

Content Analysis For Online MultiMedia

Afshin Moin (afshinm)

Peng Seng Kuok (pkuok)

Abhishek Bharani (abharani)

1 INTRODUCTION

Content sharing has been the key ingredient of the biggest success stories in the last decade. Social networks like Facebook and LinkedIn, video sharing services like YouTube, Ecommerce websites like Amazon and online travel agencies like TripAdvisor and Expedia are all examples of online platforms that provide a convenient way for their users to exchange content, goods or services. It is obvious that the success of such online platforms depends on how effectively they link the end users to relevant content. This latter cannot be achieved without proper content classification and user feedback analysis. This motivated us to apply machine learning techniques to automate the process of user feedback analysis and content classification. We work with a real sample dataset taken from YouTube [3]. Nevertheless, same techniques can be applied to any other database with similar structure like Yelp and Netflix. In Section 3, we discuss our approach to data labeling and feature extraction. In Section 4, we review our approach to sentiment and category classification of YouTube comments. Results of the experiments for different classification techniques are compared with each other in Section 5. Also, calibration curves are presented and effect of dimensionality reduction on accuracy is examined. Section 6 concludes the report.

2 RELATED WORK

In [1] recommendation of tags is done automatically by training a classifier that maps audiovisual features from millions of you-tube video to supplied tags by an uploader of video. The system learns a vocabulary of tags and suggests tags which are relevant to the video. A detailed analysis of usefulness of comment is presented in [2] including influence of comment sentiment on rating of the comment using SentiWordNet thesaurus. They also predict the community acceptance of an unrated comment using svm classification and term based representation of comments to categorize a comment as accepted or not accepted. Other work on sentiment classification and opinion mining [3] deals with assigning positive, negative or neutral sentiment to a text using text orientation and linguistic features. In our project we do sentiment analysis using four different classification models. In addition to this we do Category prediction based on comments and tags. [1]

3 DATASET AND FEATURE EXTRACTION

In this section we discuss the properties of the dataset we use for our experiments as well as feature extraction techniques adopted.

3.1 YouTube Dataset

We apply classification techniques on a sample dataset taken from YouTube [3]. This dataset includes 691407 comments for 200 most trending YouTube videos in US in a two weeks period. The same source also offers a similar dataset for the most trending videos in GB. We will use a fraction of this dataset as our test set for sentiment analysis. The data includes comment data and video data. The comment data maps video IDs to comments and the number of likes and dislikes for the corresponding comments. The video data contains video title, channel title, category, tags, number of views,

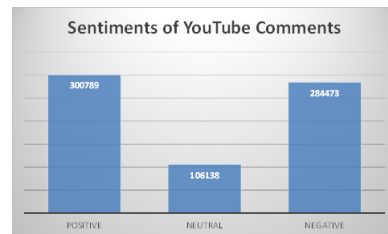


Figure 1: YouTube Dataset Statistics

likes and dislikes and user reviews. The US videos belong to 15 categories.

3.2 Label Generation

Our goal is to classify YouTube comments into positive, negative or neutral sentiment classes. However, there are no labels showing the sentiment of the comments. Considering the large number of comments, it was not possible for us to manually label a reasonable fraction of this dataset. Consequently, we used the TextBlob [2] library for Python to generate the labels. TextBlob applies NLP (Natural Language Processing) techniques to estimate the **polarity** of each comment independently from the rest of the comments, and only based on its content. For us, it replaces the burdensome task of manual labeling. TextBlob generates polarity scores in the range of $(-1, 1)$. We convert them into categorical data using Equation 1. Statistics on the size of each class is shown in Figure 1.

$$\text{sentiment} = \begin{cases} -1 & \text{polarity} < 0 \\ 0 & \text{polarity} = 0 \\ 1 & \text{polarity} > 0 \end{cases} \quad (1)$$

Words that happen more frequently in positive and negative comments are shown as word clouds in Figures 2 and 3.

In Section 5.3, we do category classification based on tags and comments. The category labels are already available in video data.

