

Machine Learning Engineer Nanodegree

Capstone Report

Patrick O'Sullivan

June 20th, 2019

Business Public Sentiment

Baseline Model

Naïve Bayes

I used the Naïve Bayes model as my baseline model as it is a quick and easy way to predict classes. They are based on a statistical classification technique known as Bayes Theorem. They are naïve because they assume that the variables are independent from each other.

The calculate the posterior probability of a certain event A to occur given some probabilities of prior events.

$$P(A|R) = \frac{P(R|A)P(A)}{P(R)}$$

P(A): Prior probability of a class.

P(R|A): Likelihood the probability of predictor given class.

P(R): Prior probability pf predictor.

P(A|R): Posterior probability of class A given the predictor R.

Step to preform

To use the Naïve Bayes algorithm to solve classification problems like deciding if text is positive or negative.

- Convert the dataset into a frequency table.
- Find the probabilities of the events to occur.
- Computes the posterior probability of each class.
- The class with the highest posterior probability is the prediction.

Strengths & Weaknesses

Strengths:

- Easy and quick to implement.
- If the conditional independence holds then it will converge quickly.
- Need less training data.
- Scalable.

[illegible]

With the large amount of training data and the various characters people use in tweets we have a lot of strange words in the column. It also means we have a lot of unique words.

Count Vectorizer

I use the `sklearn.feature_extraction.text.CountVectorizer` to

- Separate the string into individual words and give each word (column) an integer ID.
- Count the occurrences of each word (column) in each tweet (row).
- Convert all words to lower case.
- Ignore all punctuation.
- Ignore all stop words.

Split Dataset

I split the dataset into four buckets

```
X_train: Training data for the text (tweet) column.  
y_train: Training data for the label column.  
X_test:  Testing data for the text (tweet) column.  
y_test:  Testing data for the label column.
```

Number of rows in the total set: 17952

Number of rows in the training set: 13464

Number of rows in the test set: 4488

After splitting the dataset, we need to convert the data into the matrix format just like we did in the sample above using `CountVectorizer`. For the training data we need to fit first, so the model can learn a vocabulary and then transform to the matrix view. For the testing dataset we only need to transform to the matrix view.

Model

For this project I used the multinomial Naïve Bayes implementation because it is suitable for classifications with discrete features like word counts for text classification.

```
def fitNaiveBayes(self):  
    self.naive_bayes = MultinomialNB()  
    self.naive_bayes.fit(self.training_data, self.y_train)
```

Now that I've trained the model with the training dataset I can use the testing dataset to make predictions.

```
def predict(self):  
    self.predictions = self.naive_bayes.predict(self.testing_data)
```

Evaluation

And final now that we have predictions I can evaluate the model to check the accuracy of my baseline model. There are several ways to evaluate a model.

Accuracy: is the fraction of predictions our model got right (e.g. what proportion of twits we predicted/classified as positive or negative were actually positive or negative).

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision: tells us what proportion of positive identifications was correct (e.g. what proportion of actually positive twits we predicted as positive twits).

$$Precision = \frac{TP}{TP + FP}$$

Recall: tells us what proportion of positives was identified correctly (e.g. what proportion of all positive twits did we predict as positive)

$$Recall = \frac{TP}{TP + FN}$$

F1 Score: is a weighted average of precision and recall ranging from 0 to 1.

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall}$$

TP: True Positive
TN: True Negative
FP: False Positive
FN: False Negative

Scores

Accuracy score: 0.7667112299465241
Precision score: 0.7834830979888746
Recall score: 0.7719224283305227
F1 score: 0.7776598003822467

Conclusion

References

1. Sentiment140: <http://help.sentiment140.com/for-students/>
2. Twitter Search API:
<https://developer.twitter.com/en/docs/tweets/search/api-reference/get-search-tweets>
3. Kaggle: <https://www.kaggle.com>