

THE INTERNATIONAL JOURNAL OF SCIENCE & TECHNOLEDGE

Mining Students Social Data to Understand Their Learning Experience

Janani Sankari.M

Assistant Professor, Department of Information Technology, Jeppiaar Engineering
College, Chennai, India

Shoba L.K

Assistant Professor, Department of Information Technology, Jeppiaar Engineering
College, Chennai, India

Abstract:

Students informal conversation on social network (Eg. Facebook, Twitter) can be used to know their educational experiences such as opinions, concerns and learning process. These data from such sites will help us to know about students learning and their difficulties. The complexity is that we need human interpretation in analysis of the data. In this paper, we developed a workflow to integrate both qualitative analysis and large-scale data mining techniques. We focused on engineering students' Twitter posts to understand issues and problems in their educational experiences. We found engineering students encounter problems such as heavy study load, lack of social engagement, and sleep deprivation. Based on these results, we implemented a multi-label classification algorithm to classify tweets reflecting students' problems. We then used the algorithm to train a detector of student problems from about 35,000 tweets. This work, for the first time, presents a methodology and results that show how informal social media data can provide insights into students' experiences.

Keywords: Education, computers and education, social networking, web text analysis

1. Introduction

Social media sites such as Twitter, Facebook, and YouTube provide great venues for students to share joy and struggle, vent emotion and stress, and seek social support. On various social media sites, students discuss and share their everyday encounters in an informal and casual manner. Students' digital footprints provide vast amount of implicit knowledge and a whole new perspective for educational researchers and practitioners to understand students' experiences outside the controlled classroom environment. This understanding can inform institutional decision-making on interventions for at-risk students, improvement of education quality, and thus enhance student recruitment, retention, and success. The abundance of social media data provides opportunities to understand students' experiences, but also raises methodological difficulties in making sense of social media data for educational purposes. Just imagine the sheer data volumes, the diversity of Internet slangs, the unpredictability of locations, and timing of students posting on the web, as well as the complexity of students' experiences. Pure manual analysis cannot deal with the ever-growing scale of data, while pure automatic algorithms usually cannot capture in-depth meaning within the data. Traditionally, educational researchers have been using methods such as surveys, interviews, focus groups, classroom activities to collect data related to students' learning experiences. These methods are usually very time-consuming, thus cannot be duplicated or repeated with high frequency. The scale of such studies is also usually limited. In addition, when prompted about their experiences, students need to reflect on what they were thinking and doing sometime in the past, which may have become obscured over time. The emerging field of learning analytics and educational data mining has focused on analyzing structured data obtained from course management systems (CMS), classroom technology usage, or controlled online learning environments to inform educational decision-making. However, to the best of our knowledge, there is no research found to directly mine and analyze student-posted content from uncontrolled spaces on the social web with the clear goal of understanding students' learning experiences. The research goals of this study are 1) to demonstrate a workflow of social media data sense-making for educational purposes, integrating both qualitative analysis and large-scale data mining techniques as illustrated in Fig. 1; and 2) to explore engineering students' informal conversations on Twitter, in order to understand issues and problems students encounter in their learning experiences.

We chose to focus on engineering students' posts on Twitter about problems in their educational experiences mainly because:

- Engineering schools and departments have long been struggling with student recruitment and retention issues. Engineering graduates constitute a significant part of the nation's future workforce and have a direct impact on the nation's economic growth and global competency. The workflow we developed for making sense of social media data integrates qualitative analysis and data mining algorithms. The width of gray arrows represents data volumes – wider indicates more data volume. Black arrows represent data analysis, computation, and results flow. The dashed arrows represent the parts that do not concern the central work of this paper. This workflow can be an iterative cycle.

