# Newspaper Article Classification Contest

Final Report

1 **Terry Li** **Tao Yu**
2 *tianweil@andrew.cmu.edu* *taoyu@andrew.cmu.edu*

3 ## Abstract

4 This is the final report for the Newspaper Article Classification Contest
5 project. The content includes how we retrieve the features from articles,
6 which classifier we used, how to train the classifier and the performance of
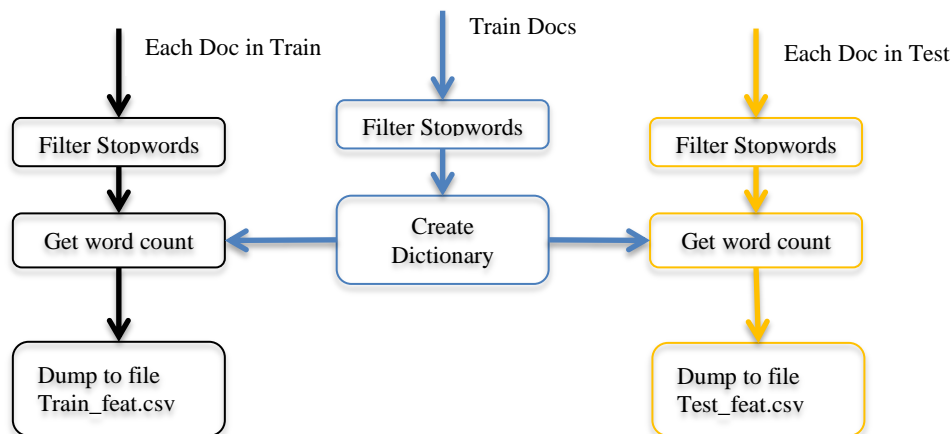7 our classifier.

8

9 ## 1 Project members

10 Terry Li (Andrew ID: tianweil) and Tao Yu (Andrew ID: taoyu)

11

12 ## 2 Feature Extraction

13

14 ### 2.1 Basic Approach

15 We use the word count as the features of our classifier. The process to extract features for
16 train articles and test articles is as below.



17

18 For each article:

19    1. Parse the articles and separate it to a bunch of words

20    2. Remove the stop words

21    3. Count the word count for each word in the dictionary for each article

22    4. Use the word count as the features for classification

23 After filtering the stop words, we have 38800 words left. Then we use these words to classify
24 the articles.

## 2.2 Word Stemming

To reduce the number of features, we leveraged the public Porter Stemming Algorithm [6] to pre-process the documents. After stemming, the number features reduced to 25313.


## 2.3 TF-IDF

To reduce the number of features, we also implemented the tf-idf feature selection algorithm. The original tf-idf algorithm does not take the class information into account. It only calculates the score of a certain word in a bunch of documents. The original equation is

$$idf(word) = \log(\frac{doc\ count}{doc\ count\ that\ conatins\ word})$$

To introduce the class information into this equation, we used a modified version of tf-idf:

$$idf(word|c) = \log(\frac{doc\ count * doc\ count\ in\ class\ c\ and\ contains\ word}{doc\ count\ that\ conatins\ word})$$

For the tf, the equation is the same:

$$tf(word|c) = \frac{word\ count\ in\ class\ c}{all\ word\ count\ in\ class\ c}$$

Because the word count for a certain word in a certain class could be 0, and the doc count that contains a certain word in certain class also could be 0, we introduced add-1 smooth in the process of calculating both tf and idf.

After calculating tf-idf for every feature, we selected 2500 features with the highest scores to be the features of our classifier.


## 3    Naive Bayes Classifier

For Naïve Bayes classifier,

$$C_{NB} = argmax\ P(x_1, x_2, \dots x_n|c)P(c) = argmax\ P(c) \prod P(x_i|c)$$

We use the word count of each to calculate the $P(x_i|c)$ for every class c.

In the beginning, the equation we used is:

$$P(x_i|c) = \frac{N_{ci}}{N_c}$$

In which $N_{ci}$ is the number of times word $x_i$ appears in the documents in class c. $N_c$ is the total number of words appears in documents in class c.

Because a certain word may have zero count in a class, then the $P(x_i|c)$ is 0. This will cause the posterior probability will become 0. Because there are always some words have 0 count in some classes. It makes all posterior probability for every class is 0, which makes our classifier useless.

To solve this problem, we introduced add-1 smooth to the probability calculation. That is we add 1 count for every feature (different word). To make the total probability of all feature still 1, we need to add the number of features to the total word count.

70   Then the equation becomes:

$$P(x_i \,|c) = \frac{N_{ci} + 1}{N_c + n}$$

72
73   In which n is the number of features.

74
75   For the Prior probability of each class, we just use MLE to calculate them.

76
77   After get the $P(c)$ and $P(x_i|c)$ for every word and every class, we finished training the
78   classifier.

