

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Twitter US Airline Recommendation Prediction

Xiaotong Duan, Tianshu Ji, Wanyi Qian

Abstract

The goals of this project are 1) build a model for six major U.S. airlines that performs sentiment analysis on customer reviews so that the airlines can have fast and concise feedback, 2) make recommendations on the most important aspect of services they could improve given customers' complains. In this project, we performed multi-class classification using Naive Bayes, SVM and Neural Network on the Twitter US Airline data set from Kaggle. Significant accuracy has achieved, which shows that our models are reliable for future prediction.

1 Introduction

In recent years, Twitter has become the de facto online customer service platform. Thus, a company's image on Twitter is of central importance and this is especially true for airlines given that many tweets are travel-related in nature. In fact, research has shown that responding to tweets has revenue generating potential, drives higher satisfaction than other customer service channels, and perhaps most importantly, satisfied Twitter users spread the word. In this project, we use tweets gathered from Twitter to learn about people's flight experiences and give airline companies suggestions on how to make their trip more enjoyable.

The data set contains about 15,000 tweets, collected from February 2015 on various airline reviews. Every review is labeled as either positive, negative or neutral. First, we want to build a model to perform sentiment analysis on the data set. Second, more interestingly, we want to assign a reason to each negative response, such as late flight, lost luggage, etc. In our data set, about 80% of the negative reviews has a negative reason label, yet the rests are labeled as "can't tell". Our goal is to assign a label to this unspecified group. By knowing every review's negative reason, we can give specific suggestions to different airline companies on how to improve their service.

2 Background and Related Work

Nowadays, developing and testing different models for a natural language processing problem is an interesting and challenging task. However, due to the nature of the problem, the accuracy of sentiment analysis on single sentence like movie reviews never reaches above 80% for the past 7 years [1]. Looking at last years project on twitter [2], their accuracy was 59.32% to 63.71%, depending on different models. In our project, we achieved near 20% more than their result, which is a significant improvement.

Since tweets texts are usually short and verbal, the same problem presents in our data set as well. However, even though the tweets are short, there are strong indicative words. Specific words can be used as indicators for spam/ham emails and achieve good test accuracy. Therefore, we believe that tweets review, without many negating negatives, can be predicted well using the frequency vector representation. To prove this, we will use Recurrent Neural Network model and the GloVe word vector [3] to compare the result.

3 Approach

3.1 Dataset

The sentiment analysis labels are positive(20%), negative(60%), and neutral(20%). The negative reason labels are bad flight(7.45%), canceled flight(9.62%), customer services issues(39.77%), damaged luggage(0.84%), flight attendant complaints(6.05%), flight booking problems(6.19%), late flights(1.99%), long lines(19.97%), and lost luggage(8.23%).

3.1.1 Preprocess of Dataset

In the preprocessing step, non-English word, symbols and website links are eliminated. Then the whole data set is randomly separated into training set (10000 samples, 70%) and test set (4636 samples, 30%).

3.1.2 Dictionary

The dictionary is made based on the training data and all sentences are broken down into list of words: (1) Delete common words such as a, an, to, of, on etc. with high frequency but little semantic usage. (2) Stem words, such as "thanks" and "thank" as one word. (3) Delete low frequency words that appear once to reduce the size of dictionary for calculation efficiency.

3.1.3 Frequencies Matrix

A feature matrix is built to convert the textual information into numerical information. In the feature matrix, the number of rows indicates the number of samples, the number of columns is the length of the dictionary, and each element indicates whether the specific word has appeared in the current review, 1 for existence and 0 for absence.

To get a sense of correlation presented in our feature matrix, i.e. "bad" and "suck" may have a higher chance to present together, we perform PCA to capture the variance. The result shows that for the first component, variance explained is 2.3%, and for the next nine components, the variance explained is all around 1.0%. This shows that there isn't significant correlation between words and to achieve better accuracy, we include all the words in the dictionary. We propose that the lack of correlation comes from the nature of the text data. Most of them are very short sentences and extremely verbal.