- Based on understanding of issues and problems in students' life, policymakers and educators can make more informed decisions on proper interventions and services that can help students overcome barriers in learning.
- Twitter is a popular social media site. Its content is mostly public and very concise (no more than 140 characters per tweet). Twitter provides free APIs that can be used to stream data. Therefore, we chose to start from analyzing students' posts on Twitter. In this paper, we went through an exploratory process to locate the relevant data and relevant Twitter hashtags (a Twitter hashtag is a word beginning with a # sign, used to emphasize or tag a topic). We collected 25,284 tweets using the hashtag #engineeringProblems over a period of 14 months, and a second dataset of 39,095 tweets using the geo-code (longitude and latitude) of Purdue University, West Lafayette. This corresponds to step 1. Three researchers conducted an inductive content analysis on samples of the #engineeringProblems dataset, which corresponds to steps 2 and 3 in Fig. 1.
- We found that major problems engineering students encounter in their learning experiences fall into several prominent categories. Based on these categories, we implemented a multi-label Naïve Bayes classification algorithm. We evaluated the performance of the classifier by comparing it with other state-of-the-art multi-label classifiers.
- Second, the paper provides deep insights into engineering students' educational experiences as reflected in informal, uncontrolled environments. Many issues and problems such as study-life balance, lack of sleep, lack of social engagement, and lack of diversity clearly emerge. These could bring awareness to educational pedagogy, policy-making, and educational practice. The remainder of this paper is organized as follows: the next section reviews theory of public discourse online, related work on text classification techniques used for analyzing tweets, and data-driven approaches in education. Section 3 describes the data collection process (step 1 in Fig. 1). Section 4 details the inductive content analysis procedures and categories identified (steps 2, 3, and 4). Section 5 details the implementation of the Naïve Bayes multi-label classifier and the evaluation results (step 5). In section 6, we show the comparison results of Naïve Bayes classifier with the popular classifier—Support Vector Machines (SVM) and one of its variations Max-Margin Multi-Label (M3L) classifier. This is an additional evaluation of the classifier in step 5. In section 7 we apply the Naïve Bayes classifier to the Purdue dataset in order to demonstrate its application in detecting students' problems at a specific university (steps 6 and 7). Section 8 discusses the limitations and possible future work, and section 9 concludes this study.

2. Related Work

2.1. Public Discourse on the Web

The theoretical foundation for the value of informal data on the web can be drawn from Goffman's theory of social performance. Although developed to explain face-to-face interactions, Goffman's theory of social performance is widely used to explain mediated interactions on the web today. One of the most fundamental aspects of this theory is the notion of front-stage and back-stage of people's social performances. Compared with the frontstage, the relaxing atmosphere of back-stage usually encourages more spontaneous actions. Whether a social setting is front-stage or back-stage is a relative matter. For students, compared with formal classroom settings, social media is a relative informal and relaxing back-stage. When students post content on social media sites, they usually post what they think and feel at that moment. In this sense, the data collected from online conversation may be more authentic and unfiltered than responses to formal research prompts. These conversations act as a zeitgeist for students' experiences. Many studies show that social media users may purposefully manage their online identity to "look better" than in real life. Other studies show that there is a lack of awareness about managing online identity among college students, and that young people usually regard social media as their personal space to hang out with peers outside the sight of parents and teachers. Students' online conversations reveal aspects of their experiences that are not easily seen in formal classroom settings, thus are usually not documented in educational literature. The abundance of social media data provides opportunities but also presents methodological difficulties for analyzing large-scale informal textual data. The next section reviews popular methods used for analyzing Twitter data.

2.2. Mining Twitter Data

Researchers from diverse fields have analyzed Twitter content to generate specific knowledge for their respective subject domains. For example, Gaffney analyzes tweets with hashtag #iranElection using histograms, user networks, and frequencies of top keywords to quantify online activism. Similar studies have been conducted in other fields including healthcare, marketing, athletics, just to name a few. Analysis methods used in these studies usually include qualitative content analysis, linguistic analysis, network analysis, and some simplistic methods such as word clouds and histograms. In our study, we built a classification model based on inductive content analysis. This model was then applied and validated on a brand new dataset. Therefore, we emphasize not only the insights gained from one dataset, but also the application of the classification algorithm to other datasets for detecting student problems. The human effort is thus augmented with large-scale data analysis. Below we briefly review studies on Twitter from the fields of data mining, machine learning, and natural language processing. These studies usually have more emphasis on statistical models and algorithms. They cover a wide range of topics including information propagation and diffusion, popularity prediction, event detection, topic discovery and tweet classification, to name a few. Amongst these topics, tweet classification is most relevant to our study. Popular classification algorithms include Naïve Bayes, Decision Tree, Logistic Regression, Maximum Entropy, Boosting, and Support Vector Machines (SVM). Based on the number of classes involved in the classification algorithms, there are binary classification and multi-class classification approaches. In binary classification, there are only two classes, while multi-class classification involves more than two classes. Both binary classification and multi-class classification are single-label classification systems. Single-label classification means each data point can only fall into one class where all classes are mutually exclusive. Multi-label classification, however, allows each data point to fall into several classes at the same time. Most

