



## UNDERSTANDING STUDENTS' LEARNING EXPERIENCE BY DATA MINING OF SOCIAL MEDIA

<sup>1</sup>Ms. Pranjali S. Jadhav, <sup>2</sup>Dr. Shirish S. Sane,  
KKWIEER, Nasik

### Abstract

The computer and Education aims to increase knowledge and understanding of way in which digital technology can enhance education. As we know, The Social Media in 21th century is very popular the Twitter and Facebook are widely used. They shared their educational experience such as their ideas, opinions, their feelings and new things about learning process. Such uninstrumented environments will offer is effective for the data of student learning .This type knowledge| of knowledge of information} is analyze and such variety of data analyzing is difficult. The students' quality expertise is mirrored in social media sites which needs human interpretation. During this paper, we have a tendency to implement 2 things analysis and huge scale data processing. We have a tendency to provide example as engineering students' Twitter posts to grasp drawback occurred in their academic expertise. initial we have a tendency to get analysis and given sample taken from concerning twenty five,000 tweets associated with engineering students school life like significant study load, the social engagements etc. looking on the results, we have a tendency to apply one rule multi-label classification to classify Tweets of Students' issues.

**Index Terms**—Education, computers and education, social networking, web text analysis

### I. INTRODUCTION

The process of collection of data, searching and analyzing a huge amount of data in a database, so as to find patterns or relationships is the application of data mining to detect fraud.

Technically data mining is the process to find correlations or patterns among dozens of fields in large relational databases.

Automated prediction of drifts and behaviors: Mining helps to automate the process of searching predictive information in a large database. Questions that normally require extensive hands-on analysis can now be easily answered from the data. A typical example of a prognostic problem is targeted marketing. Our purpose is to achieve in depth and exquisite understanding of students' experiences especially in their learning-related issues and difficulties. To determine in a tweet, what are the student's problem, is a more complicated task than to determine the sentiment of a tweet even for a human judge. Therefore, this study requires a qualitative analysis, and is impossible to do in a fully unendorsed way. Sentiment analysis is, therefore, not applicable to this study. In this study, we will implement a multi-label classification model where we allows one tweet to fall into numerous categories at the same time. Our work extends the scope of data-driven approaches in education such as learning analytics and educational data mining. Usually, educational researchers are 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, so cannot be cloned or repeated with high frequency.

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 drawbacks are 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. Problems, and other psychological and physical health problems. Our work is only the first step towards revealing actionable insights from student-generated content on social media in order to improve education quality Heavy study load leads to many consequences including lack of social engagement, sleep problems, and other psychological and physical health problems. Our work is only the first step towards revealing actionable insights from student-generated content on social media in order to improve education quality. We extend the proposed algorithm which analysis the student's learning experiences by giving solutions to their problems. The suggested solution is forwarded to the student's individual email-ids to attain the privacy of student and for improving security a novel secure algorithm called BIRCH is proposed. Finally we get the feedback from the students about solution provided and comparison graph is generated.

### RELATED WORKS

The research goals of this study are

- 1) To demonstrate a work-flow of social media information sense-making for academic functions, group action each chemical analysis and large-scale data processing techniques and
- 2) To explore engineering students' informal conversations on Twitter, so as to grasp problems and issues students encounter in their learning experiences. These studies have additional stress on applied mathematics models and algorithms. They cowl a large vary of topics quality prediction, event detection, topic discovery and tweet classification. Amongst these topics, tweet classification is most relevant to the present study. well-liked classification algorithms embody Naïve mathematician, call Tree, provision Regression, most Entropy, Boosting, and Support Vector Machines (SVM). Most existing studies found on tweet classification area unit either binary classification on relevant and inapplicable content, or multi-class classification on generic categories like news, events, opinions, deals, and personal messages.

Sentiment analysis is extremely helpful for mining client opinions on merchandise or corporations through their reviews or on-line posts. It finds wide adoption in selling and client relationship management (CRM). We selected to specialize in engineering students' posts on Twitter regarding issues in their academic experiences principally because:

1. Engineering faculties and departments have long been fighting student achievement and retention problems.

Engineering graduates represent a major part of the nation's future manpower and have an immediate impact on the nation's economic process and world competence.

2. Supported understanding of problems and issues in students' life, policymakers and educators will create additional hip to choices on correct interventions and services which will facilitate students overcome barriers in learning.
3. Twitter may be a standard social media web site. Its content is generally public and extremely epigrammatic (no quite one hundred forty characters per tweet). Twitter provides free arthropod genus which will be wont to stream knowledge. Therefore, we tend to selected to begin from analyzing students' posts on Twitter.

