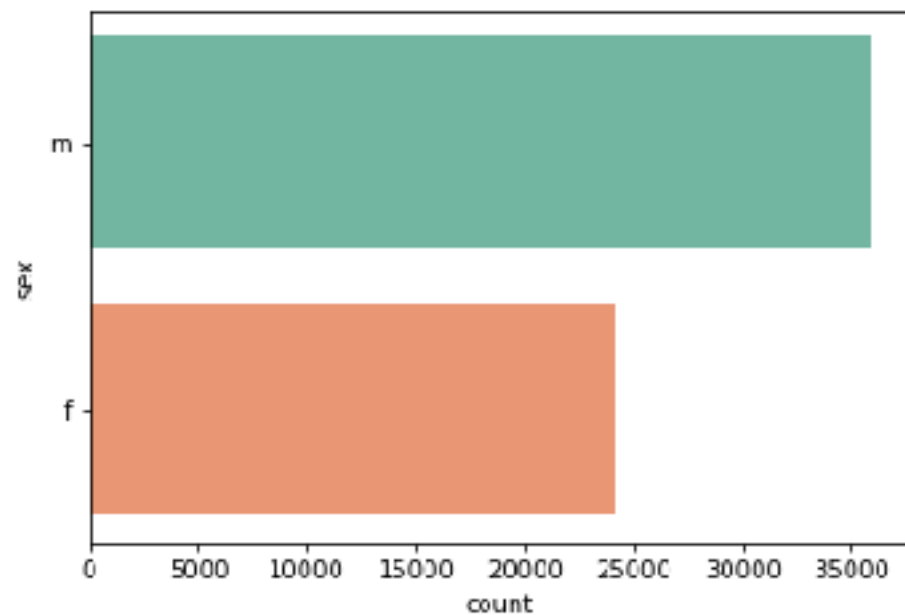# OKCUPID

Date a Scientist Project
By Teodora

This is a Codecademy capstone project for the data science path, where we have to analyse data from the online dating application called OKCupid.

The objective is to try and predict the astrological sign of the OKCupid users by using the other variables provided, and applying Machine Learning algorithms.
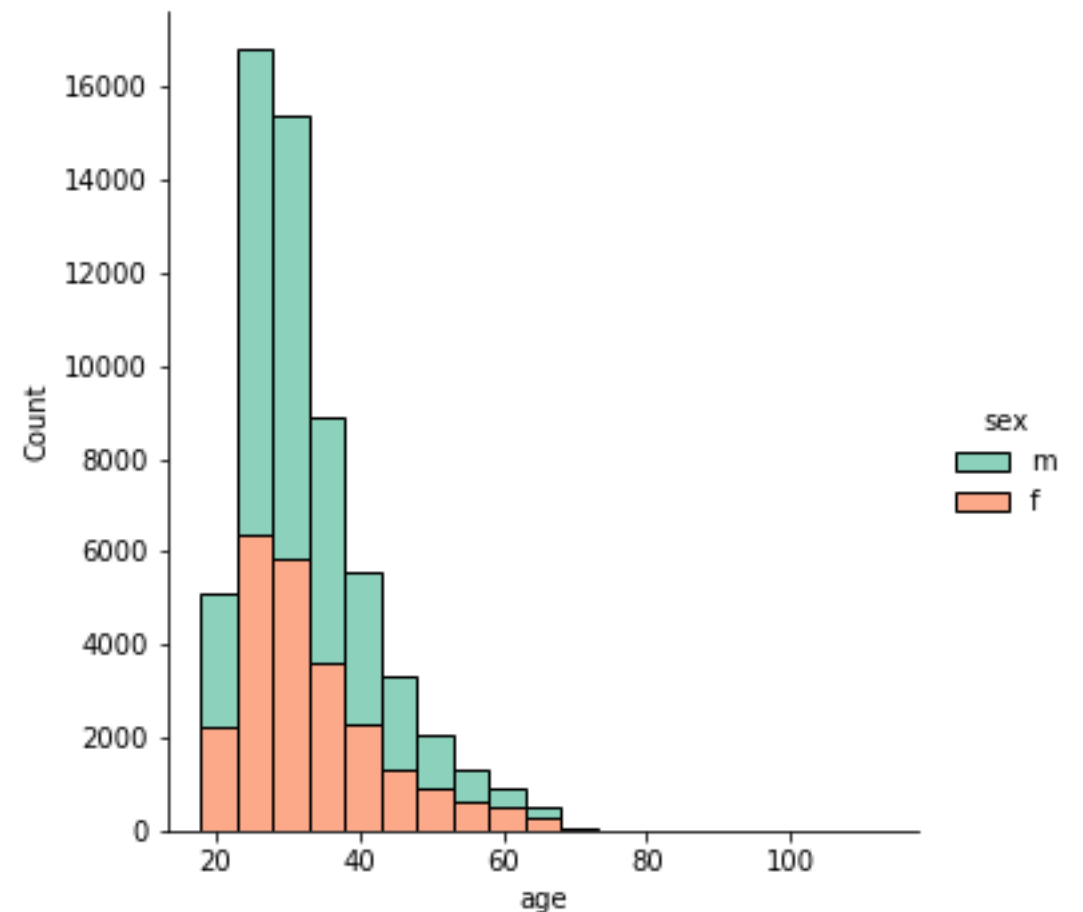
I have used K Nearest Neighbors, Decision Tree Classifier and Support Vector Machine.

