# Understanding Students' Learning Experiences through Mining User-Generated Contents on Social Media

Tran Thi Oanh[1,*], Nguyen Van Thanh[2]

[1]*VNU International School, Building G7-G8, 144 Xuan Thuy, Cau Giay, Hanoi, Vietnam*
[2]*E-learning Training Center, Hanoi Open University,*
*B101 Nguyen Hien, Hai Ba Trung Dist, Hanoi, Vietnam*

**Abstract:** This paper presents a work of mining informal social media data to provide insights into students' learning experiences. Analyzing such kind of data is a challenging task because of the data volume, the complexity and diversity of languages used in these social sites. In this study, we developed a framework which integrating both qualitative analysis and different data mining techniques in order to understand students' learning experiences. This is the first work focusing on mining Vietnamese forums for students in natural science fields to understand issues and problems in their education. The results indicated that these students usually encounter problems such as heavy study load, sleepy problem, negative emotion, English barriers, and carreers' targets. The experimental results are quite promising in classifying students' posts into predefined categories developed for academic purposes. It is expected to help educational managers get necessary information in a timely fashion and then make more informed decisions in supporting their students in studying.

*Keywords:* Students' learning experience, mining social media, students' forums, understand students' issues.

## 1. Introduction

Learning experience refers to how students feel in the process of getting knowledge or skill from studying in academic environments. It is considered to be one of the most relevant indicator of education quality in schools/universities [1]. Quality educational provision and learning environment can render most rewarding learning experiences. Student experience has thus become a central tenet of the quality assurance in higher education. Getting to understand this is an effective and important way to improve educational quality in schools/universities. This helps policy makers and academic managers can make more informed decisions, make more proper interventions and services to help students overcome their barriers in learning, provide a more valid range of activities to support enhancements to the student learning experience and provides guidance and resources for learning and teaching.

To identify students' learning experiences, the widespread used methods is to undertake a number of surveys, direct interviews or observations that provide important opportunities for educators to obtain student feedback and identify key areas for action.

_____
* Corresponding author. Tel.: 84-1662220684.
  Email: oanhtt@isvnu.vn

Unfortunately, these traditional methods are usually very time-consuming, thus cannot be duplicated or repeated with high frequency. Their scalability is also limited to a small number of participants. Moreover, they also raise the question of accuracy and validity of data collected because they do not accurately reflect on what students were thinking or doing something at the time the problems/issues happened. This is due to the time of taking survey is far from that experience, which may have become obscured over time. Another drawback is that the selection of the standards of educational practice and student behavior implied in the questions is also criticized in the surveys [2]. Therefore, in strategic approaches, institutions should also gather data from external data sources to develop intelligence on students' learning experiences.

Nowadays, social media provide great venues for students to share their thoughts about everything in their daily life. On these sites, they could discuss and share everything they may encounter in an informal and casual way. These public data sets provide vast amount of implicit knowledge for educators to understand students' experiences besides the above traditional methods. However, these data also raise methodological difficulties in making sense for educational purposes because of the data volumes, the diversity of slang languages used on the Internet, the different time and locations of students' posting as well as the complexity of students' experiences. To the best of our knowledge, so far in Vietnam, there is no study that directly mines and analyzes these student-generated contents on social webs towards the goals of understanding students' learning experiences.

In this paper, we present a research of using new technologies which allow for data mining and data scraping to extract and comprehend students' learning experiences through their digital footprints on social webs. To deal with the task, we illustrate a workflow of making sense of these social media data for educational purposes. More specifically, we chose to

focus on identifying issues or problems students encounter in their learning experiences. In summary, the main contributions of this paper are:

● Performing a qualitative method to analyze informal social data from students' digital footprints. Then, building a dataset for the purpose of understanding students' learning experiences.

● Developing a framework using data mining techniques to automatically detect students' issues and problems in their study at universities.

● Conducting experiments to prove the effectiveness of the proposed methods.

The rest of this paper is organized as follows: Section 2 presents related work. In Section 3, we describe how to collect raw data from social sites. Section 4 shows a qualitative analysis of the dataset to develop a set of categories that natural science students may encounter in their study. Section 5 describes a framework for mining social data in order to understand students' learning experiences. Section 6 shows experimental results and some findings of this work. Finally, we conclude the paper in Section 7 and discuss some future research directions.