## II LITERATURE SURVEY

### A. eMUSE

It provides integrated access to all the Web 2.0 tools selected by the instructor for the course at hand: common access point, detailed usage instructions, summary of the latest activity. It retrieve students' actions with each tool and store them in a local database and offer a summary of each student's activity, including graphical visualization, evolution over time, comparisons with peers, as well as aggregated data eMUSE compute a score based on the recorded student activity (following instructor-defined criteria) and provide basic administrative services (authentication service, enroll students to the course, edit profile etc.).

The first step towards the creation of eMUSE was to pick out the foremost appropriate net 2.0 tools to be integrated into the system, that meet 2 requirements: i) have a demonstrated education worth (according to case studies rumored within the literature); ii) provide technical support for mashup integration (well documented and maintained genus Apis, RSS feeds etc.). the combination of the web 2.0 tools into the

platform was done by means that of mashups, guaranteeing a light-weight design, with loosely-coupled parts. A mashup represents a mix of knowledge and/or functionalities from 2 or additional external sources to make a replacement web application[8].

### B. MOOC

#### MOOC

Advanced instructional technologies square measure developing apace and on-line MOOC courses have become additional prevailing, making Associate in Nursing enthusiasm for the ostensibly limitless knowledge driven potentialities to have an effect on advances in learning and enhance the training expertise. For these potentialities to unfold, the experience and collaboration of the many specialists are necessary to boost knowledge assortment, to foster the event of higher prognosticative models, and to assure models square measure explainable and unjust. the large knowledge collected from MOOCs has to be larger, not in its height (number of students) however in its width—more meta-data and data on learners' psychological feature and self-regulatory states needs to be collected in addition to correctness and completion rates. This additional detailed articulation can facilitate open up the black box approach to machine learning models wherever prediction is the primary goal. Instead, a knowledge-driven learner model approach uses fine grain data that's planned and developed from psychological feature principles to make instructive models with sensible implications to boost student learning.

Using data-driven models to develop and improve instructional materials is basically completely different from the instructor-centered model. In data-driven modeling, course development and improvement is predicated on data-driven analysis of student difficulties and of the target expertise the course is supposed to produce; it's not based on teacher self-reflection as found in strictly instructor-centered models. To be sure, instructors will and will contribute to deciphering information and creating course design choices, however ought to ideally do thus with support of psychological science expertise. Course improvement in data-driven modeling is additionally supported course-embedded *in vivo* experiments (multiple educational styles haphazardly appointed to students in natural course use, additionally known as "A/B testing") that value the impact of other course styles on

sturdy learning outcomes. In courses supported scientific discipline and data-driven modeling, student interaction is less targeted on reading or being attentive to an instructor's delivery of data, however is primarily regarding students' learning by example, by doing and by explaining. additionally to avoiding the pitfall of developing interactive activities that don't offer enough helpful information to reveal student thinking, MOOC developers and knowledge miners should avoid potential pitfalls within the analysis and use of information. One such pitfall is the application of sophisticated statistical and machine learning techniques to educational data without understanding or contributing to relevant cognitive and pedagogical principles. This "black box model" approach focuses on improving prediction without regards to understanding what is happening cognitively (i.e., inside the the box). Using data-driven learner models to improve courses contrasts with the instructor-centered model in three key ways. First, course development and improvement is based not solely on instructor self-reflection, but on a data-driven analysis of student difficulties and of the target expertise the course is meant to produce. Second, course improvement is based on course-embedded *in vivo* experiments that evaluate the effect of alternative course designs on robust learning outcomes. Third, course interaction is not centrally about instructor's delivery knowledge, but about student learning by example, by doing and by explaining.

### B. C. Learning Analytics and Educational Data Mining

Learning Analytics and educational data mining are data driven approaches emerging in education. These approaches analyze data generated in educational settings to understand students and their learning environments in order to inform institutional decision-making. First data analyzed using these approach typically are structured data including administrative data, students' activity and performance data from CMS (Course Management System). In prediction, the goal is to develop a model which can infer a single aspect of the data (the predicted variable, similar to dependent variables in traditional statistical analysis) from some combination of other aspects of the data (predictor variables, similar to independent variables in traditional statistical analysis). Structure discovery algorithms attempt to find

structure in the data without an a priori idea of what should be found, a very different goal than in prediction. In prediction, there is a specific variable that the EDM/LA researcher attempts to model; by contrast, there is not a specific variable of interest in structure discovery. In relationship mining, the goal is to discover relationships between variables in a data set with a large number of variables. Broadly, there are four types of relationship mining: association rule mining, correlation mining, sequential pattern mining, and causal data mining [6]. EDM and LA methods have similarly been useful in understanding student learning in various collaborative settings. Collaborative learning behaviors have been analyzed in order to determine which behaviors are characteristic of more successful groups and more successful learners, in multiple contexts, including computer-mediated discussions online collaboration using software development tools, and interactive table top collaboration.

