

DATA MINING

ASSIGNMENT 3

Ali Gholami

Department of Computer Engineering & Information Technology
Amirkabir University of Technology

<https://aligholamee.github.io>
aligholami7596@gmail.com

Abstract

Naive Bayes is a family of algorithms based on applying Bayes theorem with a strong (naive) assumption (bias), that every feature is independent of the others, in order to predict the category of a given sample. They are probabilistic classifiers, therefore will calculate the probability of each category using Bayes theorem, and the category with the highest probability will be output. Naive Bayes classifiers have been successfully applied to many domains, particularly Natural Language Processing (NLP).

Keywords. *Natural Language Processing, Text Processing, Text Classification, Spam Filter, Naive Bayes Classifier, Weka.*

1 Text Preprocessing

1.1 Separating Train and Test Data

In the first step, we need to properly separate the train samples and their labels from each other. To achieve this, we have implemented the function *extract labels* to do that for us. Columns are based on the current dataset which is called *spam collection dataset*.

```
def extract_labels(data):  
    return (dataset['v2'], data['v1'])
```

1.2 Extracting Features of Texts

In this section, we'll be using *Bag of Words* model from *Scikit-learn* to turn the text content into numerical feature vectors. This model does the following steps to turn texts into feature maps.

1. Assign a fixed ID to every word occurring in any document (for instance by building a dictionary from words to integer indices).
2. Each unique word in our dictionary will correspond to a feature (descriptive feature).

Scikit-learn has a high level component which will create feature vectors for us *CountVectorizer*.

```
count_vect = CountVectorizer()
X_train = count_vect.fit_transform(train_data)
```

It worths mentioning that the parameters and return value of the *fit_transform* function is as follows:

Parameters

`raw_documents` : iterable

An iterable which yields either `str`, `unicode` or `file` objects.

Returns

`X` : array, [n_samples, n_features]

Document-term matrix.

1.3 Issue With Occurrences

Occurrence count is a good start but there is an issue: longer documents will have higher average count values than shorter documents, even though they might talk about the same topics.

To avoid these potential discrepancies it suffices to **divide** the number of occurrences of each word in a document by the total **number of words in the document**. We'll delve into the formal representation of the *TF-IDF*. It is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.

1.3.1 Term Frequency

We denote the raw count of a word in a document as $f_{t,d}$. It show the number of times that term t has occurred in the document d . There are also other approaches to denote the number of occurrences of a term in a document and they are trying to illustrate the weighted frequency of each term in a document. Some of them are:

- Boolean Frequencies
- Logarithmically Scaled Frequency
- Augmented Frequencies

1.3.2 Inverse Document Frequency

The inverse document frequency is a measure of how much information the word provides, that is, whether the term is common or rare across all documents. It is the logarithmically scaled inverse

fraction of the documents that contain the word, obtained by dividing the total number of documents by the number of documents containing the term, and then taking the logarithm of that quotient.

$$idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

Where N is the number of documents in the database D. The denominator illustrates the number of documents in which they include the term t .

1.3.3 TF-IDF

A high weight in tf-idf is reached by a high term frequency (in the given document) and a low document frequency of the term in the whole collection of documents.

$$tfidf(t, d, D) = tf(t, d) * idf(t, D)$$

So far we have defined metrics that can be conducted to improve the features extracted from the emails in the first step. We'll change the transformer to use the *tf-idf* metric to extract features. Using this approach will result in a more precise text classification. This will **normalize the frequency of redundant words**.

```
count_vect = CountVectorizer()
X_train = count_vect.fit_transform(train_data)

tf_idf = TfidfTransformer()
X_train = tf_idf.fit_transform(X_train)
```

2 Text Classification

Now that we have our features, we can train a classifier to try to predict the category of a post as *Spam* or *Ham* (Not Spam). There are various algorithms which can be used for text classification. We will start with the most simplest one called Naive Bayes. For this task, we only need the two class Naive Bayes classifier. We have also included the testing section of the code. The result will be available in the *predict* variable.

```
clf = MultinomialNB().fit(train_data, train_labels)
test_data = count_vect.transform(test_data)
test_data = tf_idf.transform(test_data)
print("Accuracy: ", accuracy_score(test_labels, predicted))
```

In order to use *Recall* and *Precision* metrics of Scikit-learn, we have to convert the labels into binary format. We'll use *Label Binarizer* to achieve this.

```
lb = preprocessing.LabelBinarizer()
predicted_binarized = lb.fit_transform(predicted)
test_labels_binarized = lb.fit_transform(test_labels)
```

```
print("Recall: ", recall_score(test_labels_binarized, predicted_binarized))
print("Precision: ", precision_score(test_labels_binarized, predicted_binarized))
```

The final scores for this prediction are given below.

```
Accuracy:  0.9623655913978495
Recall:    0.7307692307692307
Precision: 1.0
```

3 Bonus: Weka Tools

For the final section of this report, we'll be introducing some useful tools of *Weka* software.