3.2 Models

1. Naive Bayes with multinomial event model from sklearn is used. Input is the frequency vector and Laplace smoothing is used.

$$\hat{\theta}_{yi} = \frac{N_{yi} + \alpha}{N_y + \alpha n} \quad (1)$$

2. Support vector machines with linear kernel and RBF kernel are used in this project. SVM uses the same input and implementation package as Naive Bayes.

$$K(u, v) = u^T v \quad (2)$$

$$K(u, v) = \exp\left(-\frac{\|u - v\|^2}{2\sigma^2}\right) \quad (3)$$

3. Neural Network

Tensorflow is used in implementation. Input is the frequency vector that represents a review. The output is a vector with probabilities for different classes and the highest is selected as prediction. Label is a one-hot vector that represents the class. Loss function is cross entropy plus a regularization term. The vanilla Neural Network that we use:

$$h = Wx + b \quad (4)$$

$$\hat{y} = \text{softmax}(Wh + b) \quad (5)$$

108 4. Recurrent Neural Network

109 A Bi-directional Gated Recurrent Unit Network (GRU) can capture the structure features
 110 of a sentence. Also, it solves the vanishing gradient problem which many recurrent neural
 111 network models have. Bi-directional GRU is commonly used in text analysis, which we
 112 want to compare with our models. Package scikit is used for implementation. In GRU, word
 113 vectors, instead of frequency vector, will be used and we choose glove.twitter.27B.zip.[3]
 114 These are pre-trained word vectors that are trained on twitter data set. The math for GRU
 115 is shown as follows:

116 For right direction $\vec{h}_t^{(i)}$

117
$$\vec{z}_t^{(i)} = \sigma(\vec{W}_i^{(z)}x_t^{(i)} + \vec{U}_i^{(z)}h_{t-1}^{(i)})$$
 (6)

118
$$\vec{r}_t^{(i)} = \sigma(\vec{W}_i^{(r)}x_t^{(i)} + \vec{U}_i^{(r)}h_{t-1}^{(i)})$$
 (7)

119
$$\tilde{h}_t^{(i)} = \tanh(\vec{W}_i x_t + r_t \circ \vec{U}_i h_{t-1})$$
 (8)

120
$$\vec{h}_t^{(i)} = z_t^{(i)} \circ h_{t-1}^{(i)} + (1 - z_t^{(i)}) \circ \tilde{h}_t^{(i)}$$
 (9)

121 Similarly for $\overleftarrow{h}_t^{(i)}$

122 Output

123
$$y_t = \text{softmax}(U[\vec{h}_t^{(top)}; \overleftarrow{h}_t^{(top)}] + c)$$
 (10)

131 **4 Experiments & Results**

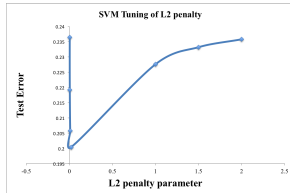
132 **4.1 Sentiment Analysis**

133 **4.1.1 Naive Bayes Classification With Laplace Smoothing**

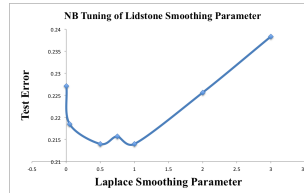
134 α (smoothing parameter) is tuned. The lowest test error (0.214) is achieved when α is 0.5 or 1.
 135 $\alpha = 1$ is used in further experiments.

136 **4.1.2 Support Vector Machine With Linear Kernel**

137 Before tuning the regularization, SVM with linear kernel, RBF kernel results in 0.23, 0.21 test error,
 138 respectively. Therefore, SVM with RBF kernel is excluded from future tuning due to the higher
 139 initial test error. L2 regularization is used to avoid overfitting. According to the graph shown below,
 140 the lowest test error (0.200) is achieved when L2 regularization is 0.02.



142 (a) SVM Tuning Result



144 (b) NB Tuning Result

145 Figure 1: Sentiment Analysis test error

146 **4.1.3 One layer neural network**

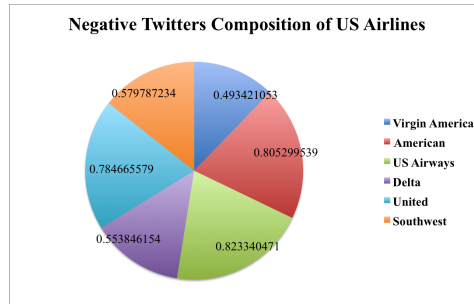
147 In each stochastic gradient descent step, only a batch of 100 samples is used in SGD to increase the
 148 training speed. Tuning parameters are learning step and regularization term, which are 0.01 and 0
 149 respectively. The best test error for this one layer neural network is 26.3374%

162 **4.1.4 Bi-directional Gated Recurrent Unit Network**

163
164 Word vectors from GloVe with a dimension of 50 will be used in GRU. A 2-layer GRU has a test
165 error of 26.5%. A 3-layer GRU has a test error of 25.6%. Learning step is 0.01; l2 is 0.