#### *D. 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 #Iran Election using histograms, user networks, and frequencies of top keywords to quantify online activism. Analysis methods used usually includes qualitative content analysis, linguistic analysis, network analysis, and some simplistic methods such as word clouds and histograms. The classification model based on inductive content analysis applied and validated on dataset. Therefore, it emphasize not only the insights gained from dataset, but also the application of the classification algorithm to other datasets for detecting students' problems. The human effort is thus augmented with large-scale data analysis. Popular classification algorithms include naïve bayes, decision tree, logistic regression, maximum entropy, and boosting and support vector machine. 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 multiclass classification involve more than two classes. Both binary classification and multiclass 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[9]. 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. Many methods have been developed to mine sentiment from texts. For example, both Davidov et al. and Bhayani et al. use emotions as indicators to provide noisy labels to the tweets thus to minimize human effort. However, only knowing the sentiment of student-posted tweets does not provide much actionable knowledge on relevant interventions and services for students. The 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 [7]. Therefore, it requires a qualitative analysis and is impossible to do in a fully unsupervised way. Multilabel classification, however, allows each data point to fall into several classes at the same time.

### **I. PROPOSED SYSTEM**

The analysis goals of this study square measure

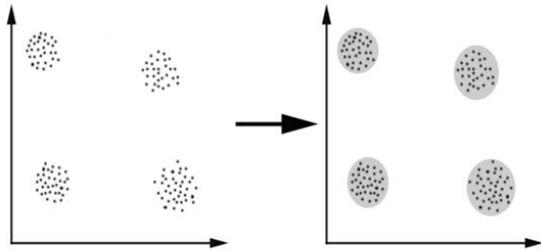
- 1) To demonstrate a advancement of social media knowledge sense-making for academic functions, integrating each analysis and large-scale data processing techniques as illustrated in Fig. 1; and
- 2) To explore engineering students' informal conversations on Twitter, so as to grasp problems and issues students encounter in their learning experiences.

We selected to target engineering students' posts on Twitter regarding issues in their educational experiences principally because:

1. Engineering faculties and departments have long been fighting student accomplishment and retention problems. Engineering graduates represent a big a part of the nation's future personnel and have a right away impact on the nation's economic process and world ability
2. supported understanding of problems and



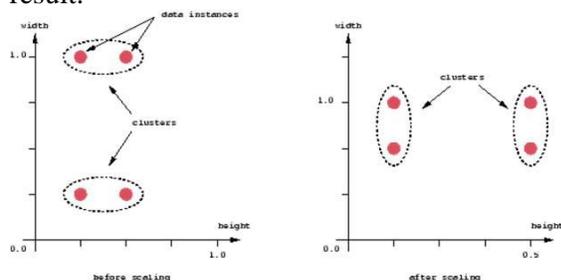
are kept in a cluster whereas dissimilar objects are in totally different teams. It is the most necessary unattended learning drawback. It deals with finding structure in a assortment of untagged information. For higher understanding please visit Fig 2.



**Fig 2:** showing four clusters formed from the set of unlabeled data

For clustering algorithm to be advantageous and beneficial some of the conditions need to be satisfied.

1) Scalability - Data must be scalable otherwise we may get the wrong result. Fig 3 shows simple graphical example where we may get the wrong result.



**Fig 3:** showing example where scalability may leads to wrong result

- 2) Clustering algorithm must be able to deal with different types of attributes.
- 3) Clustering algorithm must be able to find clustered data with the arbitrary shape.
- 4) Clustering algorithm must be insensitive to noise and outliers.
- 5) Interpret-ability and Usability - Result obtained must be interpretable and usable so that maximum knowledge about the input parameters can be obtained.
- 6) Clustering algorithm must be able to deal with data set of high dimensionality.

## V. IMPLEMENTATION

### 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. 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. From the limited number of relevant tweets we retrieved, we found a Twitter hashtag #engineering Problems occurring most frequently. Students used the hashtag #engineering Problems to post about their experiences of being engineering majors. This was the most popular hashtag specific to engineering students' college life based on the data retrieved using the Boolean terms. Using Radian6, we streamed tweets containing this hashtag for about 14 months (421 days) from November 1st, 2011 to December 25th, 2012. In total, we collected 25,284 tweets with the hashtag #engineering Problems posted by from 10,239 unique Twitter accounts. Counting re-tweets, replies, and mentions, a total of 12,434 unique user accounts were involved. After removing duplicates caused by re-tweeting, there were 19,799 unique tweets in this dataset. We also identified several other much less popular but relevant hashtags such as #lady Engineer, #engineering- Majors, #switching Majors, #college Problems, and #nerd status. As a side note for future work, these hashtags can also be used to retrieve data relevant to college students' experiences. To demonstrate the application of the classification algorithm, we obtained another new dataset using the geocode of Purdue University West Lafayette (40.428317, -86.914535) with a radius of 1.3 miles to cover the entire campus. From February 5th to April 17th, 2013, we obtained 39,095 tweets using the Twitter search API. These tweets came from 5,592 unique user accounts. There were 35,598 unique tweets after removing duplicates. One reason we chose

