

# Classification of Optional Practical Training (OPT) comments using a Naive Bayes classifier

Anand  
a3anand@ucsd.edu

Sampath Krishna  
svelaga@ucsd.edu

Jorge A. Garza  
Guardado  
jgarzagu@ucsd.edu

Adithya Apuroop K  
akaravad@ucsd.edu

## ABSTRACT

This project aims to classify the optional practical training comments using a naive Bayes classifier. We demonstrate the effectiveness of the naive Bayes approach and further enhance its performance using a simplified form of an expectation maximisation algorithm. We explore how sentiments change over time, and also provide preliminary results that help in understanding how sentiments vary with ethnicity.

## 1. INTRODUCTION

OPT is a scheme in which students with F-1 visas are permitted by the United States Citizenship and Immigration Services (USCIS) to work for at most one year on a student visa towards getting practical training to complement their field of studies. On April 2, 2008, the department of homeland security(DHS) announced an extension to the OPT which was passed by USCIS as an interim rule. This rule allows students in Science, Technology, Engineering or Math (STEM) majors, an OPT extension for up to 17 months.

In August 2015, a US federal court gave its verdict on a lawsuit challenging the 17-month OPT STEM extension. The court has decided that the interim rule was deficient as it was not subjected to public notice, comments and opinions and it vacated the 2008 rule allowing the 17-month extension. However, a stay was put in place until February 12, 2016. DHS will have until then in order to take action regarding the fate of the STEM extension program. This rule was open to public comments for a one month duration, ending on Nov 18<sup>th</sup>. The comments are publicly available at [1].

## 2. THE DATA SET

### 2.1 Data collection

Data was collected from the Department of Homeland Security (DHS) web page forum [1] containing at the time, 42,925 comments. This data was obtained over a period of 30 days, ranging from 19th October to 18th November. The DHS web page provides a CSV file containing all user names and comments, but the comments are stored as a web page link. A script was written to download and then parse each web page containing the comment for each user and the resulting data was stored in a JSON file. The data we used contained the fields: 'userName', 'comment', 'docID', 'receivedDate', 'postedDate'

### 2.2 Dataset Preprocessing

As a pre-processing step, we removed all the punctuations from the words. We also changed all words to lower case, although a more rigorous model could make use of the upper case information to identify stronger sentiments. Finally, all the common stop words were removed as they convey little meaning.

### 2.3 Dataset Labeling

Since the original dataset is unlabeled, we manually labeled the first 900 comments as *support* or *oppose*. Out of these, the first 600 were used for training, comments from 601-700 constituted the validation set and 700-900 were used for testing. We used validation set to pick the best possible model from a pool of possible models.

### 2.4 Data visualisation/Exploratory analysis

Figure (see Fig. 1 and Fig. 2) contain the word clouds of the most common words (after removing stop words) on the train data for positive and negative labels. Some of the most commonly found words in supporting comments were {*benefit, support, economy, STEM, international, students, good*} etc meaning that people supporting the OPT extension feel that the extension will benefit the economy and is good for international students. While in opposing comments we found words like {*American, job, worker, student, foreign, program*} etc meaning that people opposing are concerned about the jobs being taken away by the foreign students.



Figure 1: Positive comments word cloud



Model	Train acc	Validation acc
Unigram	98%	86%
Bigram	99.5%	86%
Unigram+Bigram	99.33%	88%

**Table 1: Table showing the validation errors on the 3 schemes being considered**

model. While the naive Bayes in itself performs reasonably well, its performance can be boosted by augmenting it with a simple fix.

#### 4.1 Semi supervised estimation

It has been suggested by the authors in [3], that in cases where the number of training examples is small, the performance of the naive bayes classifier can be improved by combining it with an expectation maximization algorithm. In short, the authors suggest to do this :

1. Predict the class probabilities  $P(class|data)$  for all examples in the dataset
2. Retrain the model based on *class probabilities* estimated in the previous step

The first step above is an expectation step in disguise, and the second step corresponds to the maximisation. Although the second step requires us to retrain the model based on the probabilities in the previous step, we use a relaxed version of this step as follows:

1. Predict the class probabilities  $P(class|data)$  for all examples in the dataset
2. Retrain the model based on *class labels* estimated in the previous step

This algorithm, which we'll refer to as *classification maximisation (CM)* algorithm is a convenient approximation to the more rigorous expectation maximization. What this means, is that we use the predicted labels as the actual labels and retrain the model based on these labels until convergence. These iterations significantly improve the accuracy of the naive Bayes model by incorporating the knowledge from the large pool of unlabeled examples. Refer to Table 2 to see the performance comparison of the classification maximisation and naive Bayes algorithms. Note that the classification maximisation algorithm achieves significantly better TPR and TNR on the test set as compared to naive Bayes. Also note that TNR is a particularly important term to evaluate the performance of the classifier, as the negative examples are relatively rare and one would want to classify them correctly. After all, an "all positive" classifier would achieve an accuracy of about 85% on this dataset.