## 2. Related work

Social media has risen to be not only a personal communication media, but also a media to communicate opinions about products and services or even political and general events among its users. Many researches from diverse fields have developed tools to formally represent, measure, model, and mine meaningful patterns (knowledge) from large-scale social for the concerned domains. For example, researchers investigate the task of sentiment analysis [3], which determine the attitude or polarity of opinions or reviews written by humans to rate products or services. In healthcare, many researches [4] has shown that social media services can be used to

disclose a range of personal health information, or to provide online social support for health issues [5]. In the marketing field, researchers mine the social data to recommend friends or items (e.g. movies, music, news, books, research articles, search queries, social tags, and products in general.) on social media sites. Recommender systems [6] typically produce a list of recommendations in one of two ways – through collaborative and content-based

filtering or the personality-based approach based on the information of a user's past behavior, similar decisions made by other users as well as a series of discrete characteristics of an item. Most existing studies recast the above tasks as a classification problem. The classification can be either binary classification on relevant and irrelevant content, or multi-class classification on generic classes.

In the educational field, Educational Data Mining is an emerging discipline, concerned with developing methods for exploring the unique and increasingly large-scale data that come from educational settings, and using those methods to better understand students, and the settings which they learn in. Most studies in this field focus on students' academic performance [7, 8] using the information when students interact with the tutoring/e-learning systems. In comprehending students' posts on social sites such as Twitter [9] firstly provide a workflow for analyzing social media data for educational purposes. This study is beneficial to researchers in learning analytics, EDM, and learning technologies. Among previous study, our work is closest to this one.

In our study, we also implemented a multi-class classification model where one post can fall into multiple categories at the same time. In building dataset, we focus on mining social media for Vietnamese education. We extend understanding Vietnamese students to include informal social media data based on their informal online conversations on the Web.

## 3. Collecting data from social media sites

### 3.1. Collecting raw data

Collecting data relating to students' experiences on the social site is not an easy task because of the diversity and irregularity of languages used. We wrote a Java program to automatically crawl student-generated posts on a blog of a university, and acquired lots of posts. In principal, we could collect raw data from any social media channel which allows students to post anything they wish to. In this paper, we chose to collect data from a forum of a famous university in Vietnam (  a great forum on the web for students to post anything about their study, their life and their concerns. It is quite simple to collect raw data of students' posts on this forum by a crawling program. However, the challenge is to filter out posts referring to studying topics because of irregularity and diversity of languages used. Among lots of collected raw data, we found that only 20% posts were relevant to the students' study issues (we randomly selected 300 posts, in which 242 posts were irrelevant).

To improve the quality of raw data, we investigated the topic tree in this forum and filtered out irrelevant posts which usually fall into sub-tree topics. Finally, we got ~7000 posts, after filtering, we obtained and manually labeled 1834 posts relating to students' learning experiences.

### 3.2. Pre-processing data

*Cleaning data:* The purpose of this process is to make data clean to prepare for extracting features of classification models. In more details, we performed several pre-processing techniques as follows:

- Removing and replacing teenagers' languages which are commonly used on social media posts such as: *ak, đc, dc, ntn, ntnao, nhìu, hok, e, wa, wa', j, j`, r, k, ko bây h, bj h, t gian, hjx, sv, t7*
- Removing hashtags such as *#nhàtrọ, #tựhàoBK, …*

- Removing all words containing special symbols or not alphabetic/numeric letters. These words usually are email addresses, URL addresses, etc.

*Word Segmentation***:** The entire data after cleaning was automatically segmented on the

level of the word. This is important techniques used in Natural Language Processing in many languages whose word boundary is not separated by white spaces. An example of a Vietnamese post after word-segmented is illustrated in Figure 1.

```
lại thêm một ngày buồn ở bk và có_lẽ sẽ còn nhiều ngày
buồn nữa sẽ đến buồn xong rồi lại phải cố_gắng thôi bk
phũ_phàng như_vậy áp_lực như_vậy nhưng rồi mọi chuyện
sẽ ổn thôi mình mà phải đầu_hàng à mình mà phải
chịu_thua à không có đâu đợi đấy kì sau sẽ phục_thù
```

Figure 1. An example of Vietnamese post after segmenting words
(morphemes are concatenated by hyphen).

*Removing Stop Words:* Stop words are basically a set of commonly used words in any language. These words appear to be of little value in helping select documents matching a user need, therefore, are excluded from the vocabulary entirely. In Vietnamese, some examples of stops words are "và", "hoặc", "mỗi", "cũng", etc. We based on a typical Vietnamese stop word list \footnote{The size of this list is …} which is commonly used for many task in NLP.

## 4. A qualitative analysis on the dataset

Previous research [9] have found that in English, automatic supervised algorithms could not reveal in-depth meanings in the social media sites. This situation is also true in our context, especially when we want to achieve deeper understanding of the students' experiences. In fact, we tried to apply Z-LDA algorithms [10], one of the most typical and robust topic modelling technique, to our dataset. Unfortunately, it has only produced meaningless word groups with lots of overlapping words across different topics. Hence, we have to set a set of categories relating students' learning experiences by performing inductive content analysis on the dataset.