166
167 **4.1.5 Sentiment Analysis Result**

168 In sentiment analysis task, SVM with linear Kernel achieves the best test accuracy. Therefore, SVM
169 is recommended in this section. According to the result from linear SVM, Virgin American performs
170 the best according to its lowest negative review composition in its total reviews.
171



184 **4.2 Negative Reason Prediction**

185 In this section, the goal is to determine the most negative reason on flight services. All the negative
186 reviews have been collected for this task, and separated into labeled set and unlabeled set. We will
187 make predictions on the unlabeled set.
188

189 **4.2.1 Naive Bayes Classification**

190
191 Using Naive Bayes Classification, the test error is average to 29.26% after ten-fold cross-validation,
192 with a Laplace smooth factor of 0.5.

193
194 **4.2.2 Support Vector Machine**

195 L2 regularization is tuned. The best test error for SVM is 32.82%, when l2 regularization = 0.03.



Figure 2: Negative reason test error

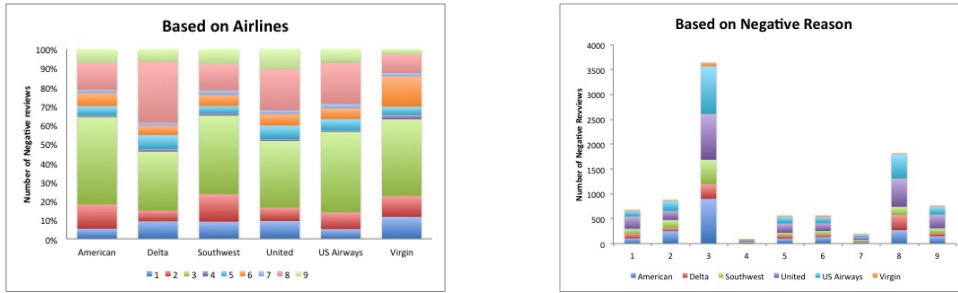
210 **4.2.3 One layer neural network**

211 Input is same as before. Label's dimension changes to 9. Learning step and regularization term are
212 0.01 and 0 respectively. The test error is 37.82%

213 **4.2.4 Negative reasons classification**

214 Given by the lowest test error, Naive Bayes is used for the prediction of unclassified data. Result is
215 shown below that most complaints are on customer service. One postulate might be due to the high

216 volume of contact. Since various reasons can lead to calling customer service. Thus correlations
 217 between classes may play a factor in determining this result.
 218



220
221
222
223
224
225
226
227
228 (a) Major negative reasons for each airlines (b) Airlines performance on service issues

229
230 Figure 3: Negative Reason Classification

231
232

Airline	1st Negative Reason	2nd Negative Reason
Virgin America	Customer Service Issue	Flight Booking Problems
United Airlines	Customer Service Issue	Late Flight
Southwest	Customer Service Issue	Late Flight
Delta	Late Flight	Customer Service Issue
US Airways	Customer Service Issue	Late Flight
American Airline	Customer Service Issue	Late Flight

233
234
235
236
237
238
239
240
241 **5 Conclusion**

- 242
- 243 • It is pleased that our vectors work. It is surprising that SVM and Naive Bayes perform
 244 better than deep learning methods. And the accuracy is very high, 80%. We think the
 245 reason behind this is that while movie reviews have a lot of sarcasm[1], which is very
 246 difficult for any model to grasp, twitter reviews are much more straight forward, and thus
 247 most of the sentiments are expressed directly at the word level. That is to say, with specific
 248 word appearance, sentiment is indicated clearly, which justifies our feature representation
 249 using frequency vector. It is possible to judge a twitter airline review's sentiment only by
 250 identifying positive words in a review. Therefore, given the nature of our data set, the task
 251 can be solved at bag-of-word level well.
 - 252 • However, it is too early to say that neural network can not perform better than bag-of-word
 253 models. The frequency vector used in vanilla neural network is so large that takes enormous
 254 time to train, roughly 6 hours for 10,000 iterations now. Therefore, clever ways of reducing
 255 frequency vector size are needed. Meanwhile, better tuning parameters can be figured out
 256 once training time is significantly decreased.
 - 257 • Another possible reason is that for recurrent neural network, GRU in our project, labeling
 258 every node is very important. While this model can achieve as high as above 80% accuracy
 259 using Stanford Sentiment Tree Bank dataset[4], Our results show that without sufficient
 260 labeling, this model is not able to achieve an accuracy above 80%, which means RNN
 261 family needs strong supervision. However, most of the online reviews and other documents
 262 only have limited labels. Better labeling algorithm on new data set should be thought about
 263 in future work.

264 **6 Reference**

- 265
266
267
268
269
1. Pang, Lee. 2005. CS224D Slides.
 2. Yuan, Zhou. Twitter Sentiment Analysis with Recursive Neural Networks.
 3. Pennington, Socher, Manning. GloVe: Global Vectors for Word Representation
 4. Stanford Sentiment Tree Bank <http://nlp.stanford.edu/sentiment/treebank>