

# 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,

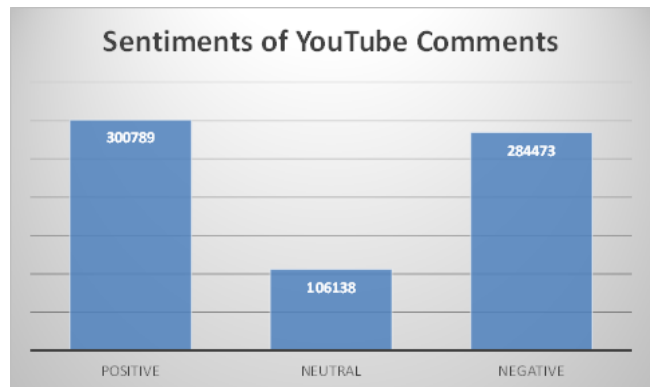


Figure 1: YouTube Dataset Statistics

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 2a and 2b.

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



possible methods: one-vs-one and one-vs-all. The one-vs-one method builds a classifier for every two classes. The class of a point is then the one chosen by the most classifiers. The one-vs-all approach builds a classifier for each class compared to all the remaining classes. The class of a data point is the one whose classifier achieves the greatest margin. In our experiment, we apply the one-vs-all method.

## 5 EXPERIMENTS AND EVALUATION

In this section, we present our experimental methodology, evaluation metrics and the results of the experiments.

### 5.1 Sentiment Analysis

In order to measure the performance of different classification techniques on the YouTube data, we split the data from US comments into train and test subsets. To generate the train set, we randomly choose 80% of the comments from each of the three categories. The remaining comments form the test set. We use 20% of the GB comments as development set. To compare the performance of different methods we use accuracy, precision and recall.

Precision-recall is useful, when doing binary classification, to better understand the output of the classifier. It is often used to determine the success rate of model prediction when the classes are imbalanced.

In our experiments, we plot compare, Accuracy, Precision and Recall to evaluate the output quality of classifiers. Since we are dealing with a multi-class classification problem, we first need to binarize the output and produce a curve for each class label. We also draw the precision-recall curve using the elements in the label indicator matrix for the binary prediction. This is also known as micro averaging. Once the precision and recall metrics are obtained, we compute the harmonic mean also known as F1 score. F1 score is defined as

$$F1 = 2PR/(P + R), \quad (6)$$

where P is the precision and R is the recall. It is used to measure the accuracy of the dev and test data. The maximum F1 score is 1 which means the model reached perfect precision and recall.

The curve displays the trade-off of the precision and recall at various thresholds. As the area of the curve increases, the precision and recall increase, where the increase of precision represents the increase of true positive rate and increase of recall represents the decrease of false negative rate.

Table 1 shows the accuracy measured on train and test sets. Naive Bayes classifier has the worst performance among all methods. This result is predictable as the conditional independence condition required by the Naive Bayes classifier is not satisfied in this dataset. In particular, the probability of two words happening in a comment given the category of the comment are not independent from one another. On this dataset, SVM classifier has the best performance. It is observed that the other three models have reasonable accuracy while LinearSVC outperforms the rest of the models. Figure 3a shows the Precision-Recall curve of Logistic Regression while Figure 3b shows that of LinearSVC. It is seen that precision and recall have an inverse relationship with each other. It happens because the more positive cases an algorithm reports, the more it is likely to report all positives points of the dataset (high recall), but it can also report more false positives (lower precision) and vice versa.

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.

Model	Train Accuracy	Test Accuracy
Logistic Regression	0.9578	0.9461
Bernoulli NB	0.7285	0.7212
Ridge Classifier	0.9419	0.9064
Linear SVC	0.9712	0.9527