In discovering these posts, we paid attention to identify what are major concerns, worries,

and issues that students encounter in their daily life and study. Firstly, two people independently investigate these posts and proposed totally 14 initial categories including: heavy study load, curriculum problems, negative emotion, credit problems, part-time jobs, studying abroad, career target, studying English, learning experiences, soft skills, choosing major fields, reference material, mental problems, and others. These two people then sit together to discuss and collapse the initial categories into seven prominent themes (as shown in Table 1). They together wrote the detailed description and gave examples for each category. Based on that, they independently labeled the dataset. Then, we measured the inter-rater agreement using Cohens' Kappa and got 0.82 F1. This rate is quite high, so the quality of the dataset is acceptable. For the posts which raters conflict on determining labels, we consulted a third person to fix their labels. After labeling, there was a total of 1834 labeled posts used for model training and testing. Table 1 gives a description of the number of instances per labels in our dataset.

Table 1. Number of posts in each category of the dataset analyzed

| No. | Labels | #instances |
|-----|--------|------------|
| 1 | Heavy Study Load | 444 |
| 2 | Negative Emotion | 141 |

| 3 | Career targets | 143 |
| 4 | English barriers | 228 |
| 5 | Material resources | 348 |
| 6 | Diversity issues | 236 |
| 7 | Others | 458 |

The description of each category is given below:

**Heavy Study Load**

Investigating students' posts let us know that classes, homework, exams, laboratories dominate students' life. Some examples include *"quá nhiều bài tập về trong một thời gian ngắn"*, *"kỳ thi sắp tới mà không nắm được chút nào kiến thức do quá khó hiểu"*, *"hắc_nghiệt quá bao năm nay mong_ước ra trường sắp được rồi còn nốt đồ_án thôi"*, *"quá_trình làm luận_văn tốt_nghiệp thật mệt_mỏi và ốm_đau tôi đã vượt qua nỗi sợ_hãi viết luận_văn tốt_nghiệp như_thế_nào"*, *"các bác ơi sao em học môn tín_hiệu và hệ_thống không hiểu gì cả làm_sao bây_giờ đây sắp thi giữa kì mà chưa được chữ gì vào đầu cả"*. In these posts, students express tiredness and stressful experiences in studying and taking examination in universities. This will lead to many bad consequences such as health problems, depression, and stress. Hence, students desire a more balanced life than their real academic environments.

**Negative Emotion**

These topics' posts are quite diverse, ranging from bad emotions of dormitories' life, homesick, disappointment, sickness, stressed with school works to bad friend relationships, student-teacher relationship, etc. Some examples include *"ừm thì chết một lúc một lúc bỗng_nhiên tim ngừng đập một lúc không phải suy_nghĩ một lúc không buồn một lúc không cảm_thấy chán_nản một lúc không cảm_thấy mình chới_với một lúc không cười một lúc không khóc một lúc không phải cô_đơn một lúc không phải ray_rứt một lúc ừm thì chỉ một lúc một lúc ngừng thở một lúc bình_yên ..."*, *"buồn vào hồn không tên thức_giấc nửa_đêm nhớ chuyện xưa vào đời đường_phố vắng đêm nao quen một người mà yêu_thương chót chao nhau chọn lời để rồi làm_sao quên biết tên người*

*quen biết nẻo đi đường về và có biết đêm nào ta hẹn_hò để tâm_tư nhưng đêm ngủ không yên ..."*. Therefore, it is very important if students could get necessary helps, emotional support for that particular situation.

**Career Targets**

Students want to choose a career that will make us happy, but how can we know what that will be? Choosing a career path (or changing one) is, for most of us, a confusing and anxiety-riddled experience. Many will tell you to "follow your passion" or "do what you love," but this is not very useful advice. Students always wonder about how their future would be. Some examples include *"em là sinh_viên khoa cơ_khí em đang rất phân_vân không biết nên chọn cơ_điện_tử hay cơ_khí động_lực cái việc chọn chuyên_ngành rất quan_trọng vì nó sẽ là sự_nghiệp sau_này của mình điều này ..."*, *"những công_việc mà sinh_viên ngành ta ra trường có_thể làm được đánh_giá về công_việc ví_dụ như thu_nhập ban_đầu thu_nhập về sau_này khả_năng thăng_tiến trong công_việc về lương_bổng về chức_tước về khả_năng chuyên_môn ..."*, *"chào các anh_chị em là sinh_viên đang học muốn đi theo ngành truyền_thông và mạng máy_tính nhưng em chưa biết rõ lắm về các công_việc sau_này sẽ làm ở ngành này mong các anh_chị biết về ngành giúp xin chân_thành cảm_ơn ..."*. Hence, if educational managers could catch these students' wonders, they could support their students in choosing the right careers that best fit students' personalities, as well as their preferences.

