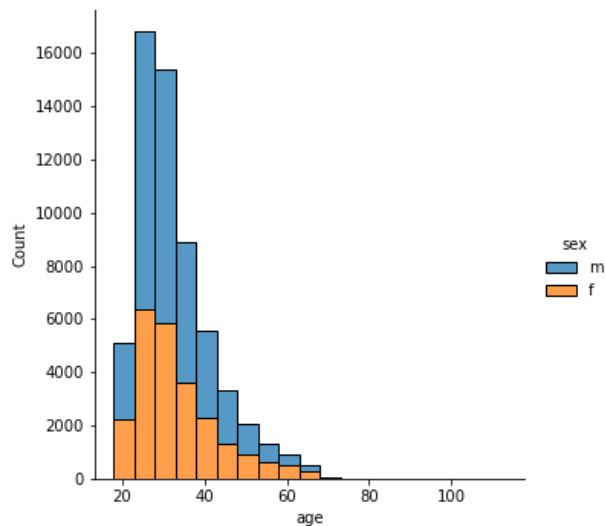Nov. 28 2020

# 1.1 Information about the data

This is a capstone project for the data science path, where we are to analyze data from an on-line dating application called OKCupid. As phones have become more computer-like there has been a rise in dating apps to find love. Many of these apps use data science techniques to recommend possible matches in hope to connect users to their "true Love". The amount of data collected via the app gives us a large amount of personal information about human romance.

The project goal is to try and perdict the astrological sign of the OKCupid users by using the other varibles provided. This is a fun use of the data since many users find astrological signs as a fun aspect of dating. Since the data provided doesn't include any birth imformation beyond age, we developed a means to predict the correct sign.

Using the data provided in "profiles.csv", I present a solution that will use descriptive statistics and data vizualization to find important features to understand the counts, relationship and distribution between the features. Supervised learning models of machine learning were used to make the predictions. Final analysis of the predictions will be checked via accuracy metrics and confusion matix.
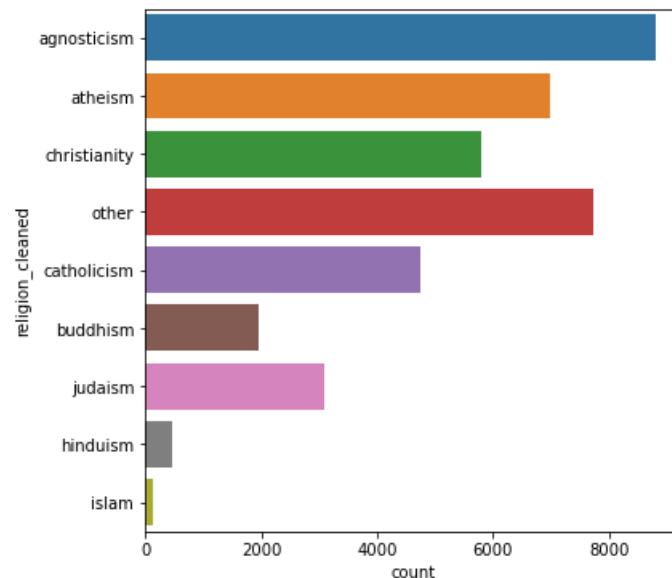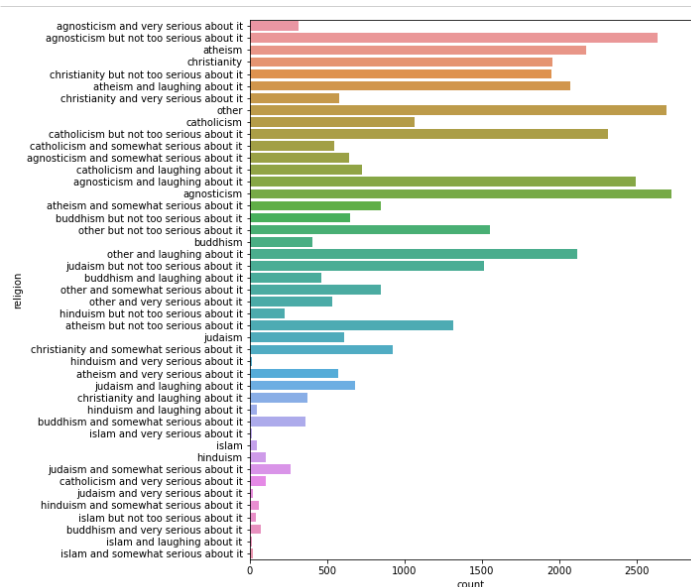
# 1.2 Dataset detail

- Count of the dataset of men and women
- Compares age between the sexes
- Most of both groups are in mid-twenties to late thirties

