
Automatically Generate Tags for Questions Posted on StackExchange.com

Zhe Yu*

Department of Computer Science
North Carolina State University
Raleigh, NC 27606
zyu9@ncsu.edu

Abstract

In this project, the problem of how to automatically generate tags for questions posted on StackExchange websites is discussed. There are certain needs for this task since not all the users can tag their questions accurately due to the lack of understanding of the question or the website. An evaluation of $Fscore_M$ is utilized to meet the specific requirement of the task. After carefully studied the current multi-classification approaches in text mining, term frequency is selected for featurization, three different methods, SVM, Naive Bayes, and Decision Tree, are chosen to compare with each other in fulfilling the requirements. By analysing the results, imbalance is found to be the major problem for this task. (Not finished yet:)An over-sampling technique (SMOTE) is introduced to balance the data. The final result with SVM + SMOTE is satisfying.

1 Introduction

Now a days, more and more people post their questions and find their answers on StackExchange.com. With the help of tags, every user can find the questions they are most confident to solve. However, not all the users can tag their questions accurately due to the lack of understanding of the question or the website. This certain needs of automatically generate tags motivates the author to conduct this experiment. Without loss of generality, the experiment is conducted on one single site of StackExchange.com, which is <http://anime.stackexchange.com/>.

One significant characteristic of the data is imbalance. There are only 22 documents with the most minority tag "resources" while 1609 documents with the most majority tag "others". To evaluate the performance, we use the $Fscore_M$ [1] across the tags to equally reflect the performance across each tag. The $Fscore_\mu$ [1] is also calculated as a counter example since it only reflects the performance of majority tags. Our target is to achieve highest $Fscore_M$.

(Under construction:) In order to achieve a high $Fscore_M$, data needs to be balanced before the training of the classifier. Synthetic minority over-sampling technique (SMOTE) [2] is implemented to balance the data. An increase in $Fscore_M$ is reported while the $Fscore_\mu$ decrease.

1.1 Data

All the data used in this project is in a single file called anime.txt. It is collected from <http://anime.stackexchange.com/> and formatted by RAISE Lab. It contains 4827 documents with tags. For each document, the first tag is treated as its label. All material of this project can be found at <https://github.com/azhe825/AGT.git>. An example of document on anime.stackexchange.com is shown in Figure 1.

