# Public Opinion Analysis on Politics
## Prediction on Selection by NYT's Editors

Junqian Zhang - 944556

[1] University of Milan
[2] Data Science and Economics

**Abstract.** The Public comment the articles of New York Times to present their opinions, but not all the comments will be selected by the editors. This project analyzes the public comments on Politics to get clues to what New York Times considers worth promoting.

**Keywords:** Opinion Analysis · Aspect Extraction · Sentiment Orientation · Text Classification.

## 1 Introduction

The editors of New York Times will not pick all the comments from the public. The aim of this project is to predict whether a comment can get the editors' selection under the topic of *Politics*. This work wants to explore the relation between the Public's opinion and editors' sentiment. Opinion analysis is an active research field, and thanks to great many researchers' work, plenty of methods for opinion analysis have been proposed. This work is aspect-based analysis and tries to use the sentiment orientation of comments on different aspects to perform the binary classification.

## 2 Problem Statement and Methodology

The main challenge in this work is to find the relation between the sentiment of one group - the Public and the sentiment of the other group - the editors. The sentiment of the Public is represented by the sentiment orientation matrix, and the sentiment of the editors are represented by whether they pick a comment. The analysis regards the frequent nouns or noun phrases as aspects, which are all explicit aspects and detects the sentiment orientation of comments on these aspects to build sentiment orientation matrix. This work tries to solve the problem by performing classification using the sentiment orientation matrix. The whole procedure is implemented as the following five steps.

### 2.1 Text Cleaning

The dataset used for analysis is a handful of comments. For the following steps, all the comments are cleaned based on sentence-level. All the sentences are tokenized and words in the sentences are executed lemmatization to get uni-grams

by **Part-of-Speech** Tagging. However, at the same time, stop words are not abandoned.

## 2.2   Aspect Extraction

Under the *Politics* topics, people are discussing different sub-topics, which are called as *aspects* in this paper, such as tax. This project focus only on explicit aspects based on frequency which are nouns and noun phrases appearing in the comment body. In this case, all the noun phrases are assumed to be composed of two words, i.e., all the noun phrase are bi-grams.

Based on the uni-grams, bi-grams are built at sentence-level to find candidate noun phrases. This project assumes that the candidate noun phrases should be in the pattern of two nouns, or stop word plus noun or adjective plus noun. So, to efficiently get the bi-grams, bi-grams without nouns inside are abandoned. To obtain the valid noun phrases and rank them, I calculate Pointwise Mutual Information (**PMI**) scores for all the bi-grams. PMI values take into account of the correlation between the two words inside the noun phrase, avoiding the case where the noun phrases in fact are the aspects of the aspects. Finally, 50 bi-grams with highest PMI scores are selected but the meaningless ones among them are abandoned.

Except these selected noun phrases, all the other bi-grams are splitted into uni-grams again and corpus is recleaned the corpus by removing the stop words. With comments as documents, I calculate the term frequency–inverse document frequency (**TF-IDF**) for all the nouns and candidate noun phrases. 20 features with highest TF-IDF values are selected as aspects.

## 2.3   Aspect Categorization

The aspects found in last step are the explicit aspects for this analysis, but these aspects have other expressions because different people may have different describing habits. So the target of this step is to group aspect expressions into aspect categories.

Assuming that aspects expressions who belong to the same category, have the same context, it is necessary to compare the context of words and phrases. A **word embedding** model can be trained to represent all the uni-grams and bi-grams as vectors, and thus it is possible to compare the context between each other. These vectors can be seen as a description of the context of each element in the vocabulary. **Skip-gram** model is used to get the input for the word embedding. Through skip-gram, all the sentences in all the comments are organized into sequences. Each sequence has 3 grams and each gram concerns the 2 grams surrounded as neighbours.

As all the words and phrases are in a numeric form, the cosine similarity distance between them can be calculated to see the difference of their contexts. With the word embedding model, for each aspect category extracted in the last step, I take the top 5 words or phrases that are closest to it, which are regarded as different aspect expressions. But some expression may be not nouns or noun phrases, so these are removed from the category.

As a result, the aspects extracted from last step are seen as different aspect categories. For each aspect categories, different aspect expressions are found by word embedding model and skip-gram model.

### 2.4   Sentiment Orientation

This paper uses the sentiment orientation matrix as a tool to represent the opinions from the Public. Rows of *sentiment orientation matrix* represent comments, and columns represent aspects. The values of matrix are among **-1**, **0**, and **1**, representing *positive*, *neutral* and *negative* respectively. To get the sentiment matrix for final classification, sentiment orientation of each comment on each aspect category should be detected.

The sentences, who contains at least one aspect, are called *opinion sentences* here. To be clearer and more specific, for each opinion sentence, I tag it with aspects it contains. For each comment, if it does not contain any opinion sentence on one aspect, then **0** is assigned to it on that aspect, regarding as *neutral*. For each aspect, all the comments with opinion sentences on this aspect are collected and these comments abandon the sentences which do not contain this aspect. Since there is no sentiment orientation label and comments show only two different sentiment orientations, *positive* and *negative*, **K-Means** is implemented to cluster the comments into 2 groups. Based on the assumption that comments on the same side tend to use similar vocabulary, **Bag-of-Word** model is used to represent the comments by **TF-IDF** approach. The two comments which are closest to centroid of each cluster, are called *center comments*. The two center comments represent two sentiment polarities of the aspect. I check the text of center comments, and label them *positive* and *negative* manually. Then the comments in the same cluster are labeled after the center comment. In this way, **-1** is assigned to comments which are labeled with negative on one aspect, and otherwise **1** is assigned to comments which are labeled with positive.

