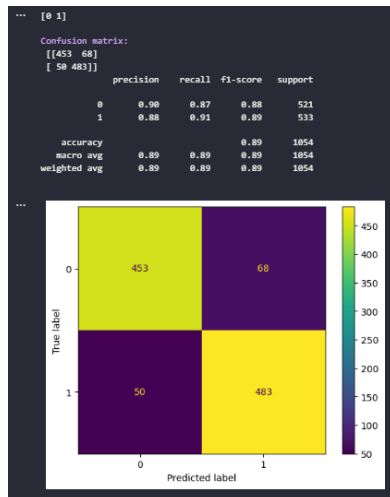


Assignment1 report

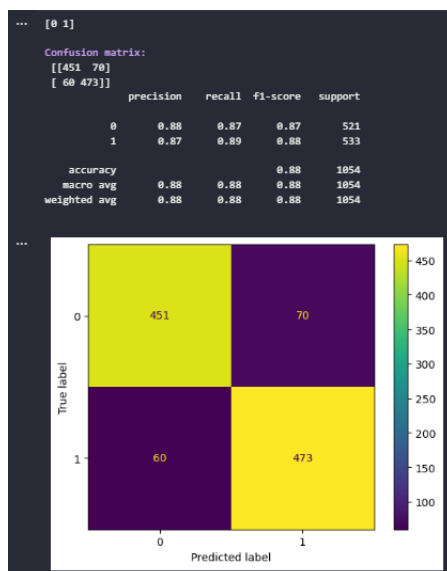
GridSearch on the Logistic Regression Model using TD-IDF Vectorizer:

```
%time
grid_result = grid_search(X_train, y_train)
print_info(grid_result)
✓ 5m 63s
```

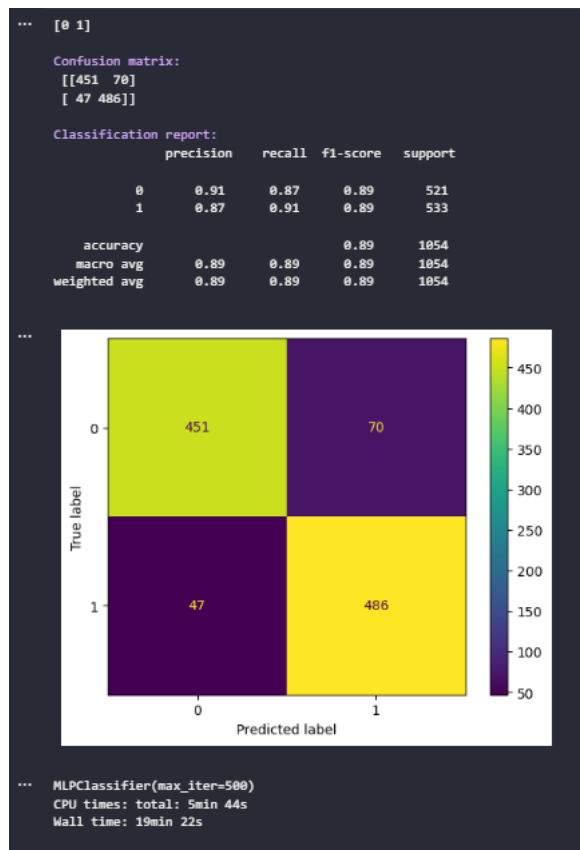


GridSearch on the Logistic Regression Model using Count Vectorizer

```
%time
grid_result_c = grid_search(X_train_c, y_train_c)
print(grid_result_c)
✓ 42m 23s
```



GridSearch on the MLP Classifier using TD-IDF Vectorizer



To conclude:

By reviewing all the three matrices altogether, we can see that by using the different models there are not a lot of differences that separates them. But by comparing the necessary computational time used for TD-IDF and CountVectorizer, and by looking at the fact that the precision (the performance of positive classifications that was actually correct) and recall (the proportion of actual positives that was identified correctly) for the models using TD-IDF performs slightly (almost insignificantly) better, it is safe to say that TD-IDF Vectorizer is the more efficient choice. For faster score interpretation we can just review the F1-score for all the confusion matrices.

Lastly, by taking into consideration that TD-IDF also takes into consideration the importance of words across documents using log, I actually assumed that it would perform better than it did against CountVectorizer.