# CSSM 502: ADVANCED DATA ANALYSIS WITH PYTHON

## DAVID CARLSON

## FINAL PROJECT

19.01.2022

EZGİ AYDIN

0078635

# TABLE OF CONTENTS:

# 1.INTRODUCTION

This present report is the outcome of the sentiment analysis performed over "Women's Clothing E-Commerce Reviews" dataset, as a final/term project of CSSM 502: Advanced Data Analysis with Python course. Mainly, the report first covers an exploratory data analysis and then continues with machine learning algorithm training.

I have chosen sentiment analysis as my project because a) While I was in university, I did work as a Marketing Insight Analyst for Ebrandvalue company, which provides advisory over "brand value" based on social media sentiment analysis. However, since I was part of marketing team and I had zero coding experience back then, I did not have chance to see/experience the kitchen of this analysis. Even though what I performed in this project is base level, I wanted to try and set this as a starting point for many to come. b) I have made some progress over the years, and I want to continue with that by learning and experiencing new things and topics. NLP is an area that I have great interest in and motivation to learn deeper about it.

Before continuing with the report, I just want to acknowledge that, this was the best class that I took this semester by far since we were treated as adults and given assignments were on spot to boost our skills in the area. Therefore, I want to take this opportunity to thank Professor Carlson for being progressive, understanding, broadminded, encouraging us to research and so many more. For working students, that means a lot.

This report consists of 10 parts, as presented in the table of contents. First, I will be starting with the exploratory data analysis to familiarize with the data set. Later, I will continue with data processing and machine learning algorithms.

## 2. SENTIMENT ANALYSIS

According to Towards Data Science website, sentiment analysis is "contextual mining of text which identifies and extracts subjective information in source material and helping a business to understand the social sentiment of their brand, product or service while monitoring online conversations.[1]" In other

---

[1] https://towardsdatascience.com/sentiment-analysis-concept-analysis-and-applications-6c94d6f58c17

words, sentiment analysis provides us with an overview of customer's opinion of subject matter. In this report's dataset, this is the positive and negative reviews of customers over women clothing.

## 3.GOAL DEFINITION AND DATASET

My goals in this project are:

1.To perform a comprehensive exploratory data analysis,

2. Observe whether customer reviews predict customer recommendation,

3.Train and test 4 different models and compare their performance on predicting positive/negative sentiments.

The dataset that I used is "Women Clothing E-Commerce Reviews[2]" gathered from Kaggle. It is a public dataset including 23000 customer reviews and ratings, with some additional demographics. Because this is real commercial data, it has been anonymized, and references to the company in the review text and body have been replaced with "retailer". The data set includes following columns:

```
Index(['Unnamed: 0', 'Clothing ID', 'Age', 'Title', 'Review Text', 'Rating',
       'Recommended IND', 'Positive Feedback Count', 'Division Name',
       'Department Name', 'Class Name'],
```

- Clothing ID: Integer Categorical variable that refers to the specific piece being reviewed.
- Age: Positive Integer variable of the reviewer's age.
- Title: String variable for the title of the review.
- Review Text: String variable for the review body.
- Rating: Positive Ordinal Integer variable for the product score granted by the customer from 1 Worst to 5 Best.
- Recommended IND: Binary variable stating where the customer recommends the product where 1 is recommended, 0 is not recommended.
- Positive Feedback Count: Positive Integer documenting the number of other customers who found this review positive.
- 8) Division Name: Categorical name of the product high level division.
- 9) Department Name: Categorical name of the product department name.
- 10) Class Name: Categorical name of the product class name.