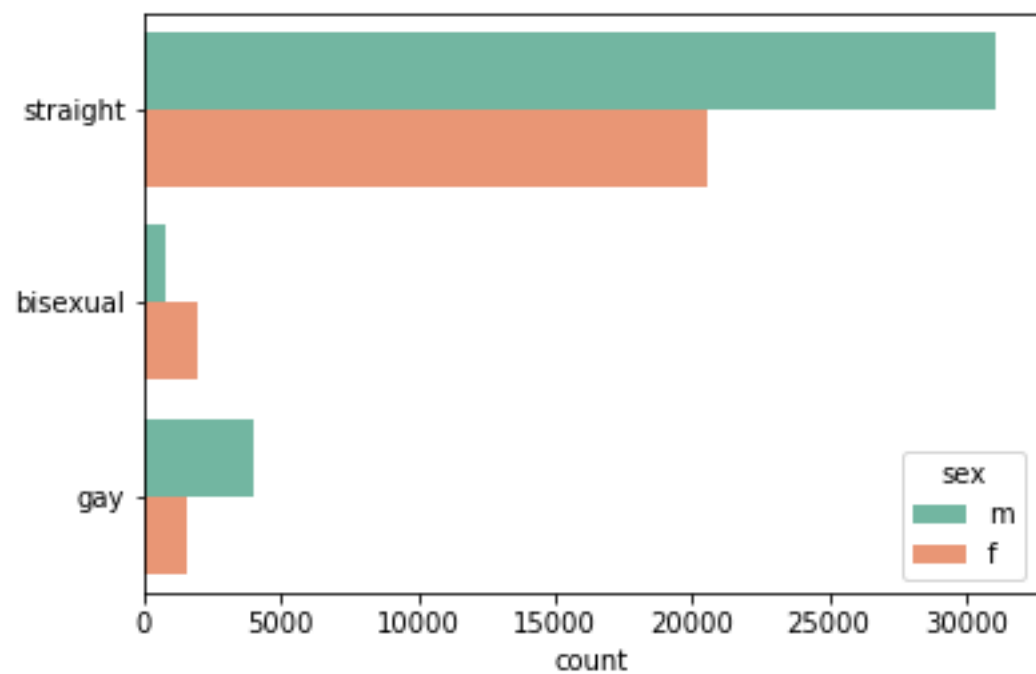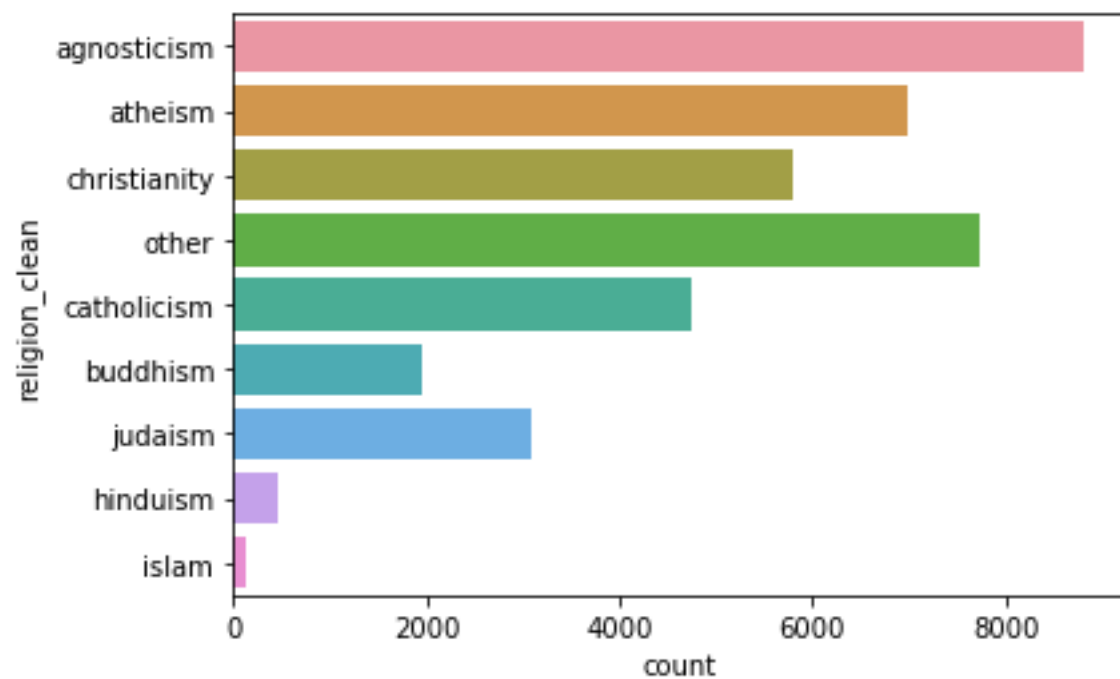But first, let's have a look at the data.