## 5. RESULTS AND OBSERVATIONS

From the predictions of the classification maximisation algorithm, we find that approximately 85.17% of the users support OPT extension while 14.83% oppose it.

### 5.1 Excerpts of comments from the labeled set

There are certain cases where our naive Bayes model fails to predict the sentiment correctly. Consider a false negative classification in our test set :

*"OPT is helping to find better workers for the jobs, not simply give the jobs to foreigners."*

The classifier recognises the words *jobs* and *foreigners* as predictive of negative sentiment but doesn't notice the negation in the original clause.

Also, a complex viewpoint expressed via contrast and juxtaposition stumps our classifier. Consider the comment :

*"Admittedly, there are Americans who can not find a job. But there are also foreign students who can not find a job. The majority of US companies already give priorities to US workers. As a result, the unemployment percentage of international students is already higher than that of native Americans. It is unfair to say that more US workers can not find a job. We should compare the percentage instead of the absolute number."*

which actually supports the OPT proposal but is predicted to be a disapproval since the bag-of-words approach lumps *jobs*, *Americans* and *workers* with negative sentiment. The (+,-) log likelihood is (-307.6, -305.9) : which indicates an edge case for our classifier.

There are a few false positives as well. Let us evaluate a straightforward negative comment which manages to fool the classifier :

*"I oppose the extension of OPT. Schools, especially public schools, welcome foreign students because they pay high tuition. And then, with extension of OPT, they earn back what they invest and maybe much more. Who is the winner? Obviously, foreign students who get more than they invested, schools which get a lot of money and companies which get a lot of comparable cheap workers. Who is the loser? Obviously not the government, the working Americans are loser, the middle class is loser. They take risk of losing their jobs but they don't get any benefit from having more and more foreign students."*

Here the classifier is incapable of parsing the topic sentence but counts phrases like *welcome foreign students* and *lot of money* towards a false positive prediction. The log likelihood for positive and negative prediction are -505 and -513 respectively. The word *oppose* is present in our list of words which appear only in negative comments. However the frequency is a weak 560 which doesn't sway the negative probability sufficiently.

However, for a majority of comments, the bag-of-words approach combined with iterative maximisation works surprisingly well. We will next look at a few comments drawn at random from the unlabeled dataset and see how our classifier performs.

### 5.2 Excerpts of comments from the unlabeled set

Here are some comments from the unlabeled set. We mention the log likelihood for positive and negative predictions alongside the comment. A higher log likelihood implies a greater probability for the classification.

Consider comments like

Model	Train acc	Validation acc	Test acc	TNR	TPR
naive Bayes(Uni)	98%	86%	90.5%	40.7%	98.2%
naive Bayes(Uni+Bi)	99.3%	86%	88.5%	29.7%	97.7%
CM Algorithm(Uni)	96.33%	97%	95.5%	<b>92.5%</b>	96%
CM Algorithm(Uni+Bi)	96.5%	96%	96.5%	<b>92.5%</b>	97.1%

Table 2: Table showing the comparison of various schemes

*"Please stand up for American citizens and say NO to this travesty."*

with a (+,-) log likelihood of (-67.3, -57.9) : clearly classifying it as a negative comment. True positive comments like

*"International students bring money, skills and jobs in USA. This rule is not taking away any job from us. In fact because of this rule more and more jobs are being created for American citizen with or without degree in STEM field."*

show a wide margin between the pos/neg log likelihoods of -485.7 and -518.7 respectively.

*"I oppose the Department of Homeland Security's proposed rule that would expand the Optional Practical Training program. This expansion would allow U.S. tech companies to hire ... a de facto shadow H-1B program, in violation of Congressional intent."*

(+, -) log likelihood = (-961.6,-819.0). is also classified correctly.

With a test accuracy of 0.965, the classification maximisation approach combined with naive Bayes performs competitively compared to other advanced techniques like neural networks or support vector machines.

### 5.3 Other interesting observations

We also tried to analyse the data based on the ethnicity of the users. We were curious to know how people supported/opposed OPT based on the country of origin. For this, we have collected the publicly available common first names and surnames of Americans and Chinese ethnicity. We couldn't get the corresponding data for India to perform such an analysis. The results are summarized in Fig 4

It was initially quite surprising that a significant fraction of Americans support the OPT. There might be several reasons for this. Firstly, the database for American names contains many foreign names as well, and this might have created conflicts with true foreigners who supported the extension. Secondly, all American names are not in the database and that might have resulted in some conflicts too. Nevertheless, upon examination of some reviews, we found that many Americans were supportive of OPT due to the positive talent it brings to the country.

We ran the classifier on the entire corpus of around 42,000 comments and plotted the distribution of sentiment over time. The graph [5] suggests that the initial sentiment was overwhelmingly positive with negative comments beginning to trickle in a week after the voting started.