### 2.5   Text Classification

The final step is to build the relation between the sentiment orientation matrix and editor's selection by classification. This step uses **Random Forest** as tool.

## 3    Experiment Results

### 3.1    Data Description

The data analyzed in the project are the comments made on the articles published in New York Times in April 2017 under the section name of *Politics* from Kaggle[1]. The data contain 38,381 comments from 51 articles. The features of data used for analysis are the followings:

- `sectionName` describes the topic of the articles commented, only *Politics* are selectd in this analysis
- `commentBody` stores the contents of the comment
- `editorsSelection` describes whether a comment is picked by editors, and is the target variable, and **0** for not picked, **1** for picked

### 3.2    Results and Evaluation

After implementing all the procedures introduced above, I successfully get 13 aspect categories and their corresponding aspects expressions, as the following listed.

**Table 1.** Aspect category and its aspect expressions (Ordered by first letter)

| Aspect Category | Aspect expression |
|---|---|
| care | care, healthcare, comprehensive, health care |
| country | country, nation, usa, planet, america |
| court | court, scotus, sc, supreme court, bench, overturn |
| democrat | democrat, dems, democrats, dem, legislator |
| gop | gop, the gop, republicans |
| house | house, visitor, occupant, wh, emanate, whitewater |
| obama | obama, gall, president obama, obamas, predecessor |
| people | people, folk, constituent, illegals, person, americans |
| president | president, potus, duly, presidency, occupant |
| state | state, governor, assistant, county, commerce |
| tax | tax, amt, income tax, simplify, monies, loophole |
| trump | trump, donald, dt, djt, donnie |
| vote | vote, reelect, runoff, slim, stein, ballot |
| world | world, region, the world', beacon', eu |

The comments with no opinion sentence are abandon, **25,898** are used for classification task. And **25,110** of them are **not** selected by editors, while **788** are selected. Due to the fact of *ubiased* labels, to prevent overfitting and bad performance on *True Positive*, this paper uses **Adaptive Synthetic Sampling**

---

[1] Data source: https://www.kaggle.com/aashita/nyt-comments

(ADASYN) to get a more balanced dataset. After adjustment, the data therefore is transformed into **25,110** not selected and **25,004** selected.

Take all the data as input, the classification result is as followed:

**Table 2.** Classification Reprort

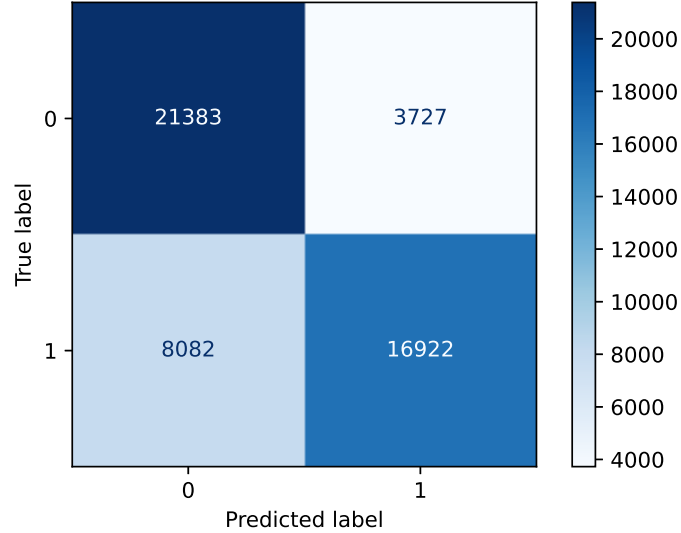| Labels | Precision | Recall | F1-Score | Support |
|--------|-----------|--------|----------|---------|
| 0 | 0.73 | 0.85 | 0.78 | 25110 |
| 1 | 0.82 | 0.68 | 0.74 | 25004 |



**Fig. 1.** Confusion Matrix

The evaluation strategy to evaluate the performance of the model is a 10-fold Cross Validation. The final score is **0.754**.

## 4   Concluding Remarks

The *Accuracy*, *Precision*, *Recall* and *F1-Score* are not low for the classification task, meaning that the sentiment orientation from the Public on different aspects has correlation with editor's selection idea. From the confusion matrix, we can

easily find that most of error come from *False Negative*. This partially results from the small sample size of selected comments. Although the input has been adapted, it is not enough to solve the unbiased problem.

This project does not concern any implicit aspects. For the future work, topic model or sentiment words for detecting implicit aspects can be implemented to improve the performance of the prediction by adding more features. Another thing needs to improve is the sentiment orientation judgement. In this work, the number of labels for sentiment orientation is small and thus it brings uncertain to the analysis. To be more precise on the polarity, other information retrieval tools can be exploited. In addition, this work does not reveal what the editors sentiment orientation on different aspects is. It is necessary to detect their opinions because that's the direct reason why they selected a comment or not. What's more, for better classification models for this problem, classifier based on Naive Bayes or support vector machines can be trained to compare.

## References

1. Mining Hu and Bing Liu, "Mining and summarizing customer reviews", *Proceedings of the 10th ACM SIGKDD International conference on knowledge discovery and data mining*, 2004
2. Bing Liu, Sentiment Analysis and Opinion Mining, Morgan  Claypool Publishers, May 2012