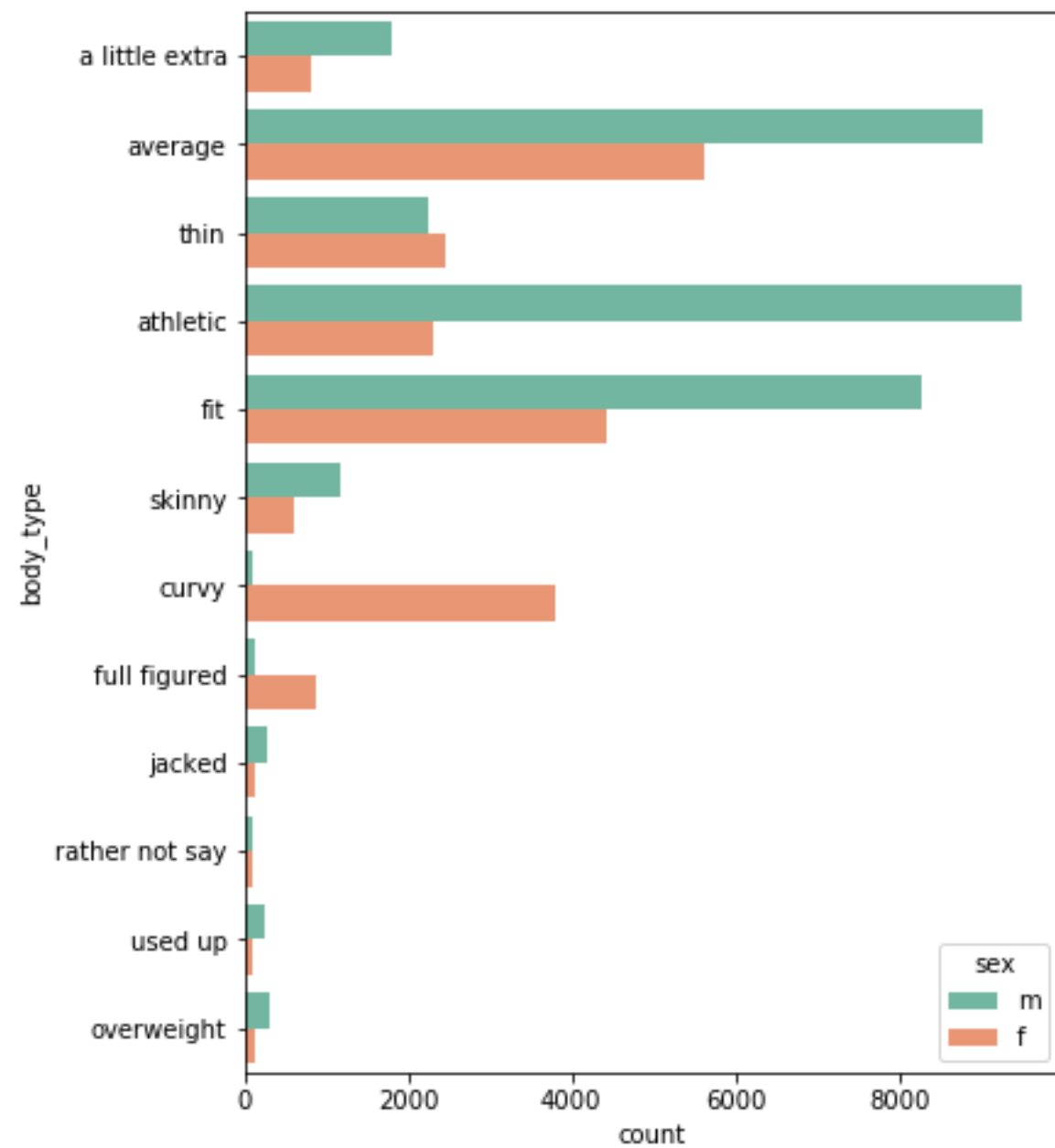We have about 10.000 more men than women in our dataset.



