

IMapBook Collaborative Discussion Classification

Uroš Polanc

Abstract

Automatic text classification has been considered as a vital method to manage and process a vast amount of documents in digital form, which is widespread and continuously increasing. That brought a wave of new tools for student education, such as the IMapBook platform. It gives students a way to communicate and discuss the books they are reading. In our assignment, we will try to classify those discussions based on the book it's relevant to, and the type of message.

Keywords

IMapBook, Text Classification, TF-IDF, Logistic Regression

Advisors : Slavko Žitnik

Introduction

With the rapid development of the internet and big data, digital information is increasing at a high rate. With this technologies and applications such as IMapBook [1], a web-based application that allows for discussions on reading materials as well as interactive games, also started emerging.

In our assignment, we will look at some IMapBook collaborative discussions on different books. For that, we first need to look at some vital information, that we extracted from the discussions and the book texts. We will create two different message classifications, based on :

- Book ID (3-class),
- Code Preliminary (16-class).

From the message types we can say that not all are important to the discussions of the books, and only the content discussion and content question are directly related to the books.

With this we were also provided with a (1) dataset that contains all message types and (2) dataset that contains only messages important to the discussion. With this we will be able to further compare how much noise impacts the classifiers and prediction.

By checking the manual classifications in the dataset, we noticed some were misspelled and some were incorrectly written. Therefore, we applied some fixes to the manual

classifications. The errors such as excess whitespace, lower and upper case, etc... were removed. Classifications where there were multiple classes, were simplified into one class, to simplify the classification process. The final fixed distribution is not intraclass uniform as can be seen in figures 3 and 2.

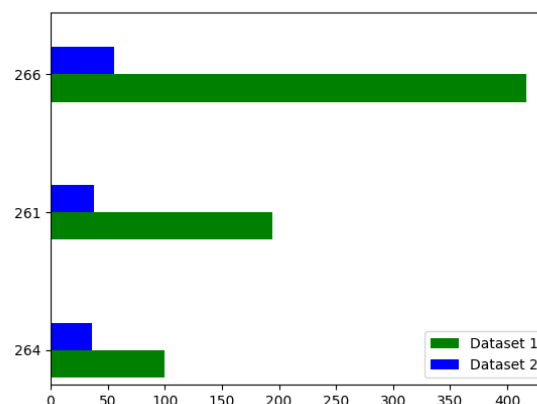


Figure 1. The distribution of book IDs in the dataset.

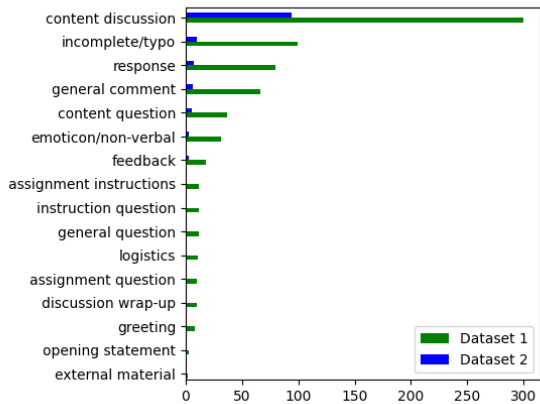


Figure 2. The distribution of message types in the dataset.

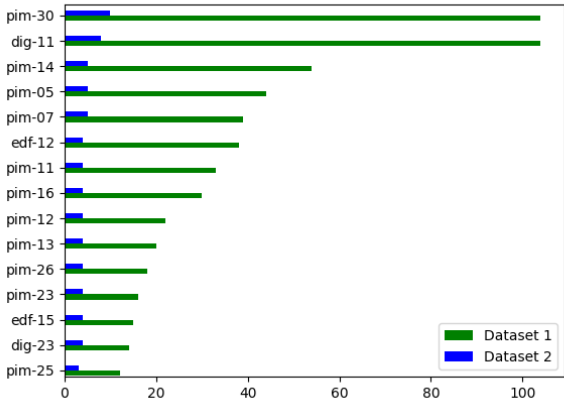


Figure 3. The distribution of active users in the dataset.

From all the graphs above we can see that in dataset 2, where there are messages concerning the discussion, we get a more uniform distribution of users as well as books. While the distribution of message types is dominated by the content discussion.

To better understand the contents we are trying to classify, we checked the top unigrams of the book texts. We also calculated the TF-IDF of the books as well as the collaborative discussions (as a whole) and displayed the top 15 stemmed words for each (see figure 4). This gives us a quick insight into the most important words of each book.



Figure 4. The top 15 words based on TF-IDF for each set.

From this we can say that the texts are contextually different, aside from a few words such as 'one', the top frequent words are more or less unique to their books.

Existing solutions

There are many approaches to automatic NLP text classification, which fall into three types of systems :

- rule-based systems,
- machine learning-based systems,
- hybrid systems.

Rule-based systems classify text into organized groups by using a set of handcrafted rules. But these systems require deep knowledge of the domain they are trying to classify. On the other hand, machine learning-based systems can learn the different associations between pieces of text by using pre-labeled examples as training data. Some of the more popular text classification algorithms [2] include Naive Bayes, Support Vector Machines, and deep learning. While hybrid systems apply both handcrafted rules and machine learning approaches. The current state-of-the-art consists of approaches such as BERT [3] or ELMo [4].

Results

We tested each of the categories for multiple classifiers using *tfidf* on both datasets. We achieved the following results, where RF is Random Forest, NB is Naive Bayes, LR is Logistic Regression, KNN is k-Nearest Neighbor and MV is Majority Voting.

Table 1. *CodePreliminary* with *tfidf* on dataset 1.