---

[2] https://www.kaggle.com/nicapotato/womens-ecommerce-clothing-reviews

## 4.EXPLORATORY DATA ANALYSIS

```
Shape:(23486, 11)
------------------------------------------------
Info:

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 23486 entries, 0 to 23485
Data columns (total 11 columns):
 #   Column                  Non-Null Count  Dtype
---  ------                  --------------  -----
 0   Unnamed: 0              23486 non-null  int64
 1   Clothing ID             23486 non-null  int64
 2   Age                     23486 non-null  int64
 3   Title                   19676 non-null  object
 4   Review Text             22641 non-null  object
 5   Rating                  23486 non-null  int64
 6   Recommended IND         23486 non-null  int64
 7   Positive Feedback Count 23486 non-null  int64
 8   Division Name           23472 non-null  object
 9   Department Name         23472 non-null  object
 10  Class Name              23472 non-null  object
dtypes: int64(6), object(5)
memory usage: 1.5+ MB
None
```

I first performed an exploratory data analysis, which includes overviewing the dataset closely by each column and understanding it in several different aspects. Detailed work can ben found on py. File. Some analysis results are as follow:

```
Missing Values:
                Missing_Number  Missing_Percent
Title                    3810             0.16
Review Text               845             0.04
Division Name              14             0.00
Department Name            14             0.00
Class Name                 14             0.00
```

```
Number of Uniques:
Unnamed: 0              23486
Clothing ID             1206
Age                       77
Title                  13993
Review Text            22634
Rating                     5
Recommended IND            2
Positive Feedback Count   82
Division Name              3
Department Name            6
Class Name                20
dtype: int64
```
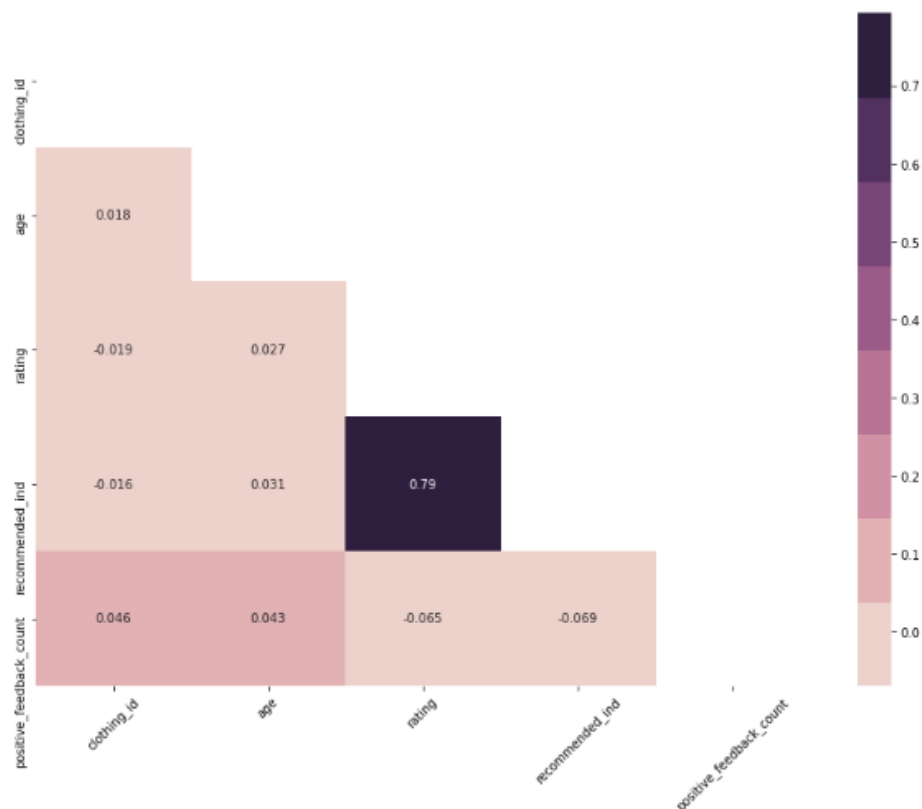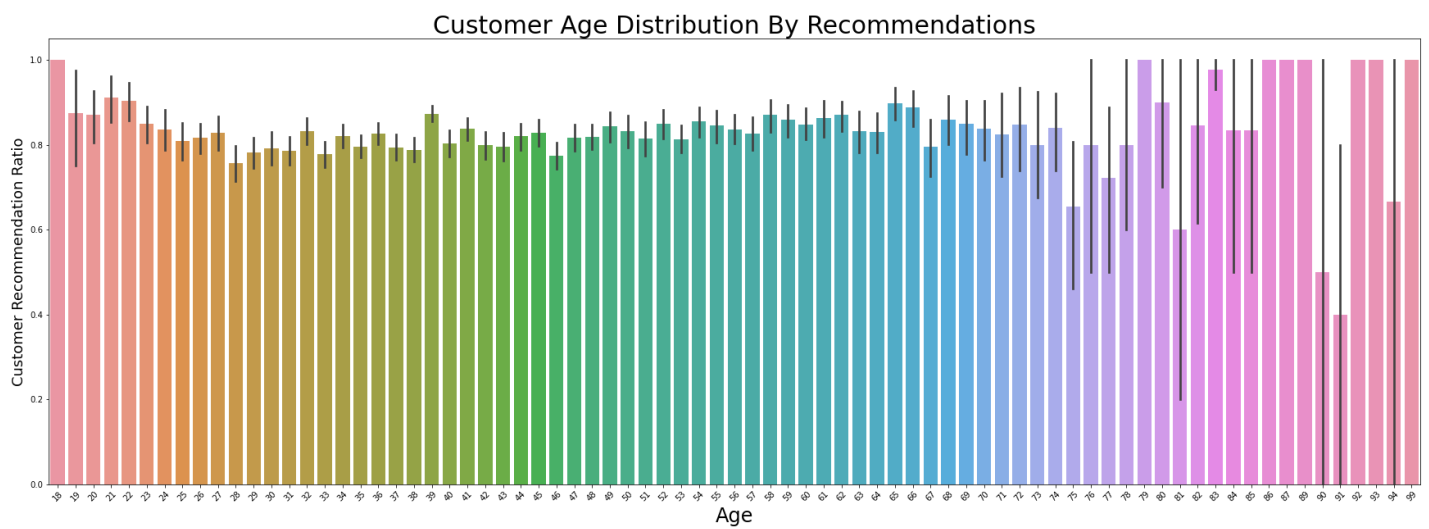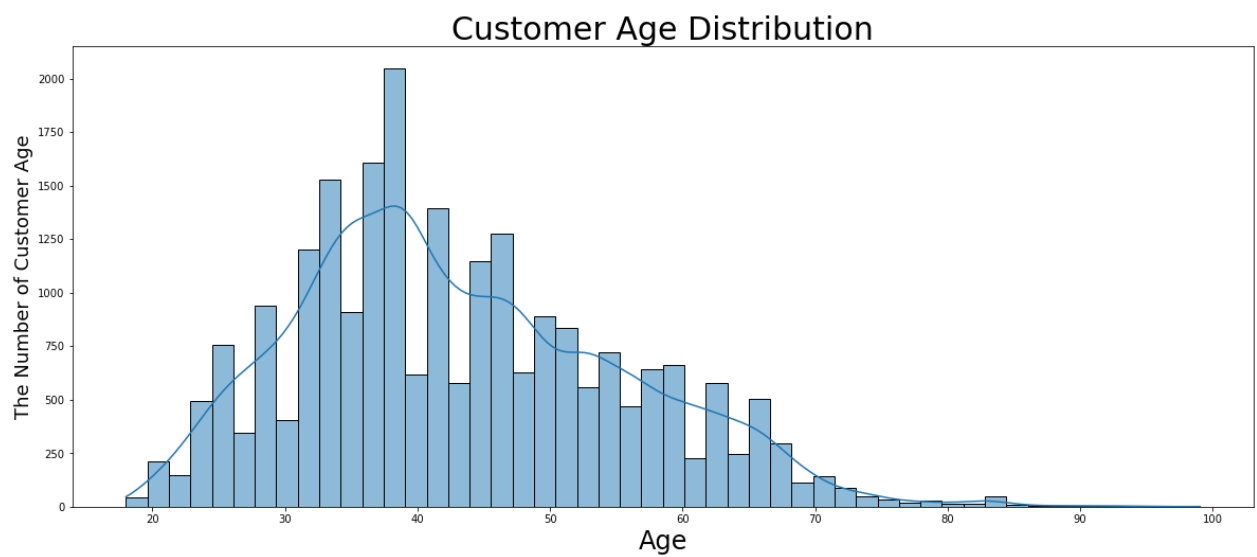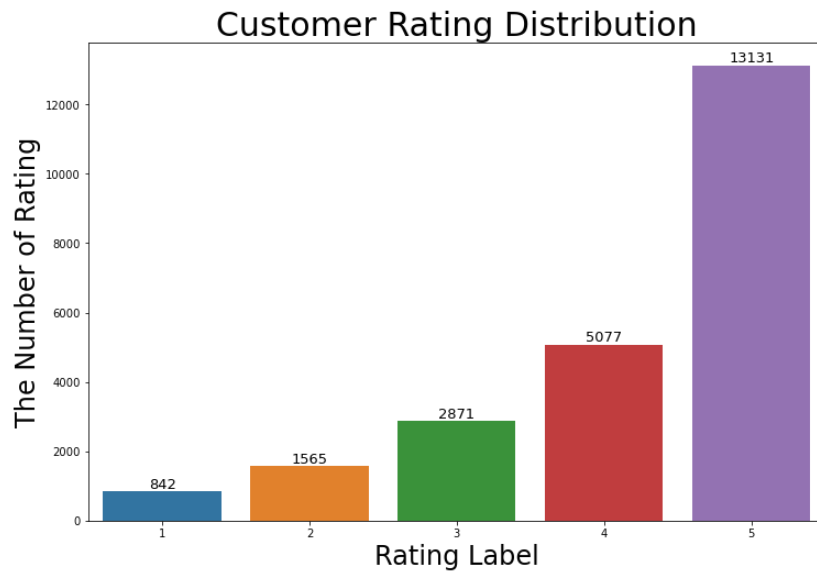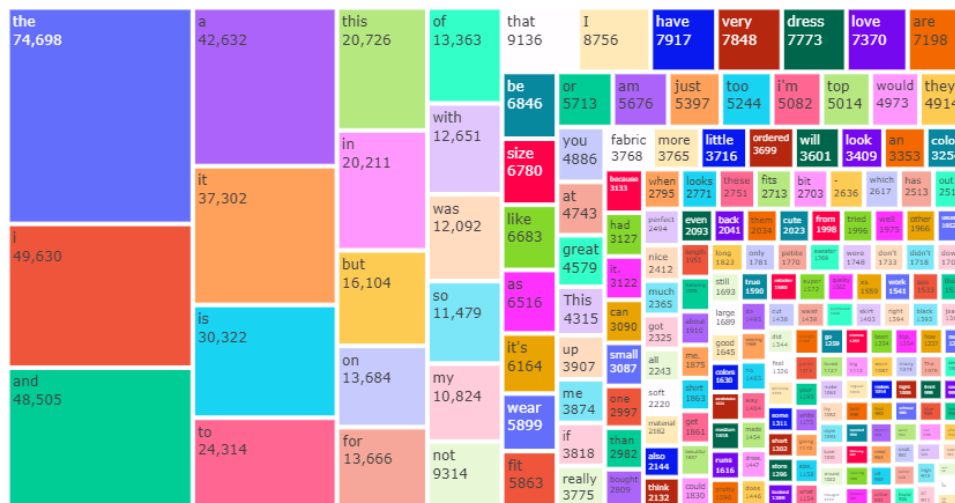
While observing the dataset, I looked for correlation between columns and the analysis resulted 0.79 correlation between Rating and Recommended IND columns. Additional to correlation heatmap, I will be sharing other exploratory data analysis results here with visualized versions.