The average age is 32 years, and as we can see, most users are in the 20 - 40 years age range.
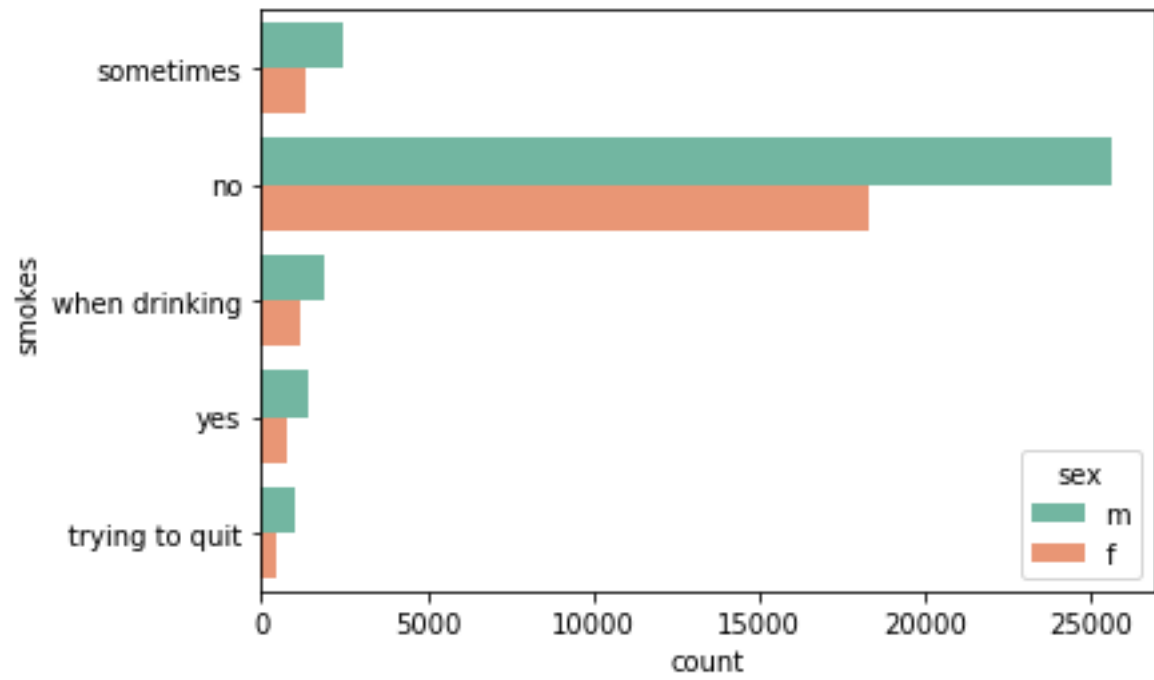