3.1 The Labor Dataset

In the following sections, we'll use the *Labor* dataset.

3.2 Selecting the Dataset

To select the Labor dataset, open the *preprocess* tab and select the dataset according to the figure below.

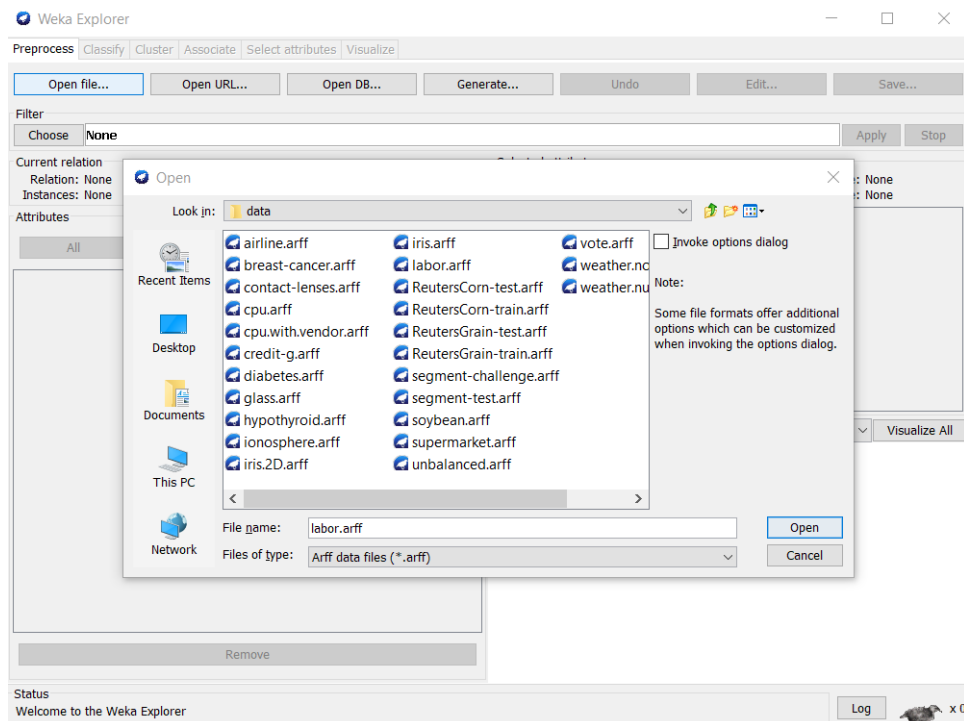


Figure 3.1: Selecting the proper dataset.

3.3 Training J48

Now we open up the *Classify* tab and select the *J48* decision tree classifier from the *trees* section. Here is the performance metrics for this classifier. The confusion matrix generated by Weka is also

```
=== Detailed Accuracy By Class ===
```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.700	0.243	0.609	0.700	0.651	0.444	0.695	0.559	bad
	0.757	0.300	0.824	0.757	0.789	0.444	0.695	0.738	good
Weighted Avg.	0.737	0.280	0.748	0.737	0.740	0.444	0.695	0.675	

Figure 3.2: Performance of J48 on Labor Dataset.

given in the figure 3.3.

```
=== Confusion Matrix ===
```

a	b	<-- classified as
14	6	a = bad
9	28	b = good

Figure 3.3: Confusion Matrix of J48 Classifier on Labor Dataset.

3.3.1 Precision

We can use the following formula to compute the *Precision* out of the given confusion matrix. *b* is considered as *positive* and *a* is considered as negative.

$$PPV = \frac{TP}{TP + FP} = \frac{28}{28 + 6} = 0.82$$

3.3.2 Recall

We can obtain *Recall* of this classifier using the following formula.

$$TPR = Recall = \frac{TP}{TP + FN} = \frac{28}{28 + 9} = 0.75$$

3.3.3 F1-Score

F1-Score can be retrieved as follows.

$$F1_{score} = \frac{2 * PPV * TPR}{PPV + TPR} = \frac{2 * 0.82 * 0.75}{0.82 + 0.75} = \frac{1.23}{1.57} = 0.78$$

3.3.4 Decision Tree Visualization

To visualize the learned decision tree, right click on the item in results list and select the *visualize tree* option.