*Unit ID: zyu9, Student ID: 200109973

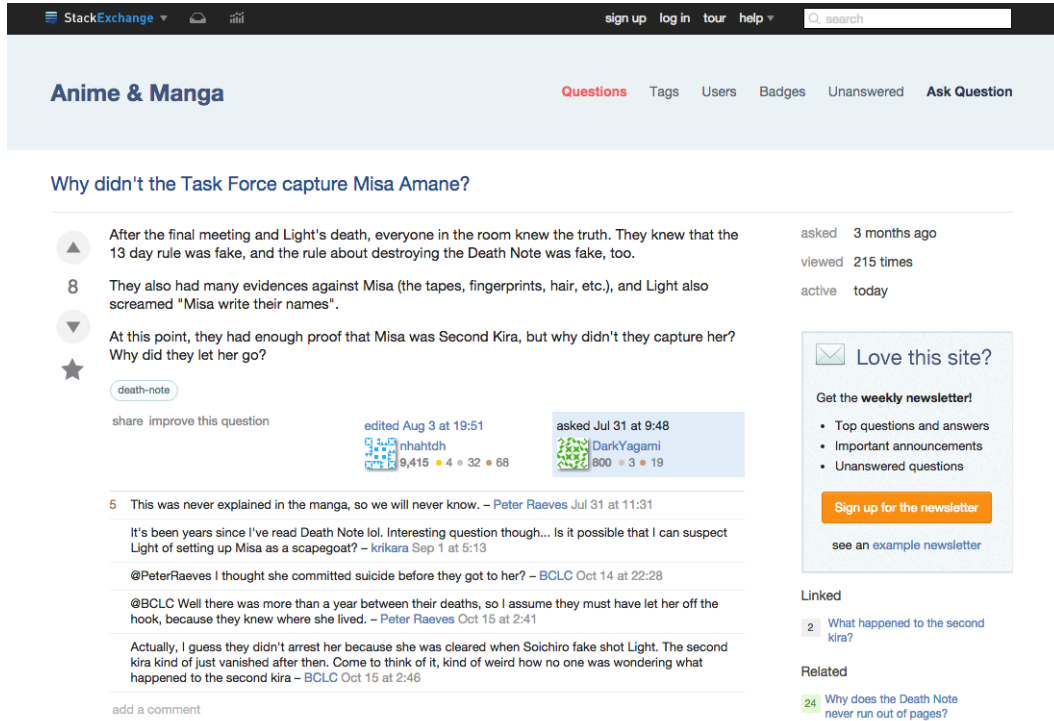


Figure 1: Example of post

The number of documents in each tag is: {'others': 1609, 'identification-request': 1099, 'naruto': 547, 'one-piece': 204, 'anime-production': 164, 'fullmetal-alchemist': 114, 'tropes': 110, 'bleach': 87, 'death-note': 86, 'fairy-tail': 78, 'dragon-ball': 66, 'code-geass': 51, 'japanese-language': 49, 'sword-art-online': 48, 'monogatari-series': 46, 'madoka-magica': 46, 'culture': 45, 'pokemon': 45, 'fate-stay-night': 43, 'shingeki-no-kyojin': 40, 'hunter-x-hunter': 38, 'anime-history': 35, 'manga-production': 33, 'from-the-new-world': 26, 'terminology': 25, 'neon-genesis-evangelion': 24, 'music': 24, 'toaru-majutsu-no-index': 23, 'resources': 22}.

1.2 Relative works

Clayton and Byrne, 2013 [3] have worked on StackOverflow tag prediction and developed an ACT-R inspired Bayesian probabilistic model. This approach achieves a 65% of accuracy by choosing the tag that has the highest log odds of being correct, given the tags prior log odds of occurrence and adjusting for the log likelihood ratio of the words in the post being associated with the tag.

Kuo, 2011 [4] has also worked on StackOverflow tag prediction. Kuo uses a co-occurrence model that predicts tags based on the relation (co-occurrence) between the words and tags. Initially built for next-word prediction in large documents, this model is adapted to the StackOverflow dataset by constraining the next word predicted to only tags. A 47% classification accuracy is achieved.

Another model called SNIF-ACT (Fu & Pirolli, 2007 [5]) also uses co-occurrences and can predict link traversal for search queries. The model predicts the most likely link that a person will click on by a search query (goal state) and fetched results.

2 Method

As suggested by [6] and [7], term frequency is calculated as feature and term frequency inverse document frequency (tf-idf) is used to select the most effective features. Inspired by CSC 522, three of the most famous classification methods, SVM, Naive Bayes, and Decision Tree, are used to conduct the experiments.

For SVM: linear kernel with primal problem solving is chosen since state-of-art accuracy on text classification is reported in [8].

For Naive Bayes: multinomial model is chosen according to the performance report in [9].

For Decision Tree: default settings are applied.

Cross validation: 5 by 5 cross validation, each time the classifier is trained on 80% of the data and tested on the rest 20%. For each classifier with each number of features, there are 25 results. The median and iqr of the 25 results represent the performance. A low iqr means that we can trust the result.

$Fscore_M$: F-score on each tag can be calculated after the trained classifier is tested on test set. The mean of F-score on each tag is then calculated to represent the overall performance of the classifier. Regardless of the population within each tag, the $Fscore_M$ represents performance on each tag equally [1].

$Fscore_\mu$: The $Fscore_\mu$ on each tag is calculated by multiply F-score on each tag by its population. Classifiers are trained to achieve high in this score by default, which is not desired in this task and therefore over-sampling methods like SMOTE is introduced in the later experiment [1].

(Upcoming:) SMOTE is implemented to over-sample documents in each tag to reach the number of majority tag. That is, SMOTE is implemented on each training set and after SMOTE, each tag will have the same number of documents.

3 Experiments

In the first experiment, SMOTE is not implemented. The performances of SVM, Naive Bayes, and Decision Tree are compared against each other by their $Fscore_M$.

(Upcoming:) In the second experiment, SMOTE is implemented to over-sample documents in each tag to reach the number of majority tag. The performances of SVM, Naive Bayes, and Decision Tree, with or without SMOTE are compared against each other by their $Fscore_M$.

3.1 Without SMOTE

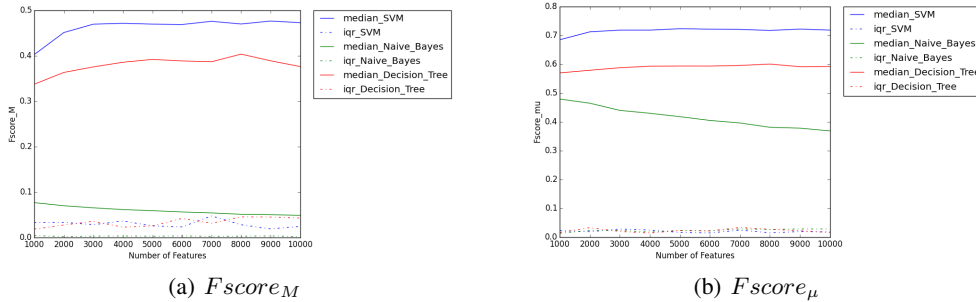


Figure 2: F-scores without SMOTE

SVM outperforms Naive Bayes and Decision Tree in both $Fscore_M$ and $Fscore_\mu$, according to Figure 2. However, even though SVM can achieve a 0.7 $Fscore_\mu$, its $Fscore_M$ is below 0.5, which is definitely not good enough.