Table 1: Sentiment classification accuracy of YouTube comments

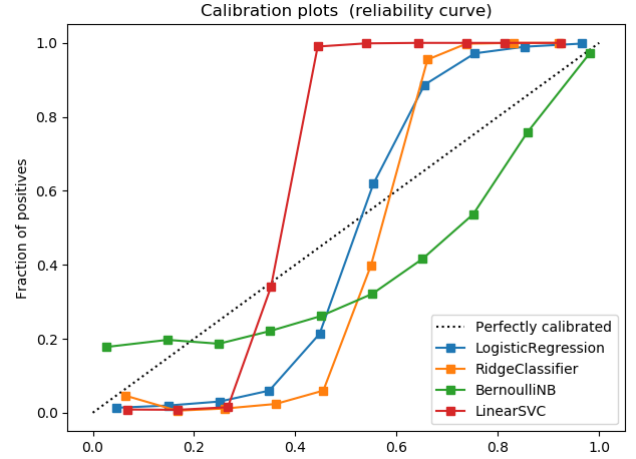


Figure 4: Calibration curves

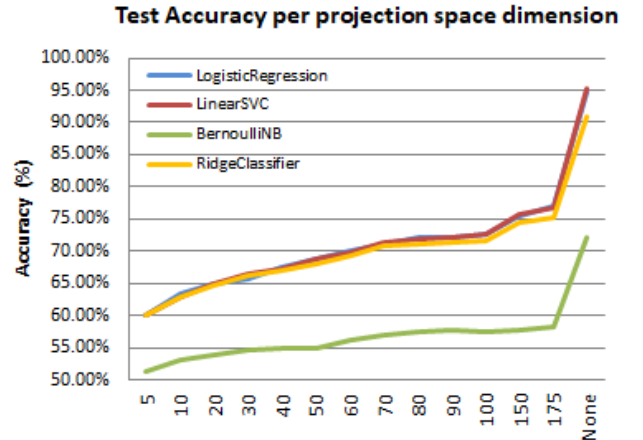
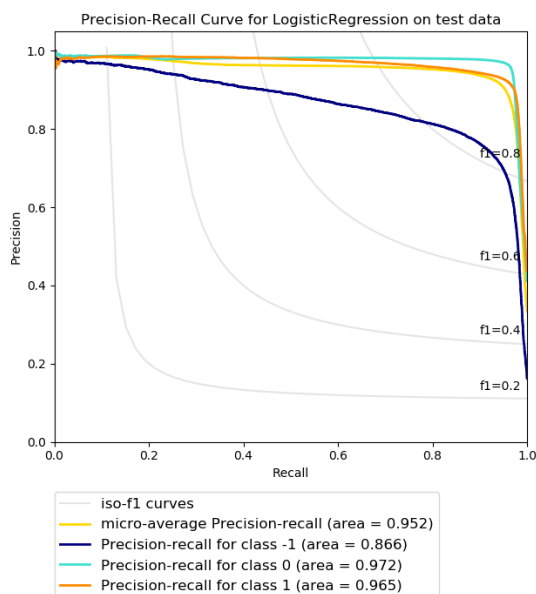
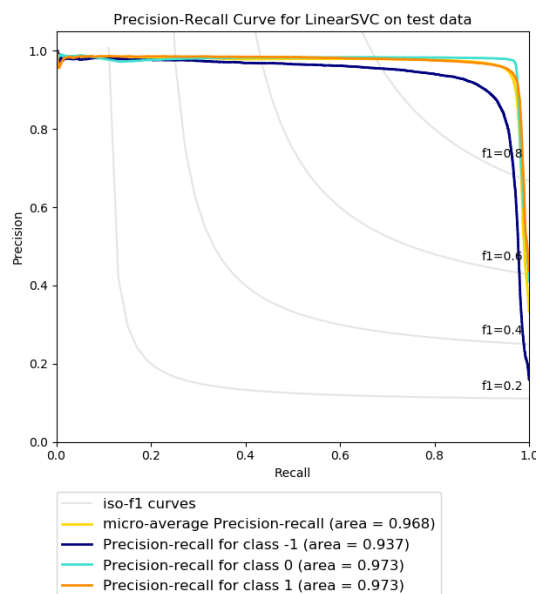


Figure 5: Accuracy vs. project space dimensions



(a) Precision vs. Recall for Logistic Regression



(b) Precision vs. Recall for Linear SVC

Figure 3: Precision vs Recall

## 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 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 6a 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 7 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.
- [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.

Model	Train Accuracy	Test Accuracy
Logistic Regression	0.6070	0.5322
Bernoulli NB	0.4554	0.4161
Ridge Classifier	0.6578	0.905
Linear SVC	0.6895	0.5780

(a) category classification accuracy using comments

Model	Train Accuracy	Test Accuracy
Logistic Regression	00.9987	0.7139
Bernoulli NB	0.4909	0.3835
Ridge Classifier	0.9955	0.7286
Linear SVC	0.9954	0.7286

(b) category classification accuracy using tags

Figure 6: Category prediction of YouTube videos

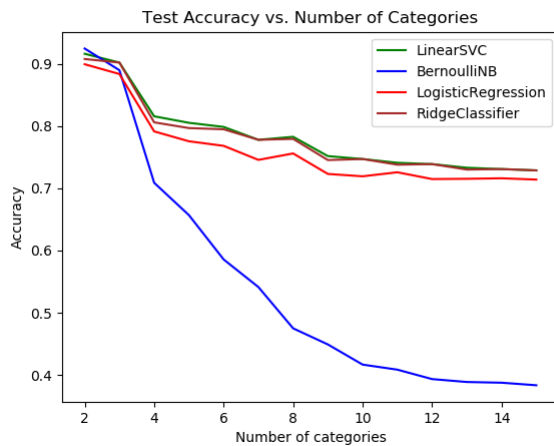


Figure 7: Accuracy versus projection space dimensions