Model	Target	F1	Accuracy
RF	Train	0.53 ± 0.02	0.60 ± 0.02
RF	Test	0.39 ± 0.06	0.50 ± 0.05
NB	Train	0.47 ± 0.01	0.57 ± 0.01
NB	Test	0.34 ± 0.07	0.47 ± 0.06
LR	Train	0.62 ± 0.01	0.68 ± 0.01
LR	Test	0.34 ± 0.07	0.47 ± 0.06
KNN	Train	0.13 ± 0.01	0.12 ± 0.02
KNN	Test	0.07 ± 0.03	0.07 ± 0.03
MV	Train	0.56 ± 0.02	0.63 ± 0.01
MV	Test	0.38 ± 0.07	0.50 ± 0.06

Table 2. *Topic* with *tfidf* on dataset 1.

Model	Target	F1	Accuracy
RF	Train	0.75 ± 0.01	0.78 ± 0.01
RF	Test	0.66 ± 0.07	0.71 ± 0.06
NB	Train	0.83 ± 0.01	0.85 ± 0.00
NB	Test	0.70 ± 0.06	0.75 ± 0.05
LR	Train	0.86 ± 0.01	0.87 ± 0.00
LR	Test	0.70 ± 0.05	0.74 ± 0.04
KNN	Train	0.49 ± 0.04	0.61 ± 0.02
KNN	Test	0.45 ± 0.05	0.59 ± 0.04
MV	Train	0.84 ± 0.01	0.85 ± 0.00
MV	Test	0.70 ± 0.06	0.75 ± 0.04

Table 3. *BookID* with *tfidf* on dataset 1.

Model	Target	F1	Accuracy
RF	Train	0.75 ± 0.01	0.78 ± 0.01
RF	Test	0.66 ± 0.07	0.71 ± 0.06
NB	Train	0.83 ± 0.01	0.85 ± 0.00
NB	Test	0.70 ± 0.06	0.75 ± 0.05
LR	Train	0.86 ± 0.01	0.87 ± 0.00
LR	Test	0.70 ± 0.05	0.74 ± 0.04
KNN	Train	0.52 ± 0.05	0.63 ± 0.03
KNN	Test	0.46 ± 0.04	0.59 ± 0.04
MV	Train	0.84 ± 0.01	0.85 ± 0.00
MV	Test	0.70 ± 0.06	0.75 ± 0.04

Now on the second dataset, which contains only comments related to the books.

Table 4. *CodePreliminary* with *tfidf* on dataset 2.

Model	Target	F1	Accuracy
RF	Train	0.61 ± 0.01	0.72 ± 0.01
RF	Test	0.61 ± 0.11	0.72 ± 0.09
NB	Train	0.61 ± 0.01	0.72 ± 0.01
NB	Test	0.61 ± 0.11	0.72 ± 0.09
LR	Train	0.65 ± 0.02	0.75 ± 0.01
LR	Test	0.61 ± 0.11	0.72 ± 0.09
KNN	Train	0.67 ± 0.01	0.76 ± 0.01
KNN	Test	0.67 ± 0.10	0.77 ± 0.08
MV	Train	0.61 ± 0.01	0.72 ± 0.01
MV	Test	0.61 ± 0.11	0.72 ± 0.09

Table 5. *Topic* with *tfidf* on dataset 2.

Model	Target	F1	Accuracy
RF	Train	0.84 ± 0.02	0.84 ± 0.02
RF	Test	0.71 ± 0.14	0.72 ± 0.14
NB	Train	0.95 ± 0.01	0.95 ± 0.01
NB	Test	0.77 ± 0.14	0.78 ± 0.13
LR	Train	0.96 ± 0.01	0.96 ± 0.01
LR	Test	0.80 ± 0.15	0.81 ± 0.13
KNN	Train	0.74 ± 0.07	0.74 ± 0.07
KNN	Test	0.69 ± 0.18	0.70 ± 0.17
MV	Train	0.95 ± 0.01	0.95 ± 0.01
MV	Test	0.77 ± 0.14	0.78 ± 0.12

Table 6. *BookID* with *tfidf* on dataset 2.

Model	Target	F1	Accuracy
RF	Train	0.84 ± 0.02	0.84 ± 0.02
RF	Test	0.71 ± 0.14	0.72 ± 0.14
NB	Train	0.95 ± 0.01	0.95 ± 0.01
NB	Test	0.77 ± 0.14	0.78 ± 0.13
LR	Train	0.96 ± 0.01	0.96 ± 0.01
LR	Test	0.80 ± 0.15	0.81 ± 0.13
KNN	Train	0.73 ± 0.07	0.73 ± 0.07
KNN	Test	0.67 ± 0.18	0.68 ± 0.18
MV	Train	0.95 ± 0.01	0.95 ± 0.01
MV	Test	0.77 ± 0.14	0.78 ± 0.12

Conclusion

In our assignment we compared multiple classifiers in IMapBooks discussion classification. Overall the final results are not that great, as at its peak we got a $0.81\% \pm 0.13\%$ test accuracy on book prediction, where the concerning part is the high deviation. This could be further improved by tweaking the parameters maybe trying a different preprocessing approach.

The results also showed that most of the time Logistic Regression performed best on both training and testing sets. An important result was also that KNN performed significantly better on the dataset containing only messages pertaining to the discussion.

All this means that proper filtering is important when using text classification, as any unwanted or not relevant data will result in significant losses to the accuracy.

References

- [1] Grandon Gill and Glenn Smith. Imapbook: Engaging young readers with games. *Journal of Information Technology Education: Discussion Cases*, 2:10, 01 2013.
- [2] Tomas Pranckevicius and V. Marcinkevicius. Comparison of naive bayes, random forest, decision tree, support vector machines, and logistic regression classifiers for text reviews classification. *Balt. J. Mod. Comput.*, 5, 2017.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. 2019.
- [4] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations, 2018.