An interesting result was obtained from trying to get the most negative words that do not occur as often or are not

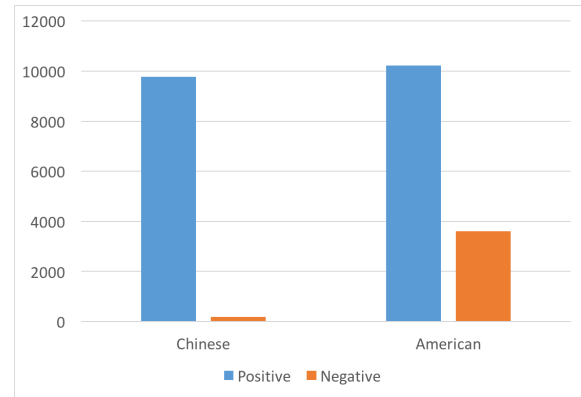


Figure 4: Sentiment breakdown by ethnicity

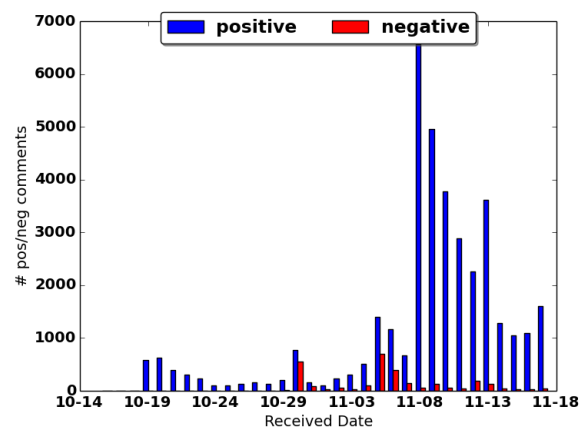
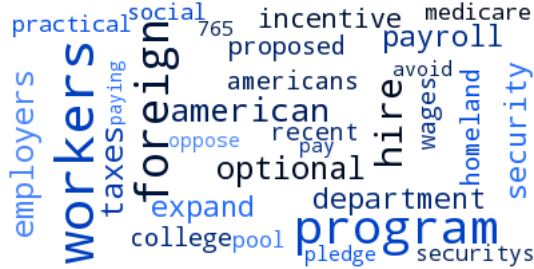


Figure 5: Sentiment breakdown over time

even listed in the positive words list. For doing so, the entire dataset was segregated based on the predicted labels, and the frequency of each negative word was subtracted from its matching positive word frequency after normalization. If a negative word occurred a large number of times in the positive word list, it was removed. Resulting words from this list can be seen in Fig. 6



**Figure 6: Negative words which do not appear as much in positive comments**

From Fig. 6 it can be noticed for example, that words like {workers, program, foreign, oppose, homeland, taxes, wages, medicare, taxes, paying} appear in the negative only words list. These words represent that most people that **oppose** are concerned mainly with **foreign** people stealing jobs from American **workers** in their own **homeland** with this new **program** and foreigners not paying **taxes**. An interesting thing, was to see the number “765” which refers to the form I-765 that needs to be filled when applying for an OPT. Another interesting thing is that words like “medicare” show up in the negative only word cloud, indicating that people opposing are also somehow concerned with this. For example, consider this comment where words like {wage, medicare, taxes, pay, social, security} appear.

*“American IT jobs should be done by natural born Americans, not foreigners, who will work for substandard wages and be exempt from the taxes that are paid to help support our economy and social security and Medicare.”*

Our code is currently hosted at [8].

## 6. ACKNOWLEDGEMENTS

This project was inspired by a similar analysis done in [2].

## 7. REFERENCES

[1] <http://www.regulations.gov/#!docketDetail;D=ICEB-2015-0002>

[2] <https://medium.com/@heretic/on-opt-optional-practical-training-10ced7051066#.goc1w933d>

[3] Nigam, Kamal; McCallum, Andrew; Thrun, Sebastian; Mitchell, Tom (2000).“Learning to classify text from labeled and unlabeled documents using EM”

[4] <http://immigrationgirl.com/breaking-news-on-opt-stem-extension-court-says-uscis-rule-allowing-17-month-stem-extension-is-deficient/>

[5] <http://vikeshkhanna.webfactional.com/opt>

[6] [http://ieeexplore.ieee.org/xpl/freeabs\\_all.jsp?arnumber=1359749&reason=concurrency](http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=1359749&reason=concurrency)

[7] <http://dl.acm.org/citation.cfm?id=288651>

[8] <https://github.com/ananducsd/opt-data-mining>

[9] Aurangzeb Khan, Baharum Baharudin, Lam Hong Lee\*, Khairullah khan.“A Review of Machine Learning Algorithms for Text-Documents Classification”

[10] <http://www.time.mk/trajkovski/thesis/text-class.pdf>