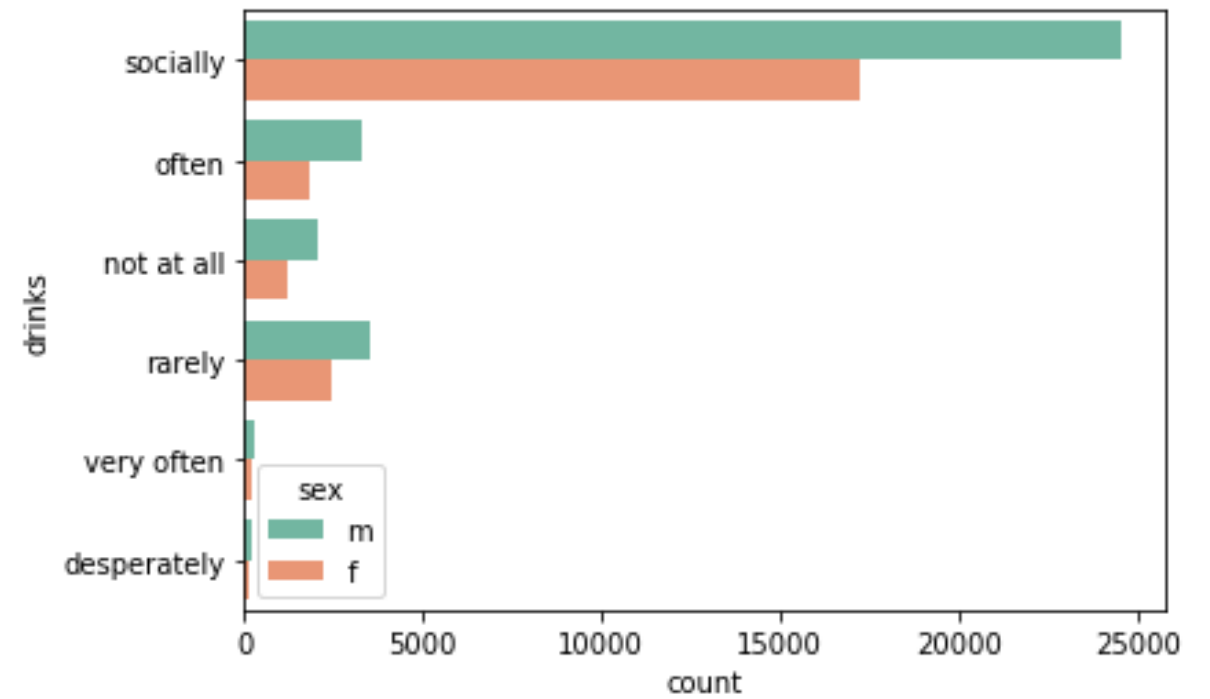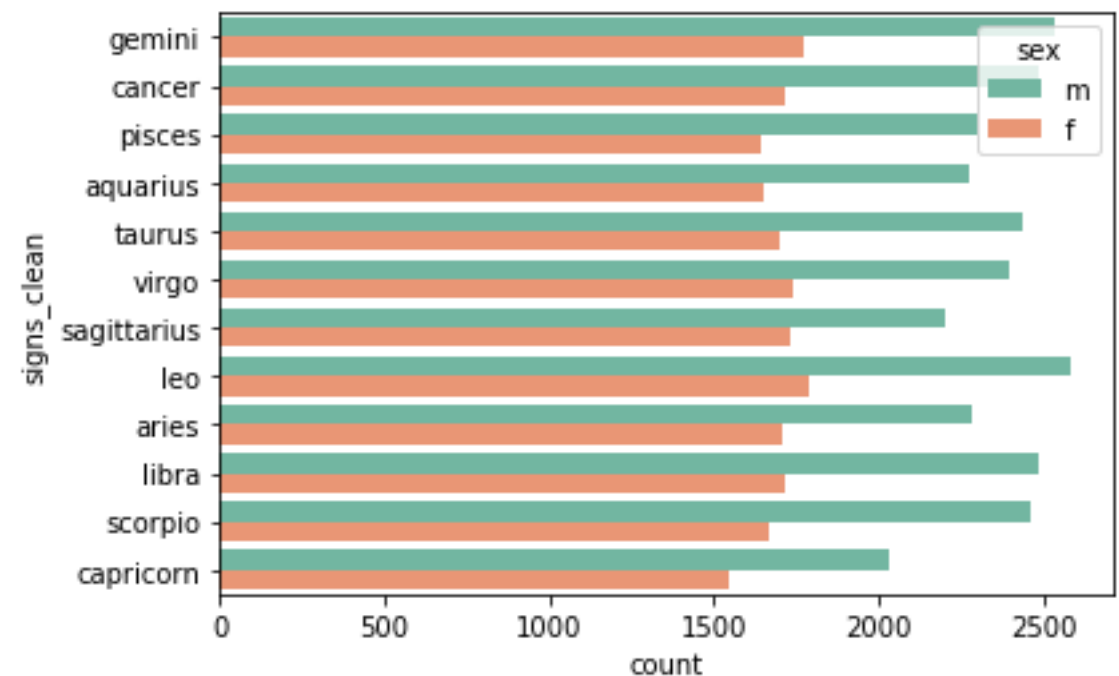
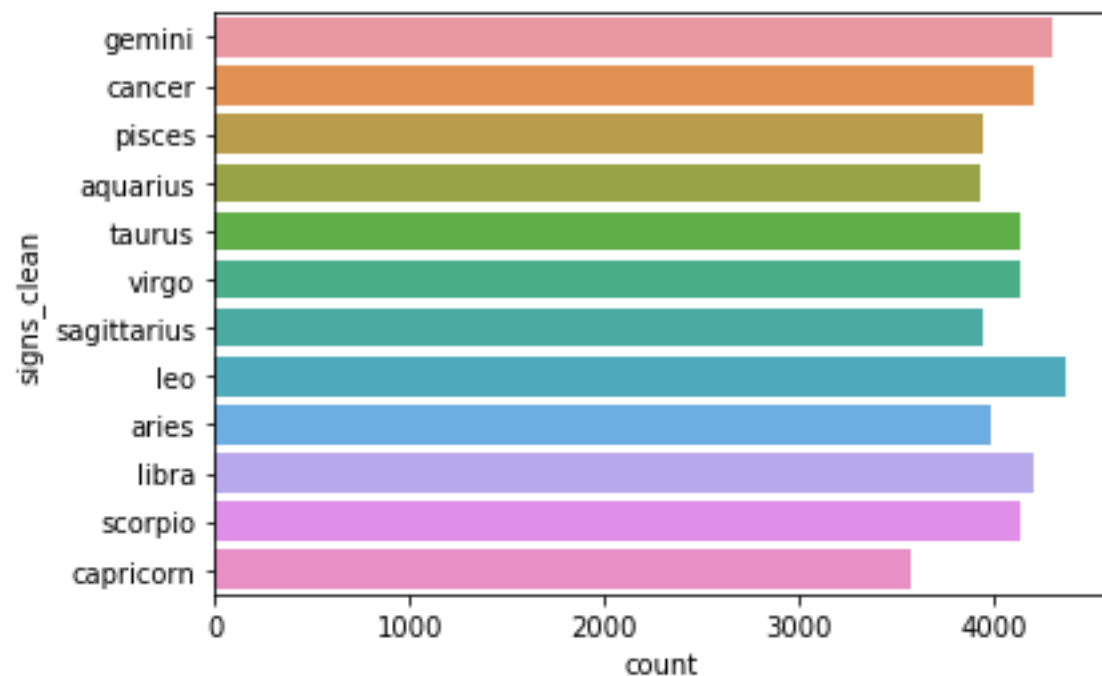Religions

