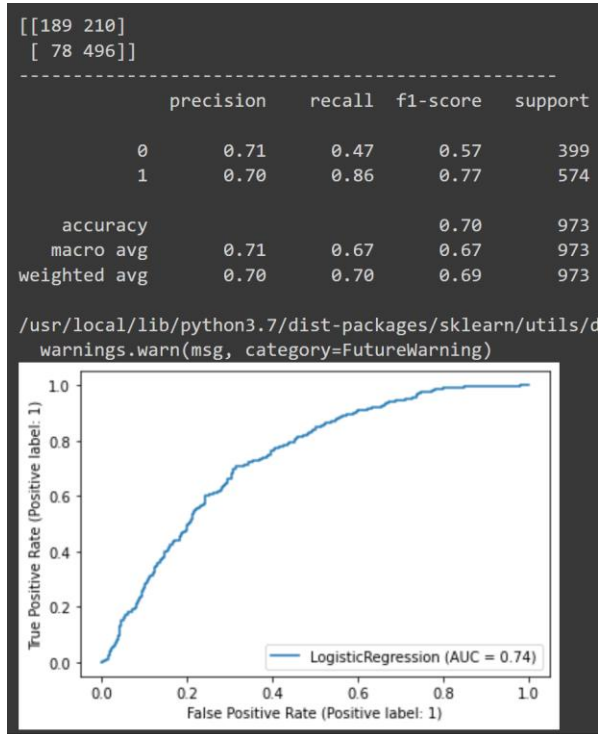


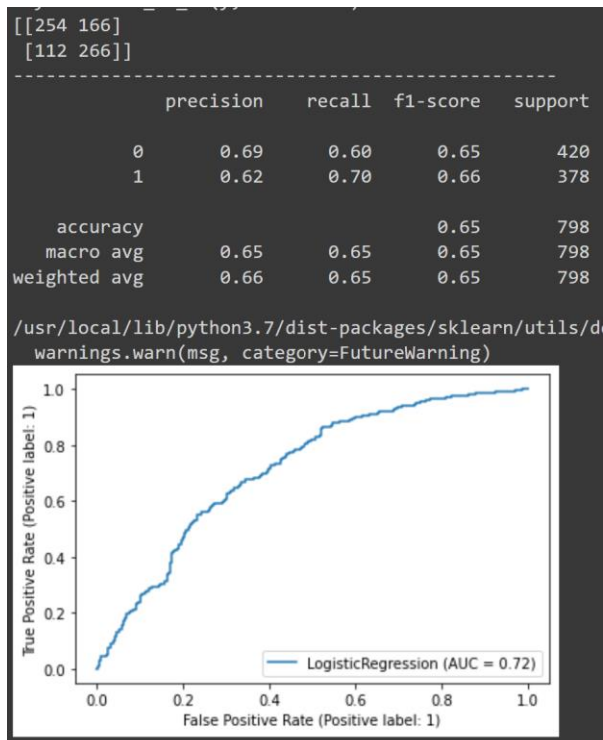
In this computer assignment, we observe a few interesting things.

First, we observe an unbalance in our dataset, where we have more positive than negative cases. For reliable results, we need to balance the sample. However, I decided to try to use the unbalanced data to train the logistic regression model. The metrics were quite surprising.



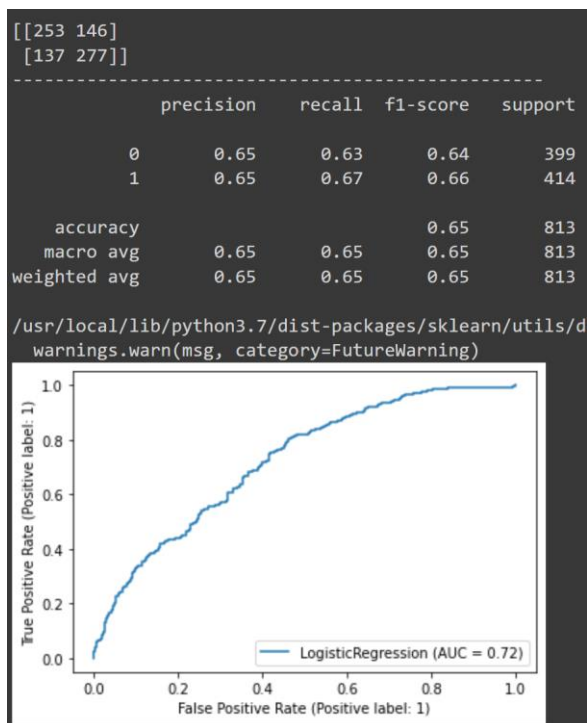
With a sample where there are more positive cases than negative ones, one would expect the machine will predict a lot more positives than negatives. However, the metrics for this model show otherwise. Somehow, the model made a lot of negative predictions than positive ones. This resulted in the model having high accuracy, positive recall, and precision. But if we look at recall for the negative predictions, we can see that the model makes a lot of false negatives predictions.

After Model 0, I did an under-sampling, and trained Model 1 using that dataset. The results were not impressive.

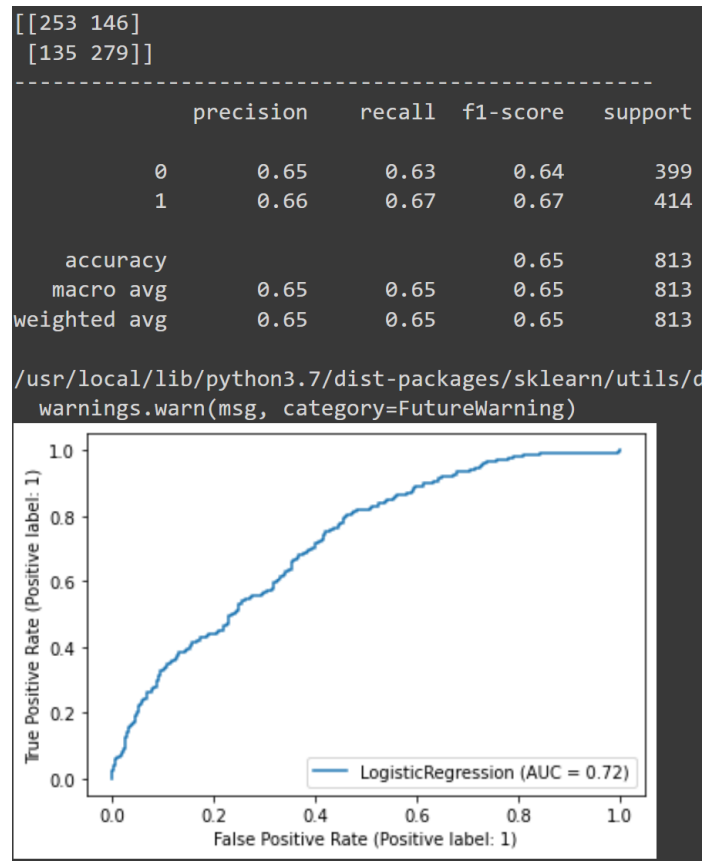


Although the model is making fewer false negative predictions, the other metrics have also declined. But this would be a more reliable model than Model 0.

Lastly, I used the SMOTE technique to over-sample the dataset, and then use the RFE technique to select my features. With this dataset, I trained Model 3. The model's metrics did not quite improve.



Even with a GridSearchCV looking for the best hyperparameters, the model's performance has only barely improved.



After hours of effort, I came to this conclusion: logistic regression might not be the best machine learning model for this specific dataset. So the next step would be trying some other models such as random forest, KNN, etc.