79
80
81   # 4      Multinomial Naive Bayes

82
83   Multinomial Naive Bayes models the distribution of words in a document as a multinomial.
84   The likelihood of a document is a product of the probability of the words that appear in the
85   document.

86
87   So the probability of one document belongs to class c is:

88

$$P(c|doc) = P(c) \prod P(x_i|c)^{f_i}$$

89
90   Where the $f_i$ is the frequency of word i in doc.

91
92   To simplify the calculation, we use the log likelihood instead of original probability.

93

$$L(c|doc) = log P(c) + \sum_{i=0}^{n} f_i * \log P(x_i|c)$$

94
95   Then we can classify the doc into the class that has the highest $L(c|doc)$

96
97
98
99   # 5      Modified Naïve Bayes Classifier

100
101   According to paper [5], the performance of Naïve Bayes Classifier is limited by the prior
102   probability $P(c)$ and the probability of word i in each class c $P(x_i|c)$. We tried the
103   modification introduced to Naïve Bayes Classifier introduced by paper [5].

104
105   ## 5.1 Complement Naive Bayes (CNB)

106
107   Instead of calculating the probability of word i in class c:

$$P(x_i \,|c) = \frac{N_{ci} + 1}{N_c + n}$$

108
109   We calculate the probability of word i that is not in class c:

$$P(x_i \,|\tilde{c}) = \frac{N_{\tilde{c}i} + 1}{N_{\tilde{c}} + n}$$

110
111   Then the log likelihood of doc belongs to class c is:

$$L(c|doc) = logP(c) - \sum_{i=0}^{n} f_i * \log P(x_i|\tilde{c})$$

## 5.2 Transformed Naïve Bayes (TNB)

Instead of using of $P(x_i|c)$ as the probability of word i belong to class c, we can also use the tf-idf score of the word i.

Then the log likelihood of doc belongs to class c is:

$$L(c|doc) = logP(c) + \sum_{i=0}^{n} f_i * \text{score}_{\text{tfidf}}(i, c)$$

## 6    Support Vector Machine

We implement a simple version of support vector machine algorithm with kernel parameters of linear, Poly and RBF. For convenience, we choose Quadprog as dual problem solver. We do experiments on linear SVM as well as with RBF kernel as below:

## 6.1 Linear SVM

SVM are shown to handle feature redundancy well, because of the reason that we have 38863 features. It is reasonable to use linear SVM. For features, we use wordcount and TF-IDF features to train the linear SVM model. Every feature will contribute to the improvement of the linear SVM model.

C is essentially a regularization parameter, which controls the trade-off between achieving a low error on the training data and minimizing the norm of the weights.

The C parameter tells the SVM optimization how much we want to avoid misclassifying each training example. For large values of C, the optimization will choose a smaller-margin hyperplane if that hyperplane does a better job of getting all the training points classified correctly. Conversely, a very small value of C will cause the optimizer to look for a larger-margin separating hyperplane, even if that hyperplane misclassifies more points. For very tiny values of C, we should get misclassified examples, often even if our training data is linearly separable.

We adjust C from 2^(-10) to 2^5, and find that the accuracy is highest when C is 2^(-10).

## 6.2 SVM with RBF kernel

For SVM with RBF kernel, we replace natural product with kernel function. So, at the end, w*would look like,

$$w^* = \sum_{i \in SV} h_i y_i \Phi(x_i)$$

and hence,

$$\langle w^*, \Phi(x) \rangle = \sum_{i \in SV} h_i y_i \langle \Phi(x_i), \Phi(x) \rangle$$

158 Similarly,

$$b^* = \frac{1}{|SV|} \sum_{i \in SV} \left( y_i - \sum_{j=1}^{N} \left( h_j y_j \langle \Phi(x_j), \Phi(x_i) \rangle \right) \right)$$

159
160
161 and our classification looks like

$$c_x = \text{sign}(\langle w, \Phi(x) \rangle + b)$$

162
163
164
165 After fine-tuned our learner parameters on evaluation data, we test our fitted SVM's
166 performance on the testing data that's previously unseen in the training and fine-tuning stages.
167 Because of time limitation, we use the grid search: Set a range of feasible values for C, for
168 instance C in [2, 2^2, 2^3, 2^4, 2^5, 2^6]. Set a range of feasible values for Gamma [1, 10,
169 100, 1000]. Then we use grid search to find the best C and Gamma.
170
171
172
173 ## 6.3 SVM Experimental Result
174

| Feature Vector Kernel | Kernel Type | Prediction Success Rate | Kernel Variable |
|---|---|---|---|
| wordcount | Linear | 66.88% | C = 2^(-10) |
| TFIDF | Linear | 66.94% | C = 2^(-10) |
| wordcount | RBF | 67.01% | G = 0.5 |
| TFIDF | RBF | 67.05% | G = 0.5 |