Purdue as an example is that it is a large public university with a strong engineering student base. Over 27% (10,533/39,000) of the students at the West Lafayette campus are enrolled in the College of Engineering during the 2012-2013 academic year. Nevertheless, the general approach we used can be applied to any institution and students in any major.

### **Inductive Content Analysis**

Because social media content like tweets contain a large amount of informal language, sarcasm, acronyms, and misspellings, meaning is often ambiguous and subject to human interpretation. Rostet. al argue that in large scale social media data analysis, faulty assumptions are likely to arise if automatic algorithms are used without taking a qualitative look at the data. We concur with this argument, as we found no appropriate unsupervised algorithms could reveal in-depth meanings in our data. For example, LDA (Latent Dirichlet Allocation) is a popular topic modeling algorithm that can detect general topics from very large scale data. LDA has only produced meaningless word groups from our data with a lot of overlapping words across different topics. There were no pre-defined categories of the data, so we needed to explore what students were saying in the tweets. Thus, we first conducted an inductive content analysis on the #engineeringProblems dataset. Inductive content analysis is one popular qualitative research method for manually analyzing text content. Three researchers collaborated on the content analysis process.

### **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 #engineeringProblems tweets, and developed 13 initial categories including: curriculum problems, heavy study load, study difficulties, imbalanced life, future and carrier 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 B and 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. We found that many tweets could belong to more than one category. For example, “*This could very well turn into an all-nighter...f\*\*\* you lab report #no sleep*” falls into heavy study load, lack of sleep, and negative emotion at the same time. “*Why am I not in business school?? Hate being in engineering school. Too much stuff. Way too complicated. No fun*” falls into heavy study load, and negative emotion at the same time. So one tweet can be labeled with multiple categories. This is a multi-label classification as opposed to a single-label classification where each tweet can only be labeled with one category. The categories one tweet belongs to are called this tweet’s labels or label set.

### **Classification Results**

From the inductive content analysis stage, we had a total of 2,785 #engineeringProblems tweets annotated with 6 categories. We used 70% of the 2,785 tweets for training (1,950 tweets), and 30% for testing (835 tweets). 85.5% (532/622) of words occurred more than once in the testing set were found in the training data set.

## **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, good education and services to them. In the future, which analysis the student’s learning experiences by giving solutions to their problems. The suggested solution is forwarded to the student’s individual email-ids to attain the Privacy of student and for improving security by a novel secure algorithm. Finally get the feedback from the students about solution provided to generate comparison graph.

## VI. FUTURE SCOPE

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. There are a number of limitations, which also lead to many possible directions for future work.

## ACKNOWLEDGMENT

The authors would like to thank the publishers and researchers for making their resources available. We also thank the college authority for providing the required infrastructure and support. Finally we would like to extend our heartfelt gratitude to friends and family members.

## VII. 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. Yasuhara, 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.
- [5] R. Ferguson, "The state of learning analytics in 2012: A review and future challenges," *Knowledge Media Institute, Technical Report KMI-2012-01*, 2012.
- [6] Xin Chen, Mihaela Vorvoreanu, and Krishna Madhavan, "Mining Social Media Data for Understanding Students' Learning Experiences", *IEEE Transaction on Learning Technologies*, 2014.
- [7] Andrei Yakushev and Sergey Mityagin, "Social networks mining For Analysis and Modeling Drugs Usage", Elsevier, 2014.
- [8] Elvira Popescu, "Providing collaborative learning support with Social Media in an integrated environment", *World Wide Web*, Springer, 2014
- [9] Joris Klerkx, Katrien Verbert and Erik Duval, "Enhancing Learning With Visualization techniques", In *Handbook of Research on Educational Communications and Technology*, Springer, 2014.
- [10] Payal S. Jain, Pallavi S. Panhale, Praneshwari A. Deokar "EFFECTIVE MINING SOCIAL MEDIA DATA FOR UNDERSTANDING STUDENTS LEARNING EXPERIENCES" *International Research Journal of Engineering and Technology (IRJET)* Jan-2016.