existing studies found on tweet classification are either binary classification on relevant and irrelevant content, or multi-class classification on generic classes such as news, events, opinions, deals, and private messages. Sentiment analysis is another very popular three-class classification on positive, negative, or neutral emotions/opinions. Sentiment analysis is very useful for mining customer opinions on products or companies through their reviews or online posts. It finds wide adoption in marketing and customer relationship management (CRM). Many methods have been developed to mine sentiment from texts. For example, both Davidov et al. and Bhayani et al. use emoticons as indicators to provide noisy labels to the tweets thus to minimize human effort. However, in the case of this paper, only knowing the sentiment of student-posted tweets does not provide much actionable knowledge on relevant interventions and services for students. Our purpose is to achieve deeper and finer understanding of students' experiences especially their learning-related issues and problems. To determine what student problems a tweet indicates is a more complicated task than to determine the sentiment of a tweet even for a human judge. Therefore, our study requires a qualitative analysis, and is impossible to do in a fully unsupervised way. Sentiment analysis is, therefore, not applicable to our study. In our study, we implemented a multi-label classification model where we allowed one tweet to fall into multiple categories at the same time. Our classification was also at a finer granularity compared with other generic classifications. Our work extends the scope of data-driven approaches in education such as learning analytics and educational data mining.

3. Data Collection

It is challenging to collect social media data related to students' experiences because of the irregularity and diversity of the language used. We searched data using an educational account on a commercial social media monitoring tool named Radian6. The Twitter APIs can also be configured to accomplish this task, which we later used to obtain the second dataset. The search process was exploratory. We started by searching based on different Boolean combinations of possible keywords such as engineer, students, campus, class, homework, professor, and lab. We then expanded and refined the keyword set and the combining Boolean logic iteratively. The Boolean search logic grew very complicated eventually, but the dataset still contained about 35% noise (during the month of November 2011, we retrieved 179 tweets, in which 63 were irrelevant to college students). Also, given that the dataset was so small, we seemed to have ruled out many other relevant tweets together with the spam and irrelevant tweets.

3.1. Development of Categories

The lens we used in conducting the inductive content analysis was to identify what are the major worries, concerns, and issues that engineering students encounter in their study and life. Researcher A read a random sample of 2,000 tweets from the 19,799 unique #engineeringProblem tweets, and developed 13 initial categories including: curriculum problems, heavy study load, study difficulties, imbalanced life, future and career worries, lack of gender diversity, sleep problems, stress, lack of motivation, physical health problems, nerdy culture, identity crisis, and others. These were developed to identify as many issues as possible, without accounting for their relative significances. Researcher A wrote detailed descriptions and gave examples for each category and sent the codebook and the 2,000-tweet sample to researchers Band C for review. Then, the three researchers discussed and collapsed the initial categories into five prominent themes, because they were themes with relatively large number of tweets.

3.2. Heavy Study Load

Our analyses show that, classes, homework, exams, and labs dominate the students' life. Libraries, labs, and the engineering building are their most frequently visited places. Some illustrative tweets are "Study over 30 hours for a test", "so much homework, so little time", and "C++CAE project due Tuesday, Mfg project Wednesday, 25 PageTech Report Wednesday + heavy homework load. Huzzah", and "homework never stops". Students express a very stressful experience in engineering. Not being able to manage the heavy study load leads to consequences such as lack of social engagement, lack of sleep, stress, depression, and some health problems.