As we can see from Figure 3, most of the tags are actually ignored (low median or high iqr) by all the classifiers due to their small population. This implies that a) $Fscore_\mu$ cannot correctly represent the true performance; b) over-sampling is essential and promising in achieving a better $Fscore_M$.

3.2 With SMOTE

Under construction.

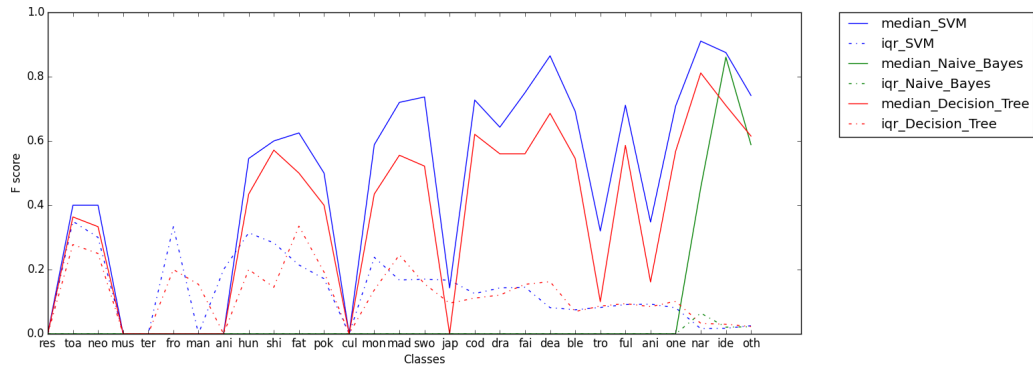


Figure 3: F-scores of each tag without SMOTE, sorted by population

3.3 Comparison between SMOTE and No-SMOTE

Under construction.

4 Conclusions and Future Works

Linear SVM outperforms Naive Bayes and Decision Tree in this task. (Under construction:)With SMOTE, linear SVM can achieve an $Fscore_M$ of ??? with feature number ???. This is the method we suggest according to this project. In our next step, different kernels of SVM can be tested to see if further improvement is available. There are other promising classifiers, artificial neural network with deep learning for example, as well. The result of this project can be implemented to automatically generate tag for posts. At least we can provide users suggestions of tag options if the result is not reliable enough.

References

- [1] Sokolova, M., & Lapalme, G. (2009). A systematic analysis of performance measures for classification tasks. *Information Processing & Management*, vol. 45(4), pp. 427-437.
- [2] Chawla, Nitesh V., Kevin W. Bowyer, Lawrence O. Hall & W. Philip Kegelmeyer. (2002) SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, pp. 321-357.
- [3] Clayton, S., & Byrne, M.. (2013) Predicting tags for stackoverflow posts. *In Proceedings of ICCM*, vol. 2013.
- [4] Kuo, D. (2011). On word prediction methods. *Technical report*, EECS Department, University of California, Berkeley.
- [5] Fu, W. T., & Pirolli, P. (2007). SNIF-ACT: A cognitive model of user navigation on the World Wide Web. *HumanComputer Interaction*, 22(4), 355-412.
- [6] Larsen, B., & Aone, C. (1999, August). Fast and effective text mining using linear-time document clustering. *In Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining* pp. 16-22. ACM.
- [7] Menzies, T. (2006). Improving IV&V Techniques Through the Analysis of Project Anomalies: Bayes networks-preliminary report. Tech. rep., West Virginia University. (<http://menzies.us/pdf/06anomalies-bayes0.pdf>).
- [8] Joachims, T. (2006, August). Training linear SVMs in linear time. *In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining* pp. 217-226. ACM.
- [9] McCallum, A., & Nigam, K. (1998, July). A comparison of event models for Naive Bayes text classification. *In AAAI-98 workshop on learning for text categorization* Vol. 752, pp. 41-48.

A Plan of Activities

Old plan:

- a) preprocess the data, generate feature matrix with term frequency or tf-idf;
- b) try different classification methods (decision tree, Naive Bayes, Adaboost, SVM, ...) to predict labels, cross validation will be applied;
- c) evaluate results from b) by comparing F_1 scores and confusion matrix;
- d) use SMOTE to improve the result.

Revised plan:

- a) preprocessing: generate feature matrix with term frequency and select features with top tf-idf scores;
- b) Decision Tree, Naive Bayes, and SVM are implemented to predict tags, cross validation is applied;
- c) $Fscore_M$ is justified for a better representation of the performance in this task. Both $Fscore_\mu$ and $Fscore_M$ are shown to compare the three classifiers;
- d) (upcoming:) use SMOTE to improve the result.

Member activities:

all the works are done by Zhe Yu.

Resources:

all the materials of this project can be found at <https://github.com/azhe825/AGT.git>.