E. Aydın
19.01.2022
CSSM 502

## Customer Rating Distribution



## Customer Age Distribution



## Customer Age Distribution By Recommendations

Top Frequent 200 Words in the Dataset (Before Cleaning)



## 5.FEATURE SELECTION & DATA CLEANING

Since I am only using 2 features from this dataset, which are "Review Text" and "Recommended IND", I have dropped other unnecessary columns. After that I have performed some missing value analysis and since the number was relatively small, I dropped those rows too.

## 6.TEXT MINING

Text is an unstructured form of data, therefore various types of noise are present in it, which causes inability to analyze without any pre-processing. The process of making this text data ready for analysis is called text preprocessing. In this project, to preprocess the data I have first removed all punctuations, numbers, stop words, outliers, tokenized text and applied lemmatizer (Lexicon Normalization). After preprocessing, I created word clouds for all, positive and negative words.
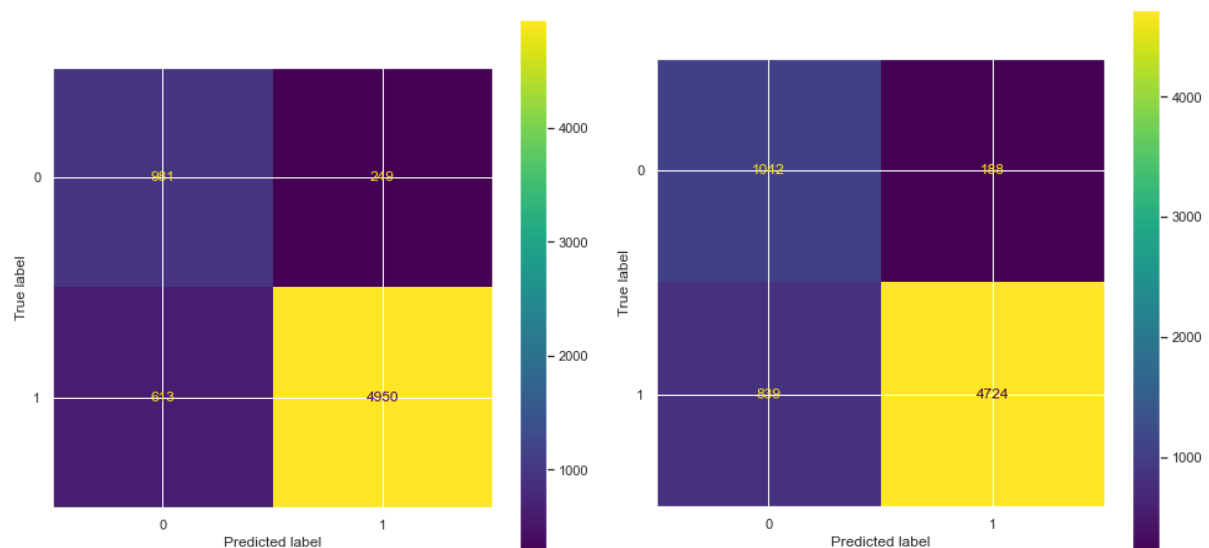
## 7.WORDCLOUD

