

Model Performance

Model Name: MultinomialNaiveBayesOnScore Test Date: 23/03/2022 15:56:00 Creator: Tobias Rothlin

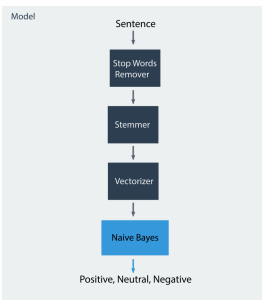


Overview

ML Principle:
Multinomial Naive Bayes

- References:**
- [MultinomialNB Explained](#)
 - [Stanford NLP Course](#)
 - [Stanford NLP Lecture](#)
 - [English Stopwords](#)

Algorithm Description:
The learning algorithm used in this classification is the Multinomial Naïve Bayes. This approach was chosen as it is easy to implement and is computational very efficient. The first step in the classification pipeline is removing all stop words for example 'i', 'me', 'my', 'myself', etc. A list of English stop word is provided by the nltk module. The stop words remover just removes every word that is in the list of stop words. Next the sentence is passed through the stemmer. Stemmers remove morphological affixes from words, leaving only the word stem. This is done with the PorterStemmer class from the nltk module. The final preprocessing step is to vectorize the sentence. This results in a bag of words representation of the sentence. First all the words must be tokenized and then counted. The result will be a numerical feature vector. To generate this vector the CountVectorizer class from sklearn is used. This class implements both tokenization and occurrence counting in a single class. With the sentence now represented in a vector the Naïve Bayes classifier can work with this vector. For the implementation of the Naïve Bayes classifier the MultinomialNB class (sklearn) is used.



Classification Pipeline

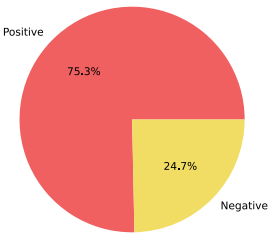
Metrics

Data: ClassifiedDataSetV1.2 with 10 folds cross validation

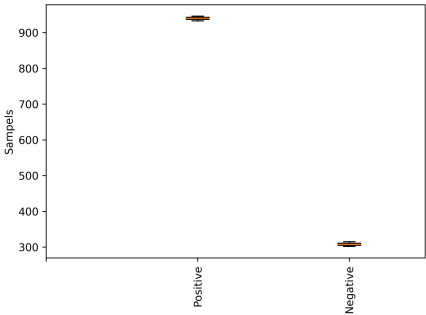
Split seed: 4.83819

Training Dataset (average)

Classes	Number of samples
Positive	940
Negative	307



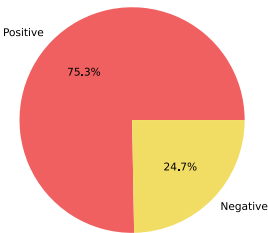
Average distribution of the samples



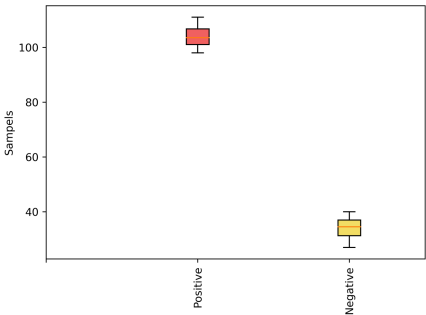
Distribution of the samples contained in each test split

Test Dataset (average)

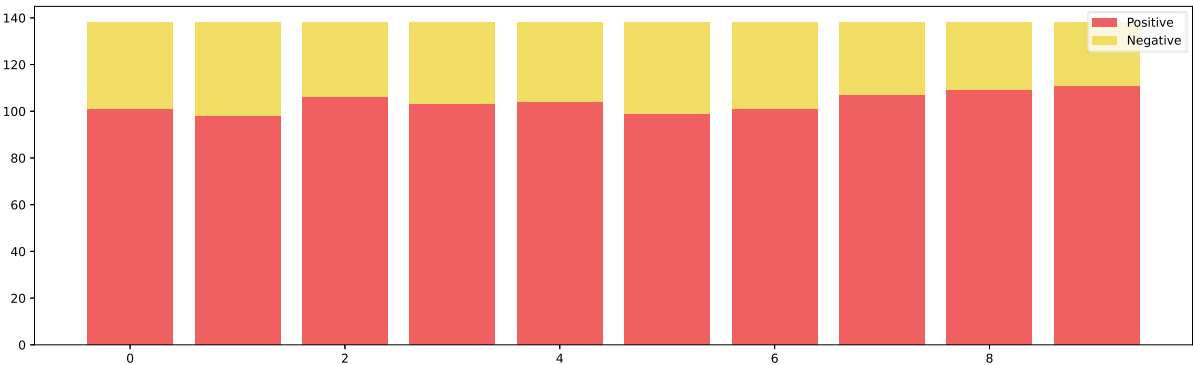
Classes	Number of samples
Positive	103
Negative	34