Sexual orientation

Body type count

*Smoking users*

*Drinking users*

*Signs*

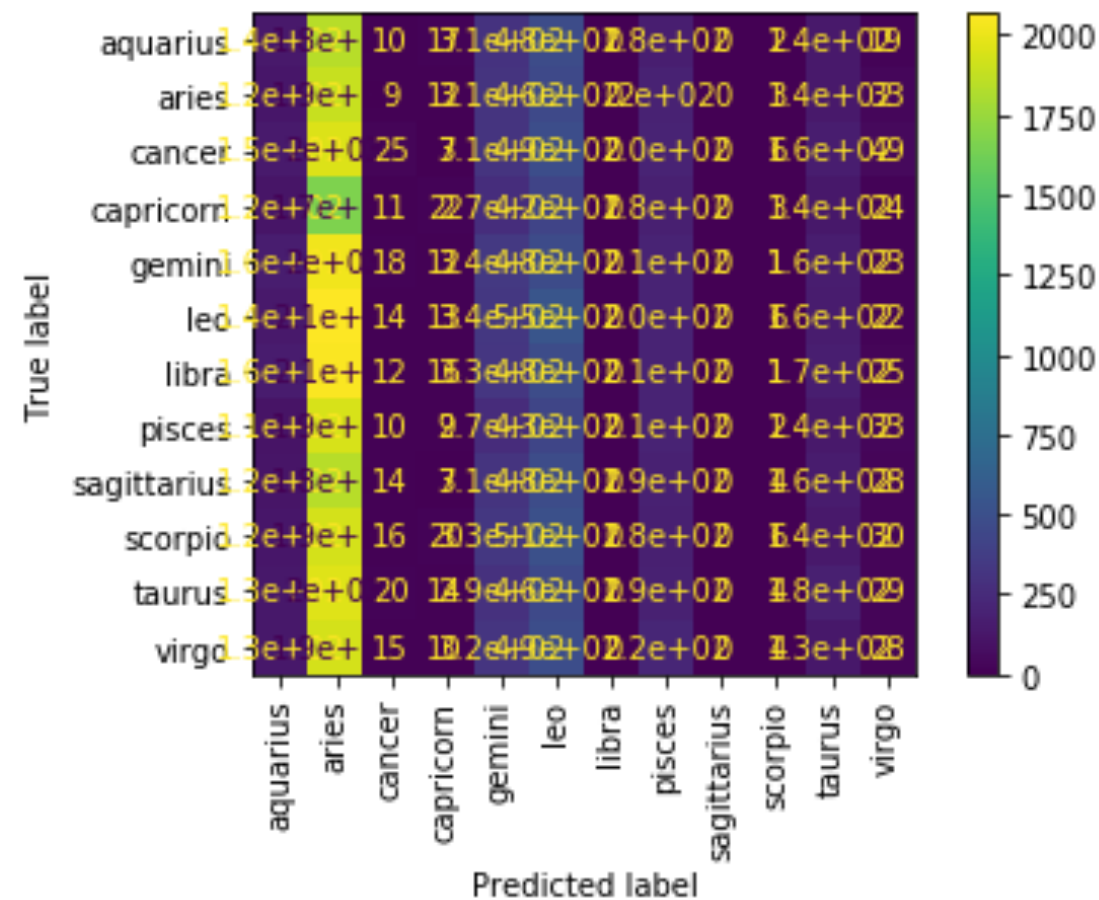As we can see, the signs column is the most proportionate in the data set, so next we will try to see if we can predict a user's sign based on other variables. First, I tried to see if we can predict a sign based on someone's smoking and drinking. I mapped the data and used the train_test_split function to train my models.

I used K Nearest Neighbors, Decision Tree Classifier, and Support Vector Machine.
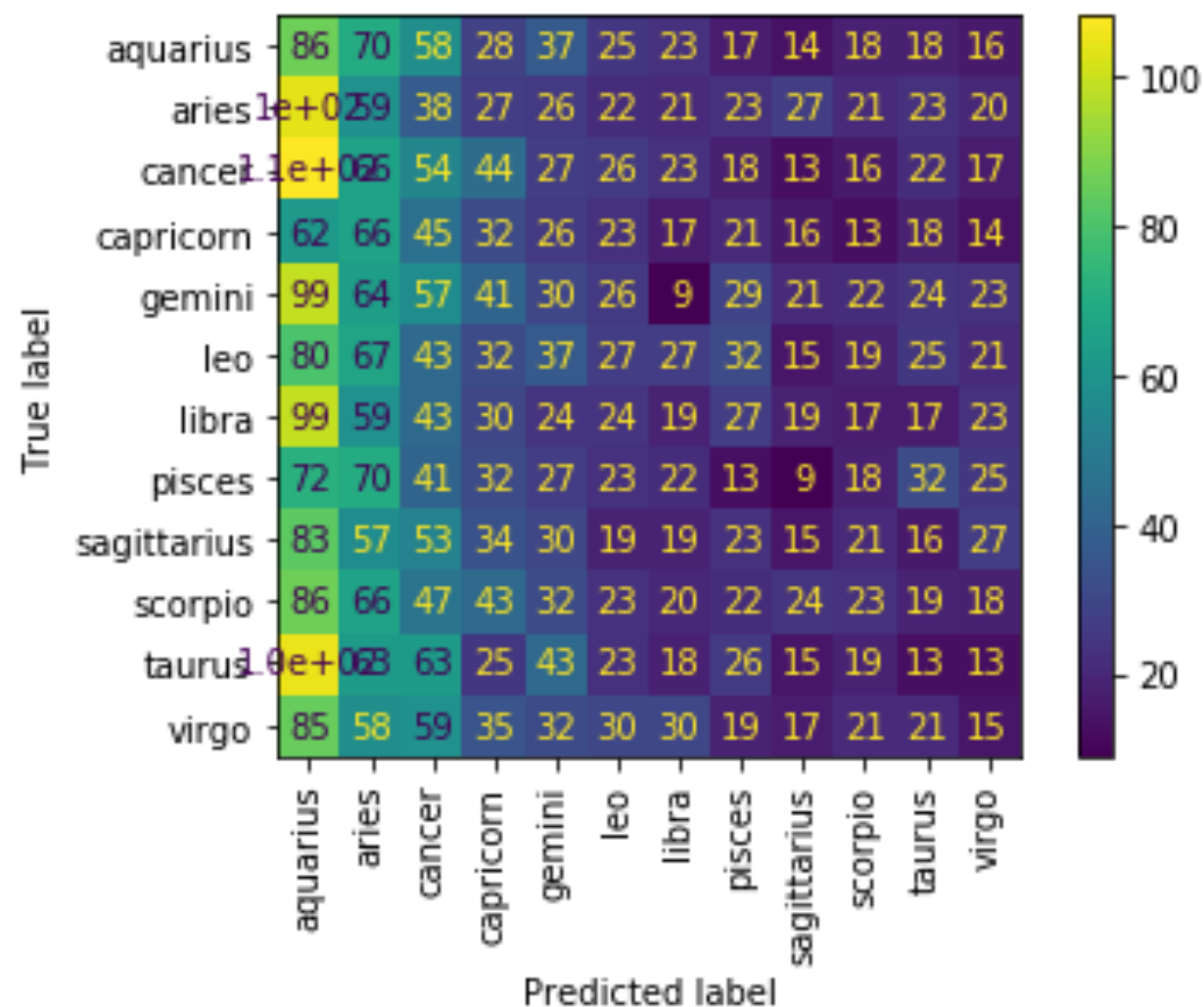
*KNN confusion matrix*