175
176 Compare wordcount with TFIDF features, TFIDF is a little better than original wordcount
177 information.
178
179 Alternatively, we can use n-cross validation to estimate our SVM's performance. If we have a
180 limited amount of annotated texts, n-cross validation is recommended as it takes advantage
181 of using all the data we have.
182
183 The experimental results show that SVM consistently achieve performance on text
184 classification tasks. SVMs eliminate the need for feature selection, making the application of
185 text classification easier. Meanwhile, SVM has the advantage of robustness.
186
187
188
189
190
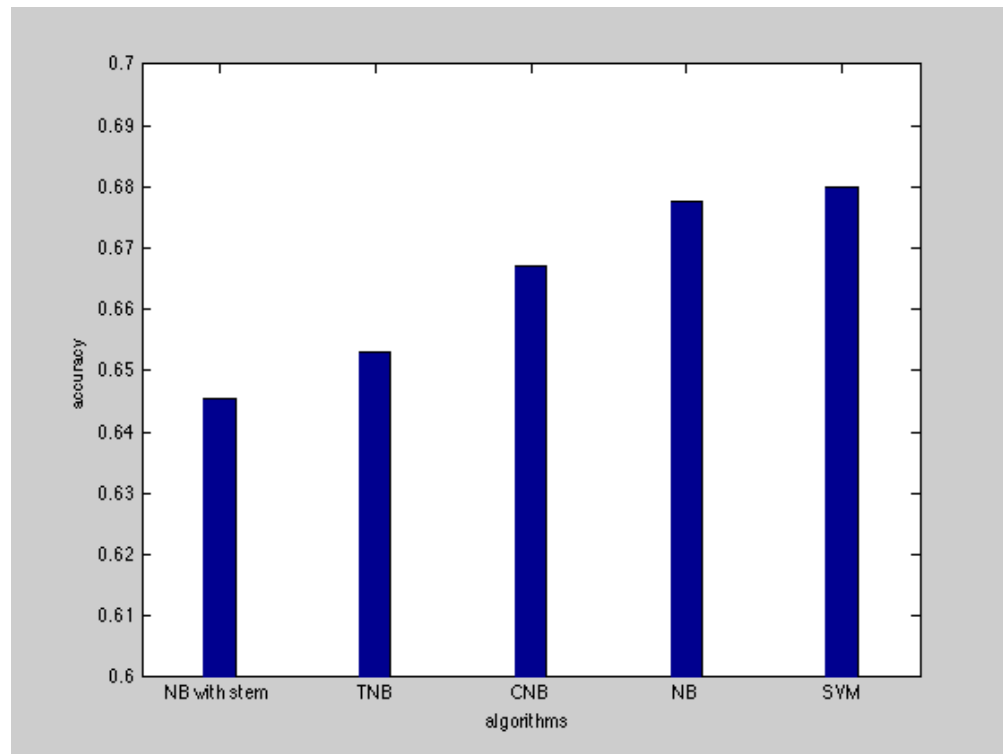191
192
193
194
195
196
197
198
199
200
201
202

203
204 # 7    Performance
205

| Feature selection method | Accuracy |
|---|---|
| Naïve Bayes with stemming | 64.55% |
| Transformed Naïve Bayes | 65.3% |
| Complement Naïve Bayes | 66.7% |
| Basic Naïve Bayes | 66.75% |
| SVM | 67.05% |

206



207
208
209
210 **Performance Analysis:**
211
212 Based on our classification result, Naïve Bayes Classifier does not have overfitting problem,
213 so the more features the better. That's why the basic Multinomial Naïve Bayes Classifier has
214 a better performance than Naïve Bayes with tf-idf.
215
216 Dealing with text classification, we have to consider very many (more than 30000) features.
217 SVMs use overfitting protection, which has the potential to handle these large feature spaces.
218
219 Although there are original 38663 features，chances are that features ranked lowest still
220 contain useful information. As a result, aggressive feature selection may lead to information
221 loss. Naïve Bayes and SVM with all the features achieve good performance.
222
223 When we convert texts into vector representations, they contain only few entries which are
224 not zero. Thus the vectors are very sparse, SVMs is well suited for problems with sparse
225 instances. Meanwhile, most of the text classification problems are linearly separable. SVMs
226 with (linear, RBF) is good at finding such a separator.
227

**References**

[1] http://en.wikipedia.org/wiki/Latent_Dirichlet_allocation

[2] http://nlp.stanford.edu/IR-book/html/htmledition/feature-selection-1.html

[3] Joachims, Thorsten. Text categorization with support vector machines: Learning with many relevant features. Springer Berlin Heidelberg, 1998.

[4] Yang, Yiming, and Jan O. Pedersen. "A comparative study on feature selection in text categorization." ICML. Vol. 97. 1997.

[5] Jason D. M. Rennie, Lawrence Shih, Jaime Teevan, David R. Karger "Tackling the Poor Assumptions of Naive Bayes Text Classifiers"

[6] http://tartarus.org/martin/PorterStemmer/index-old.html

[7] http://www.cs.cornell.edu/people/tj/publications/joachims_98a.pdf

[8] http://www.cs.cornell.edu/People/tj/svmtcatbook/