- Count showing the number of men versus women
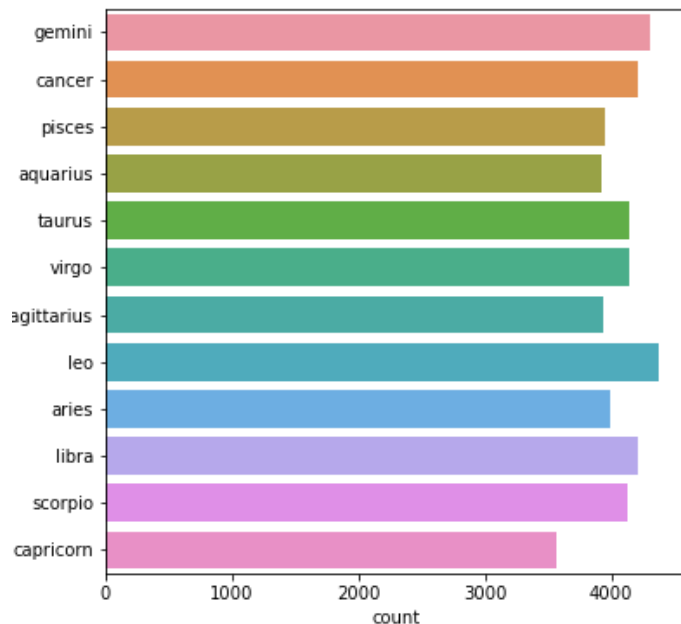- There are about 10,000 more men in the dataset

# 1.3 Dataset detail

- Religion responses
- Several had comments about strength of belief

- Count of religion response after clean up
- Comments about strength of belief were removed

# 1.4 Astrological signs

- This show the counts of the astrological signs after removing comments about interest/belief
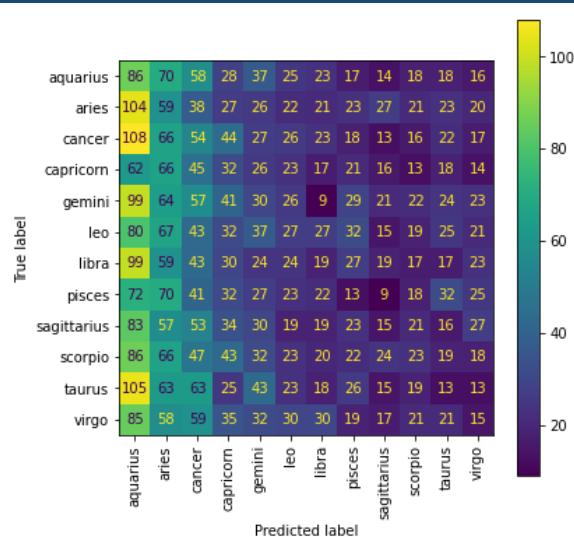- This shows that we have a proportional count on each sign

In my effort to determine the astrological sign based on variables provided, I tried two supervised machine learning approaches. The first was the K-nearest neighbor, while the second was using skLearn's DecisionTree. In the next couple of slides I will go over the details of how each faired. After looking over the variables, I chose the following to use in my machine learning models:

Body Type, diet, orientation, education, religion, sex,  & job.

I chose to leave, drinking, drug use, ethnicity, height, income, job, location, offspring, pets, smokes,  language spoken, & the essay answers out of the model.
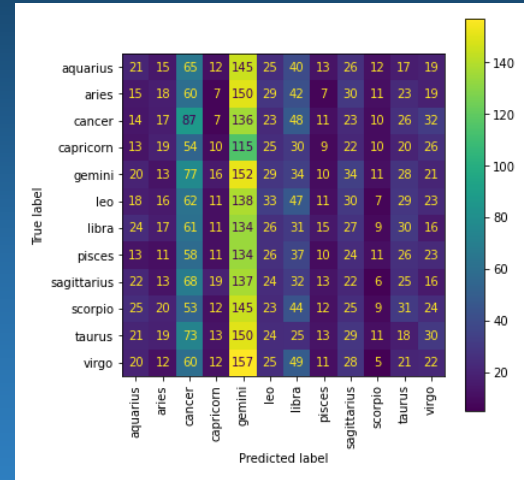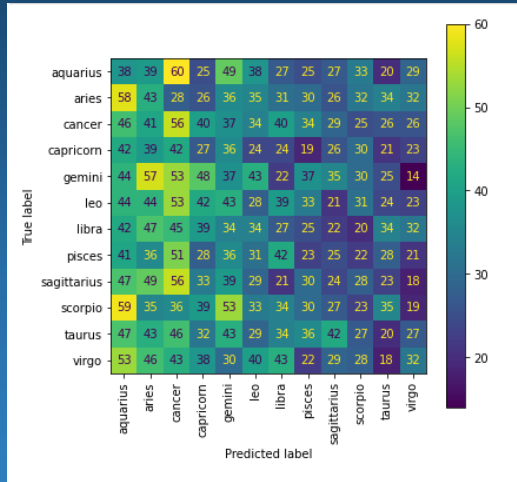
# 1.5 K-Nearest Neighbor

- The data was split into 75% train and 25% test.
- The model showed a mean accuracy of 32%.
- The confusion matrix plot below shows a rather uneven success rate.

# 1.6 DecisionTree model

- The chart on the left shows the decision-tree model used with no depth limit, resulting in a mean accuracy of 73%.
- The chart on the right shows the decision-tree model used with a depth limit of 20, resulting in a mean accuracy of 28%.
- Without limiting the tree depth (which originally was 64) the decision tree overfit the data.
- When we begin to limit the depth the accuracy dropped below the KNN method.

In conclusion this was a fun lesson to apply some basic machine learning to, even though the model ended up no better than guessing. There may have been better data in the essay answers but finding a way to quantify them proved outside my current skills.