Average distribution of the samples



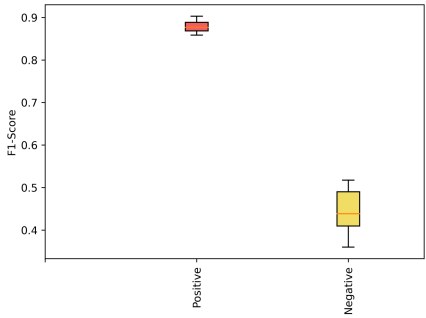
Distribution of the samples contained in each test split



Detailed training split composition

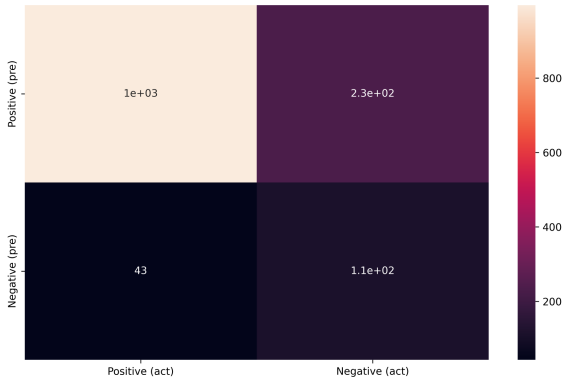
Classification Performance

Classes	Precision	Recall	F1 Score
Positive	81.24%	95.86%	87.95%
Negative	72.08%	32.55%	44.85%
Accuracy			80.22%
Macro Average	76.66%	64.21%	66.40%
Weighted Average	78.98%	80.22%	77.30%

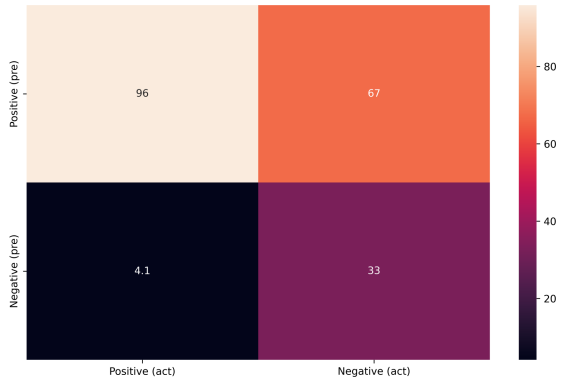


Distribution of the F1-Score

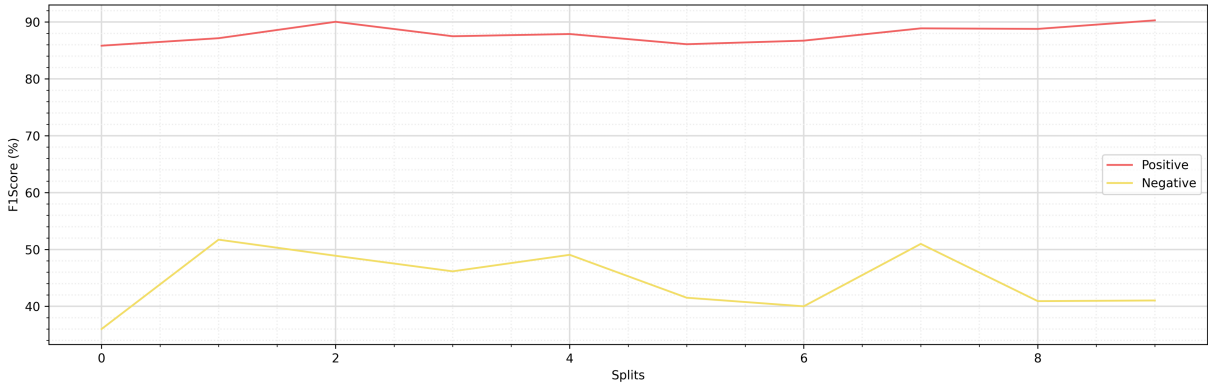
ConfusionMatrix:



Normalised ConfusionMatrix:



F1 Score by split:



F1-Score per split