3.3 Feature Extraction

For sentiment classification based on comments, we used TF-IDF [1] to convert comments into feature vectors. In this technique, TF (Term Frequency) counts the number of times a word has occurred in each document. However, the number of occurrences of words is in general higher for longer documents. Then, it is helpful to divide the number of word occurrences by the number of words in the document. The output of TF is a sparse matrix mapping documents to a normalized vector representing how many times each word has occurred in them. IDF (Inverse Document Frequency) accounts for the number of occurrences of a word in the entire document corpus. Namely, some words are more likely than others to happen in a given corpus of documents. Then, we scale TF terms by a decreasing function of the number of occurrences of each word in the whole corpus which is indeed the IDF term.

For category classification, we took two approaches. First, we used the same TF-IDF comment features. Second, we used the tags from the video data. Since video data has been gathered daily over two weeks, it may contain each video multiple times along with

| Model | Train Accuracy | Test Accuracy |
|---------------------|----------------|---------------|
| Logistic Regression | 0.709 | 0.709 |
| Bernoulli NB | 0.709 | 0.715 |
| Ridge Classifier | 0.914 | 0.905 |
| Linear SVC | 0.956 | 0.951 |

Table 1: Sentiment classification accuracy of YouTube comments

independent from one another. On this dataset, SVM classifier has the best performance.

Figure 4 shows the calibration curve of the aforementioned techniques. It is observed that Naive Bayes classifier is not well-calibrated over the range $[0, 1]$. From among the remaining three methods, Logistic Regression is well-calibrated over the range $[0, 1]$ and is in general better calibrated than the other two methods.

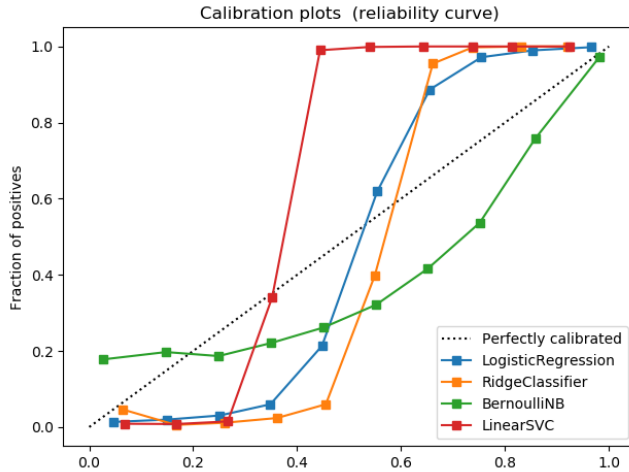


Figure 4: Word cloud for positive comments

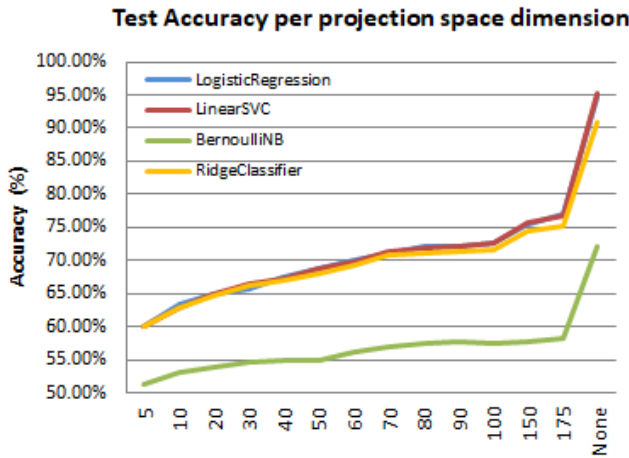


Figure 5: Accuracy vs. project space dimensions

5.2 Effect of Dimensionality Reduction

In this set of experiments, we examined how dimensionality reduction can affect the accuracy of predictions. Dimension reduction is usually applied to reduce the complexity of computations. We apply TruncatedSVD to project TF-IDF features to fewer dimensions

| Model | Train Accuracy | Test Accuracy |
|---------------------|----------------|---------------|
| Logistic Regression | 0.709 | 0.709 |
| Bernoulli NB | 0.709 | 0.715 |
| Ridge Classifier | 0.914 | 0.905 |
| Linear SVC | 0.956 | 0.951 |