## 8.TRAIN/TEST SPLIT & VECTORIZATION

Before performing any modeling, I needed to apply vectorization to convert text documents into numeric feature vectors and train test split. From the tokenization step, review text column is in token form, so I used scikitlearn count vectorizer to convert that text into a matrix of token counts. To be able to evaluate from multiple approaches, I performed 2 different vectorization methods:1. Count Vectorization, 2. TF-IDF Vectorization. The reason for this is that I read that count vectorization is based on frequency representation and a simpler form of vectorization, on the other hand TF-IDF takes weight into account. So, if there is any, we will be able to observe their differences as well as classification algorithms. After this step, I will be building models using 4 classification algorithms.
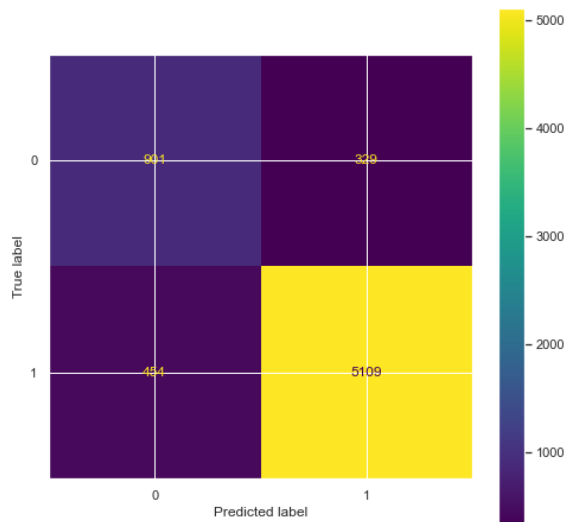
## 9.MACHINE LEARNING MODELS

After train and test split and vectorization I've created models with 4 classification algorithms, which are Logistic Regression, Naive Bayes, Support Vector Machine and Random Forest. Following are the heatmaps referencing models.
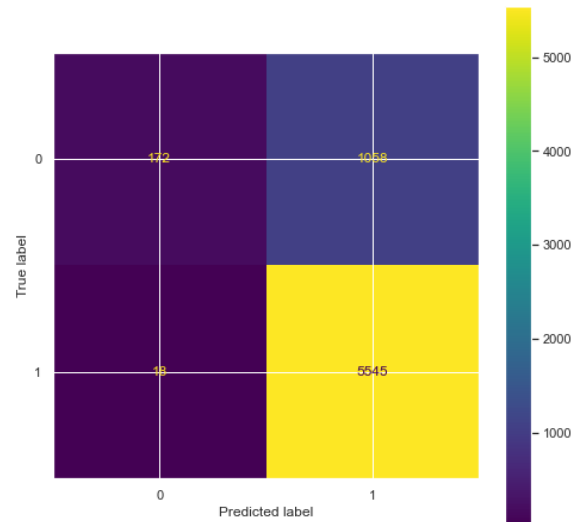
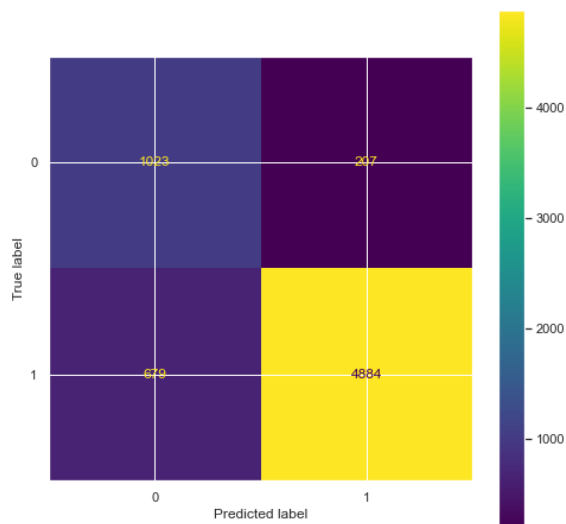1)Logistic Regression with Count Vectorizer        2) Logistic Regression With TF-IDF Vectorizer
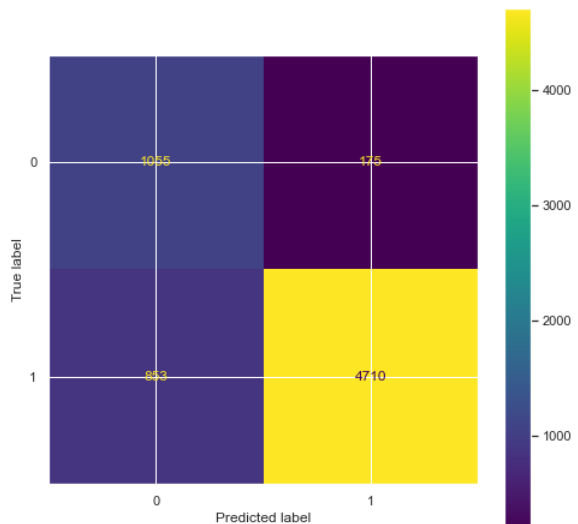
### 3) Naive Bayes with Count Vectorizer
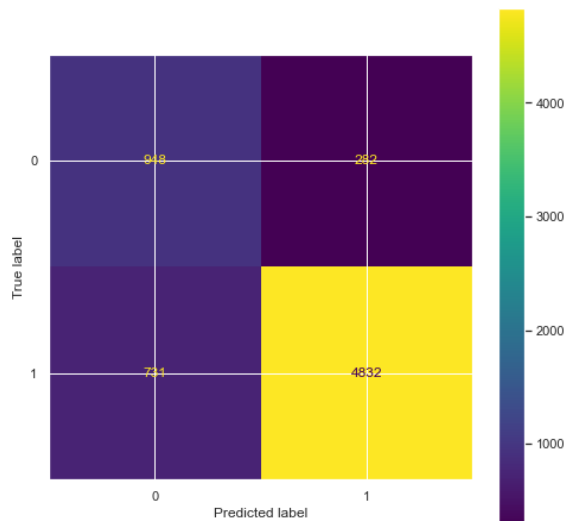