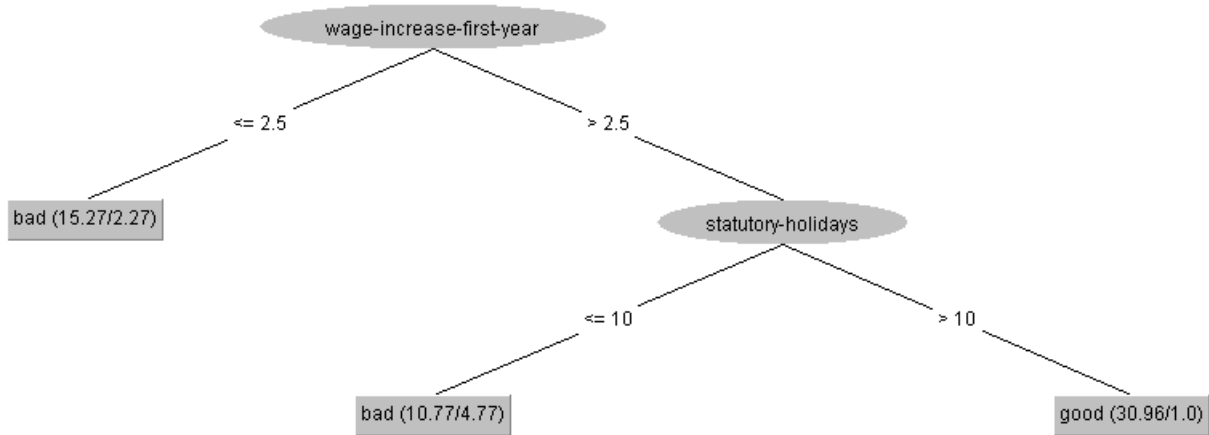


Figure 3.4: Visualization of the J48 Decision Tree.

3.3.5 Testing the Model

To keep things as simple as possible, test data in the figure 3.5 using the learned tree and find the output label of the model.

feature	value	feature	value
duration	1	shift-differential	20
wage-increase-first-year	3	education-allowance	yes
wage-increase-second-year	6	statutory-holidays	12
wage-increase-third-year	4	vacation	generous
cost-of-living-adjustment	tcf	longterm-disability-assistance	yes
working-hours	35	contribution-to-dental-plan	full
pension	ret_allw	bereavement-assistance	no
standby-pay	11	contribution-to-health-plan	half

Figure 3.5: Test Data.

The proper label will be *good*. Note that we have used the product of sum of the tree attributes to find the proper label.

3.4 Training with Decision Stump

The same evaluation results is given in the figure 3.6.

```

=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Class
              0.550   0.054   0.846     0.550   0.667     0.564   0.835   0.815     bad
              0.946   0.450   0.795     0.946   0.864     0.564   0.835   0.851     good
Weighted Avg.   0.807   0.311   0.813     0.807   0.795     0.564   0.835   0.838

=== Confusion Matrix ===

  a  b  <-- classified as
11  9  |  a = bad
 2 35  |  b = good

```

Figure 3.6: Performance of Decision Stump on Labor Dataset.

3.4.1 Precision

We can use the following formula to compute the *Precision* out of the given confusion matrix. b is considered as *positive* and a is considered as negative.

$$PPV = \frac{TP}{TP + FP} = \frac{35}{35 + 9} = 0.79$$

3.4.2 Recall

We can obtain *Recall* of this classifier using the following formula.

$$TPR = Recall = \frac{TP}{TP + FN} = \frac{35}{35 + 2} = 0.94$$

3.4.3 F1-Score

F1-Score can be retrieved as follows.

$$F1_{score} = \frac{2 * PPV * TPR}{PPV + TPR} = \frac{2 * 0.94 * 0.79}{0.94 + 0.79} = \frac{1.4}{1.7} = 0.82$$

References

- [1] Prashant Gupta, *Cross-Validation in Machine Learning*. Towards Data Science, Jun 5, 2017.
- [2] Scikit-Learn, *sklearn.tree.DecisionTreeClassifier*. <http://scikit-learn.org>.
- [3] Scikit-Learn, *Tuning the hyper-parameters of an estimator*. <http://scikit-learn.org>.