Mean accuracy for this model 0.08698097770505216, so obviously, what we are trying to achieve is not really possible. The model is only 8% accurate, which means that we should

rethink our strategy and maybe change the variables, as this is not enough data for the algorithm. Just for the sake of it, I went along and tried the other two methods as well. The mean accuracy for both of them was under 10% so I decided they are not relevant.

Next step was rethinking the strategy so I chose body type, diet, orientation, education, religion, sex and job columns as variables in trying to predict the sign.

KNN classifier had a mean accuracy of 0.32211180963969444, definitely an improvement from the first try.
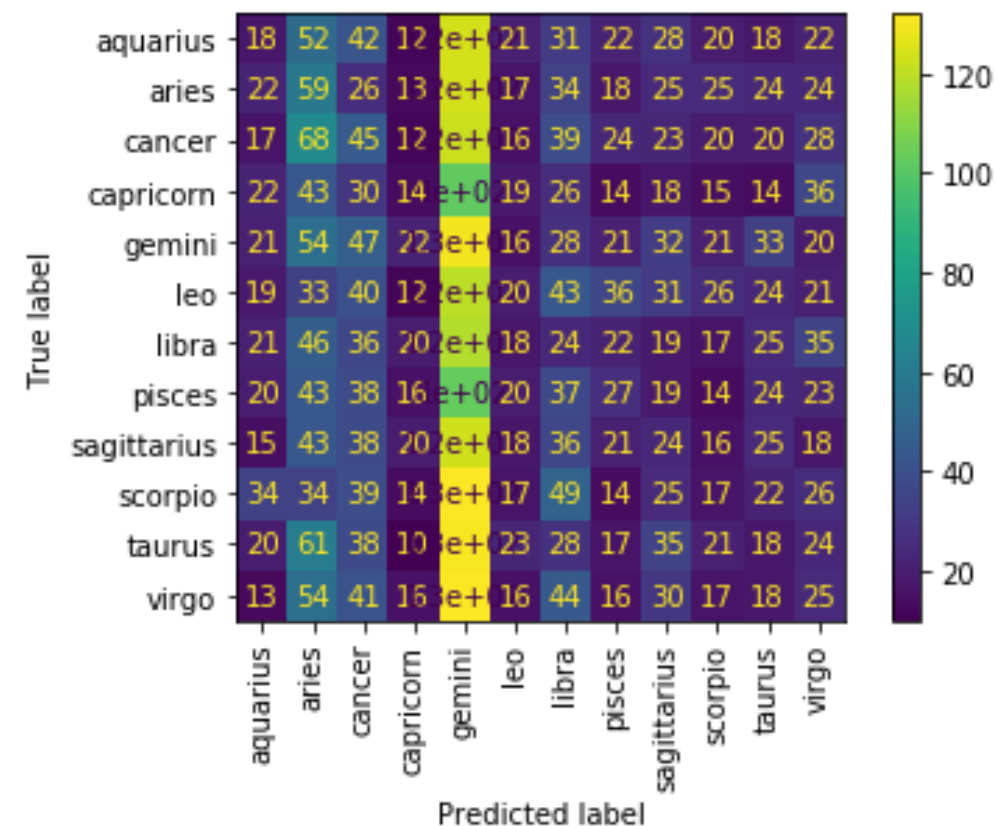