Table 2: Sentiment classification accuracy of YouTube comments

| Model | Train Accuracy | Test Accuracy |
|---------------------|----------------|---------------|
| Logistic Regression | 0.709 | 0.709 |
| Bernoulli NB | 0.709 | 0.715 |
| Ridge Classifier | 0.914 | 0.905 |
| Linear SVC | 0.956 | 0.951 |

Table 3: Sentiment classification accuracy of YouTube comments

than those of the original TF-IDF feature space. We use TruncatedSVD rather than Principal Component Analysis (PCA) because the TF-IDF features are stored in a sparse matrix by the SKLearn library which PCA cannot process. TruncatedSVD is very similar to PCA. However, the data is not centered around its mean in the former. Figure 5 shows the accuracy of sentiment analysis with TruncatedSVD applied to TF-IDF features versus number of dimensions of the projection space. It is observed that accuracy increases with the number of dimensions. Nevertheless, with 175 dimensions, it is still far from the case where no dimensionality reduction is applied. It seems that in this dataset, 175 dimensions is not enough to project the data without losing much information. We could not experiment with higher dimensions due to memory limitations of our desktops.

5.3 Category Classification

First, we report the results of experiments for category classification using TF-IDF comment features. In this method we try to detect the category of the video a comment is talking about using the TF-IDF features of the comments. Table ?? shows the accuracy of different techniques. It is seen that the results are not as good as sentiment classification. Considering there are 15 categories versus 3 sentiment classes, this problem is deemed to be harder. Yet, the probability of a successful random category guess is only $1/15$ (6.7%). Consequently, category prediction based on TF-IDF features of comments is still a big improvement over random guessing. We then conducted category classification based on video tags. Table ?? shows accuracy of category classification using tags. It is observed that tags can predict the category with much higher accuracy. Note that tags are considered a property of videos while TF-IDF features are a property of comments. In addition, tags are explicitly chosen to classify the videos. As a result, it is expected that are more directly relevant to the content of the video and are less noisy which lead to better accuracy in category prediction.

In another experiment, we examined the effect of the number of categories on the accuracy. Intuitively, the larger the number of categories, the more difficult the classification problem is expected to be. To experiment this intuition, we created a list of categories with ascending order of each category population. We considered the first 2 categories, that is, the two categories with greatest number of videos, ignored the remaining videos in train and test set and measured the accuracy. The same experiment was repeated for 3, 4, ..., 15 categories. Figure 6 shows how test accuracy decreases by increasing the number of categories.

6 CONCLUSION

REFERENCES

- [1] A. Rajaraman and J. D. Ullman. *Mining of Massive Datasets*. Cambridge University Press, 2011.

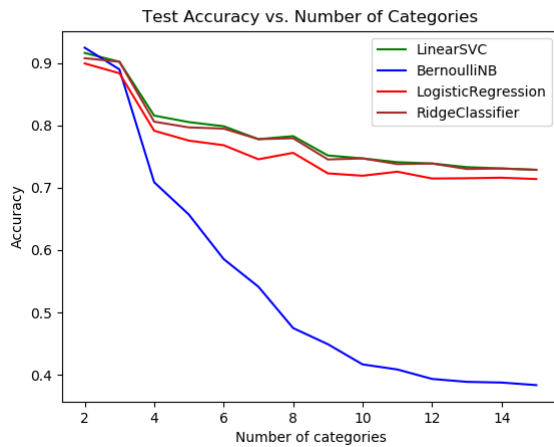


Figure 6: Accuracy versus projection space dimensions

- [2] YouTube. TextBlob Library For Python. <https://textblob.readthedocs.io/en/dev>.
- [3] YouTube. YouTube Database For Kaggle Competitions. <https://www.kaggle.com/datasnaek/youtube>.

7 CONTRIBUTIONS

All of us contributed in the discussions about what problem to target, what dataset to use and what techniques to apply. All of us helped with writing and reviewing the report. Peng searched different potential dataset candidates and summarized their advantages and shortcomings. He generated statistics about the YouTube dataset. Abhishek did sentiment polarity analysis and generated word clouds using TextBlob. Afshin implemented feature extraction and conducted classification experiments using scikit.