3.3. Sleep Problems

Our analyses find that sleep problems are widely common among engineering students. Students frequently suffer from lack of sleep and nightmares due to heavy study load and stress. For example, "Napping in the common room because I know I won't sleep for the next three days", "If I don't schedule in sleep time, it doesn't happen", and "I wake up from a nightmare where I didn't finish my physics lab on time". Chronic lack of sleep or low-quality sleep can result in many psychological and physical health problems, therefore this issue needs to be addressed.

4. Naïve Bayes Multi-Label Classifier

We built a multi-label classifier to classify tweets based on the categories developed in the previous content analysis stage. There are several popular classifiers widely used in data mining and machine learning domain. We found Naïve Bayes classifier to be very effective on our dataset compared with other state-of-the-art multi-label classifiers.

4.1. Text Pre-processing

Twitter users use some special symbols to convey certain meaning. For example, # is used to indicate a hashtag, @ is used to indicate a user account, and RT is used to indicate a re-tweet. Twitter users sometimes repeat letters in words so that to emphasize the words, for example, "huungryyy", "sooo muuchh", and "Monndayyy". Besides, common stopwords such as "a, an, and, of, he, she, it", non-letter symbols, and punctuation also bring noise to the text. So we pre-processed the texts before training the classifier:

- We removed all the #engineeringProblemshashtags. For other co-occurring hashtags, we only removed the # sign, and kept the hashtag texts.
- Negative words are useful for detecting negative emotion and issues. So we substituted words ending with “n’t” and other common negative words (e.g. nothing, never, none, cannot) as “negtoken”.
- We removed all words that contain non-lettersymbols and punctuation. This included the removal of @ and http links. We also removed all theRTs.
- For repeating letters in words, our strategy was that when we detected two identical letters repeating, we kept both of them. If we detected more than two identical letters repeating, we replaced them with one letter. Therefore, “huuungryyy” and “sooo” were corrected to “hungry” and “so”. “muuchh” was kept as “muuchh”. Originally correct words such as “too” and “sleep” were kept as they were.
- We used the Lemur information retrieval toolkit to remove the common stopwords. We kept words like “much, more, all, always, still, only”, because the tweets frequently use these words to express extent. The Krovetz stemmer in the Lemur toolkit was used to perform stemming in order to unify different forms of a word, such as plurals and different forms of a verb.

4.2. Evaluation Measures for Multi-label Classifier

Commonly used measures to evaluate the performance of classification models include accuracy, precision, recall, and the harmonic average between precision and recall – the F1 score. For multi-label classification, the situation is slightly more complicated, because each document gets assigned multiple labels. Among these labels, some may be correct, and others may be incorrect. Therefore, there are usually two types of evaluation measures – example-based measures and label-based measures. Example-based measures are calculated on each document (e.g. each tweet is a document, and also called an example here) and then averaged over all documents in the dataset, whereas label-based measures are calculated based on each label (category) and then averaged over all labels (categories).

5. Comparison Experiment: Svm and M3l

SVM (Support Vector Machines) is one of the most used and accurate classifiers in many machine learning tasks, but our comparison experiment shows that Naïve Bayes exceeds SVM in this study. We first implemented a linear multi-label SVM using the LibSVM library with the one-versus-all heuristic. We applied weight of loss parameters that are proportional to the inverse of the percentages of tweets in or not in each category to account for the imbalanced categories. However, with the same training and testing datasets as in the above section, this one-versus-all SVM multi-label classifier classified all tweets into not in the category for all categories. So we got empty label sets for all tweets. Then we applied the same training and testing datasets as above to an advanced SVM variation named Max-Margin Multi-Label (M3L) classifier. M3L is a state-of-the-art multi-label classifier. Different from the one-versus-all heuristic, which assumes label independence, this classifier takes label correlation into consideration. We used the executable file of this algorithm provided by the authors. The performance is better than the simplistic one-versus-all SVM classifier, but still not as good as the Naïve Bayes classifier. Table 4 and Fig. 3 show the evaluation measures using M3L. Because SVM is not a probabilistic model, so Table 4 does not have probability threshold values as Table 2 does.