*KNN confusion matrix*

*Decision Tree Classifier confusion matrix*

The mean accuracy for the Decision Tree Classifier is 0.7323058203204218

Compared to KNN, this is a better fit to our data, and more reliable… but is it?
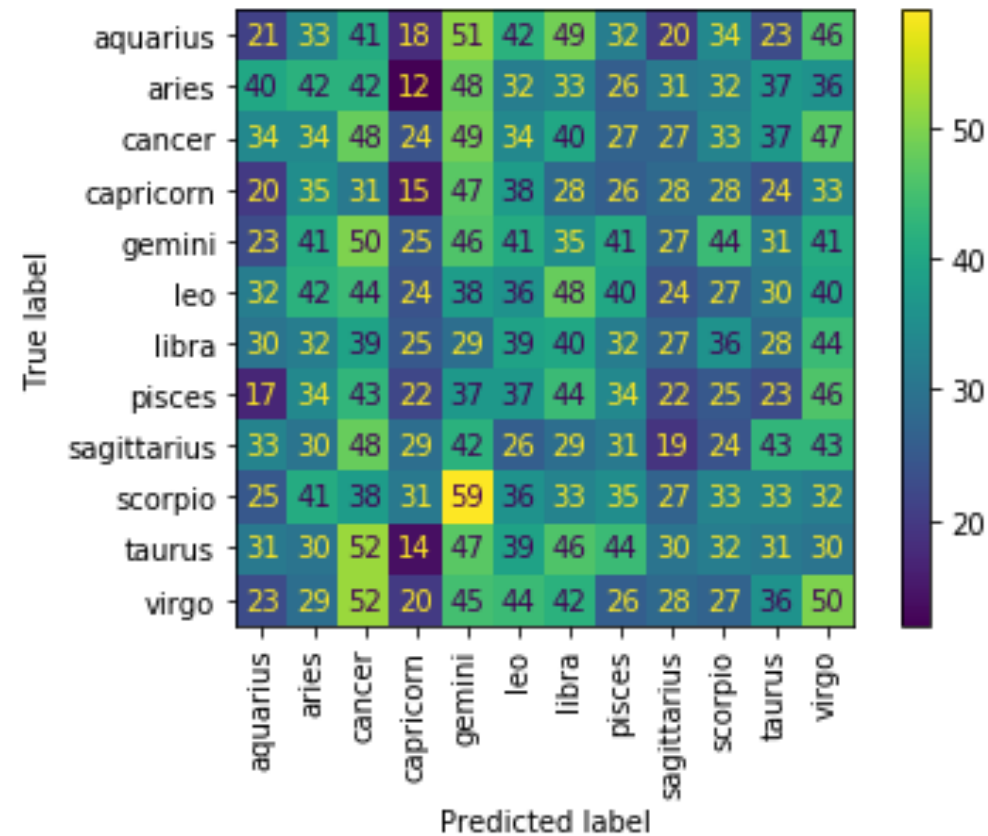
Our model's maximum depth is 64, so I tried to set the depth at 25.



*Decision Tree Classifier (25) confusion matrix*

After setting the maximum depth at 25, we get a mean accuracy of 0.38058541201919827

It is a bit above KNN, but surely if we set the depth lower, it will get lower than KNN.
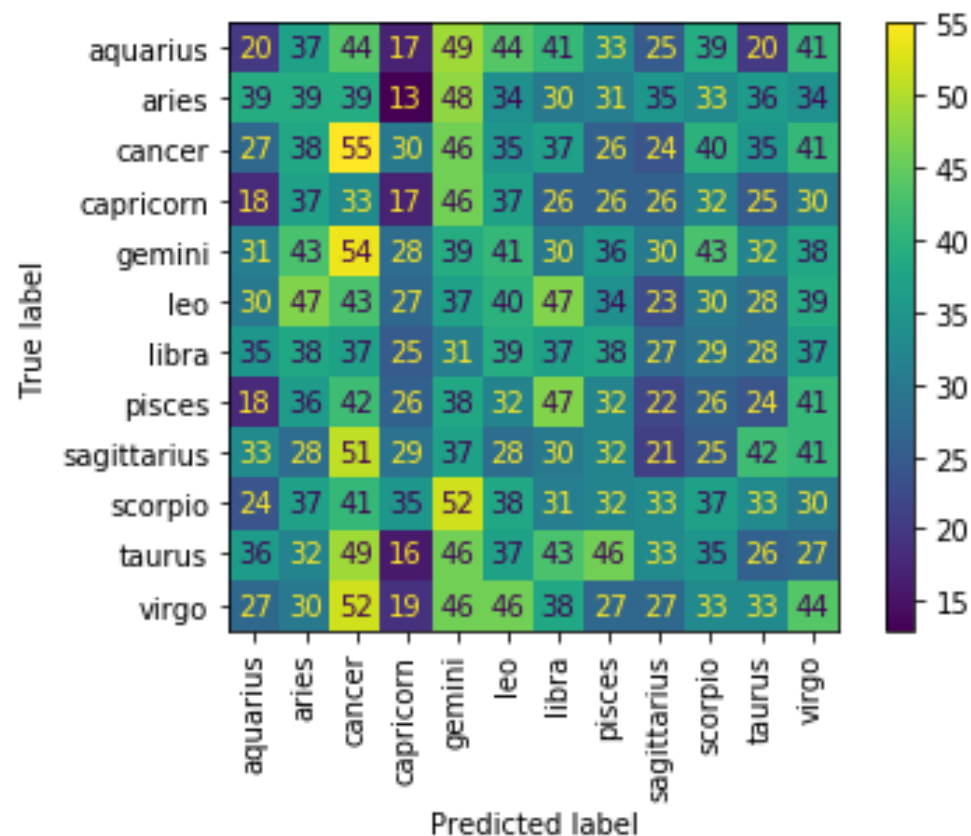
## Support Vector Machine confusion matrix

For SVM, we got a mean accuracy of 0.719461907659028

Our initial C parameter was set at 2. At C = 0.01 the data is under fit and mean accuracy is 0.09321976610559048.

I also changed C to equal 10, and our mean accuracy went up to 0.7323058203204218.

The second confusion matrix is for C = 10, so we can notice the differences.

In conclusion, it seems that Support Vector Machine is the best classifier to use for our data set, as the other classifiers did not perform up to our expectations.

Also, choosing your variables right makes a huge difference on machine learning algorithms performance.

This was a cool project to work on, and I feel like I have expanded my Python knowledge and also trained my skills.