### 4) Naive Bayes with TF-IDF Vectorizer
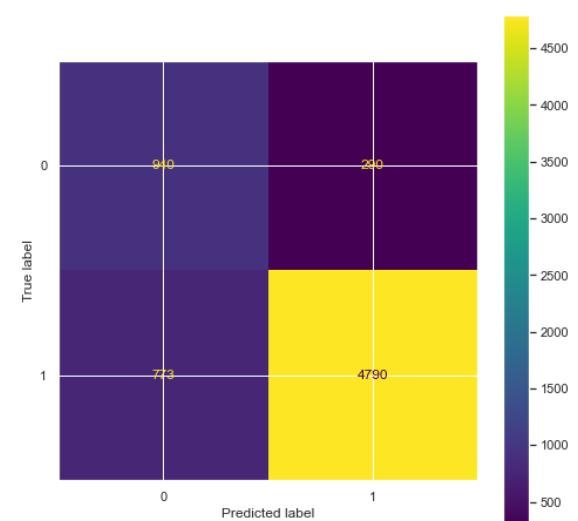


### 5) SVM With Count Vectorizer



### 6) SVM With TF-IDF Vectorizer



### 7) Random Forest with Count Vectorizer



### 8) Random Forest With TF-IDF Vectorizer

## 10.COMPARISON OF MODELS AND EVALUATION

As conclusion, since the result of each model is relatively similar to each other, I think the best model could be chosen based on which indicator we want to value most, from accuracy and precision. Additionally, I did not observe a significant affect by changing the vectorization method, since scores are also similar in that case.

E. Aydın
19.01.2022
CSSM 502

RESOURCES:

https://towardsdatascience.com/sentiment-analysis-concept-analysis-and-applications-6c94d6f58c17

https://www.analyticsvidhya.com/blog/2021/06/part-5-step-by-step-guide-to-master-nlp-text-vectorization-approaches/

https://www.guru99.com/nltk-tutorial.html

https://www.kaggle.com/edchen/text-vectorization

https://medium.com/@paritosh_30025/natural-language-processing-text-data-vectorization-af2520529cf7

https://www.dataschool.io/comparing-supervised-learning-algorithms/

https://towardsdatascience.com/how-to-turn-text-into-features-478b57632e99

https://machinelearningmastery.com/prepare-text-data-machine-learning-scikit-learn/