6. Discussion, Limitations, and Future Work

This study explores the previously uninstrumented space on Twitter in order to understand engineering students' experiences, integrating both qualitative methods and large-scale data mining techniques. In our study, through a qualitative content analysis, we found that engineering students are largely struggling with the heavy study load, and are not able to manage it successfully. Heavy study load leads to many consequences including lack of social engagement, sleep problems, and other psychological and physical health problems. Many students feel engineering is boring and hard, which leads to lack of motivation to study and negative emotions. Diversity issues also reveal culture conflicts and culture stereotypes existing among engineering students. Building on top of the qualitative insights, we implemented and evaluated a multi-label classifier to detect engineering student problems from Purdue University. This detector can be applied as a monitoring mechanism to identify at-risk students at a specific university in the long run without repeating the manual work frequently. Our work is only the first step towards revealing actionable insights from student-generated content on social media in order to improve education quality. The classifier is designed to be a multi-label classifier in order to reconcile this effect. If a tweet expresses correlation between “heavy study load” and “sleep problems”, then it can be categorized into both categories. After all, any mathematical and statistical models are a simplification of real world problems to a certain extent. The comparison experiment with M3L shows that this advanced model that accounts for label correlation does not perform as well as the simple Naïve Bayes model. Future work could specifically address the correlations among these student problems. Finally, the workflow we proposed requires human effort for data analysis and interpretation. This is necessary because our purpose is to achieve deeper understanding of the student experiences. To the best of our knowledge, there is currently no unsupervised automatic natural language processing technique that can achieve the depth of understanding that we were able to achieve. There is a trade-off between the amount of human effort and the depth of the understanding. The labels generated can be applied to any similar datasets in other institutions to detect engineering student problems without extra human effort. Often times, manual analysis is time-consuming not only because of the time spent on analysing the actual data, but also the time spent on cleaning, organizing the data, and adapting the format to fit the algorithms. We plan to build a tool based on the workflow proposed here combining social media data and possibly student academic performance data. This tool can assist in identification of students

at risk. This tool will provide a friendly user interface and integration between qualitative analysis and the classification and detection algorithms. Therefore, educators and researchers using this tool can focus on the actual data analysis and investigate the types of learning issues that they perceive as critical to their institutions and students. This tool can also facilitate collaboration among researchers and educators on data analysis. Advanced natural language processing techniques can be applied in the future to provide topic recommendations and further augment the human analysis results, but cannot completely rule out the human effort. Other possible future work could analyze students' generated content other than texts (e.g. images and videos), on social media sites other than Twitter (e.g. Facebook, Tumblr, and YouTube). Future work can also extend to students in other majors and other institutions.

7. Conclusion

Our study is beneficial to researchers in learning analytics, educational data mining, and learning technologies. It provides a workflow for analyzing social media data for educational purposes that overcomes the major limitations of both manual qualitative analysis and large scale computational analysis of user-generated textual content. Our study can inform educational administrators, practitioners and other relevant decision makers to gain further understanding of engineering students' college experiences. As an initial attempt to instrument the uncontrolled social media space, we propose many possible directions for future work for researchers who are interested in this area. We hope to see a proliferation of work in this area in the near future. We advocate that great attention needs to be paid to protect students' privacy when trying to provide good education and services to them.

8. References

1. G. Siemens and P. Long, "Penetrating the fog: Analytics in learning and education," *Educause Review*, vol. 46, no. 5, pp. 30–32, 2011.
2. M. Rost, L. Barkhuus, H. Cramer, and B. Brown, "Representation and communication: challenges in interpreting large social media datasets," in *Proceedings of the 2013 conference on Computer supported cooperative work*, 2013, pp. 357–362.
3. M. Clark, S. Sheppard, C. Atman, L. Fleming, R. Miller, R. Stevens, R. Streveler, and K. Smith, "Academic pathways study: Processes and realities," in *Proceedings of the American Society for Engineering Education Annual Conference and Exposition*, 2008.
4. C. J. Atman, S. D. Sheppard, J. Turns, R. S. Adams, L. Fleming, R. Stevens, R. A. Streveler, K. Smith, R. Miller, L. Leifer, K. Ya1939-1382 (c) 2013 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information. suhara, and D. Lund, "Enabling engineering student success: The final report for the Center for the Advancement of Engineering Education," Morgan & Claypool Publishers, Center for the Advancement of Engineering Education, 2010.