**English Barriers**

One of the main problems with Vietnamese students is language barriers, especially English. Students often feel lack of confidence in using English as a second languages to study. Some example posts include *"mấy tháng trước chuẩn_bị thi toeic tình_cờ đọc được một blog chia_sẻ kinh_nghiệm luyện nghe rất thiết_thực mình làm theo và cũng đã vượt để đủ điều_kiện ra trường chia_sẻ mọi người tham_khảo"*, *"tháng trước mình có bắt_đầu học tiếng anh theo phương_pháp effortless_english nhờ một*

*chị giới_thiệu cho ban_đầu học rất nản học được hai tháng thì bỏ khoảng hai tuần sau đó nghĩ sao lại quay lại học tiếp đến hiện_tại là khoảng gần sáu tháng rồi tuần trước mình có cơ_hội nói_chuyện với hai anh người tây làm bên cứu_trợ quốc_tế về nước_sạch...".* Understanding this point could aid managers make plans and strategies to help students overcome language barriers.

### Material resources

Students cannot receive a proper education without the right resources. Getting the suitable materials means having adequate funding, which many schools lack due to governmental budget cuts. This is an issue that is all too common among many schools in Vietnam but is continuously overlooked. Some typical example posts include *"các bác nào biết hà_nội chỗ nào bán sách dạy lập_trình phong_phú nhất không mình đang muốn kiếm tài_liệu về học mà không biết chỗ nào bán", "tổng_hợp các bộ source code đồ_án phần_mềm mức_độ khó cho anh_em tham_khảo các đồ_án được chọn_lọc một_cách kỹ_lưỡng sử_dụng các công_nghệ mới nhất thích_hợp cho anh_em làm đồ_án tốt_nghiệp", "có cao_nhân nào pro giúp_đỡ em với bài_tập lớn nhiệt động kỹ_thuật của thầy thư có tài_liệu giải bài_tập lớn của các khóa trước hoặc là ai làm được thì pm em theo địa_chỉ em cảm_ơn ạ", "có_pro nào có slide bài giảng môn đa_phương_tiện của thầy trần_nguyên_ngọc không cho mình xin với thầy khó_khăn trong việc gửi slide bài giảng quá nghe ở lớp là một chuyện nhưng muốn về nhà đọc lại cho kĩ mà không_thể có được slide của thầy khá hay và chi_tiết nên mình muốn đọc thật kĩ pro nào có thì chia_sẻ với nhé".* Therefore, universities need to know this in a timely fashion and then make plan to support students in accessing materials necessary for their study.

### Diversity Issues

There is also many posts referring to other issues such as studying abroad, lacking of soft skills, finding hostel, credit problem, etc. Some examples include *"mình đang cần liên_hệ với một bạn trong lớp này xin cho mình số đt hoặc*

*ym của bất_kỳ ai trong lớp này mình có việc rất quan_trọng nhờ giúp_đỡ xin cảm_ơn xin giúp mình với ...", "xăng tăng đột_biến vật_giá leo_thang tiết_kiệm quốc_sách một_số mẹo trong video này có_thể giúp xe bạn uống nhiên_liệu ít hơn tiết_kiệm được túi_tiền của bạn và gia_đình ...", "đúng là cuộc_sống ở nước_ngoài nhất_là ở các nước phát_triển là niềm mơ_ước của chúng_ta có_thể nói ai cũng có những nhận_xét như các bạn đã nêu nhất_là các quan_chức sau khi đi tham_quan đều cũng có những cảm_nhận như các bạn ..."*

### Others

Many posts do not have a clear meaning, or do not express the problems relating to students' learning experiences.

## 5. A Proposed method for understanding students' learning experiences using data mining techniques

Figure 2 shows the proposed framework for mining students' social data on the Web. The framework include the training phase and testing phase. In the first phase, we train a model of recognizing students' experiences automatically using data mining techniques. To train the classifying models, we utilized the dataset developed from Section 4. In the second phase, we use the trained model to classify a new post of students into predefined categories of students' issues.

To build the prediction model, we generate a multi-label classifier to classify posts based on a predefined category developed by investigating posts collected from a forum of a university. There are many common classifiers used in data mining such as SVM [11], Naïve Bayes [12], Decision Tree [13, 14], etc. These classifiers are powerful and proved to be effective in many other tasks of NLP [15]. Therefore, in experiments we also conducted a simple yet powerful machine learning method, namely Decision Tree, to estimate its performance on the task of understanding students' learning experiences.
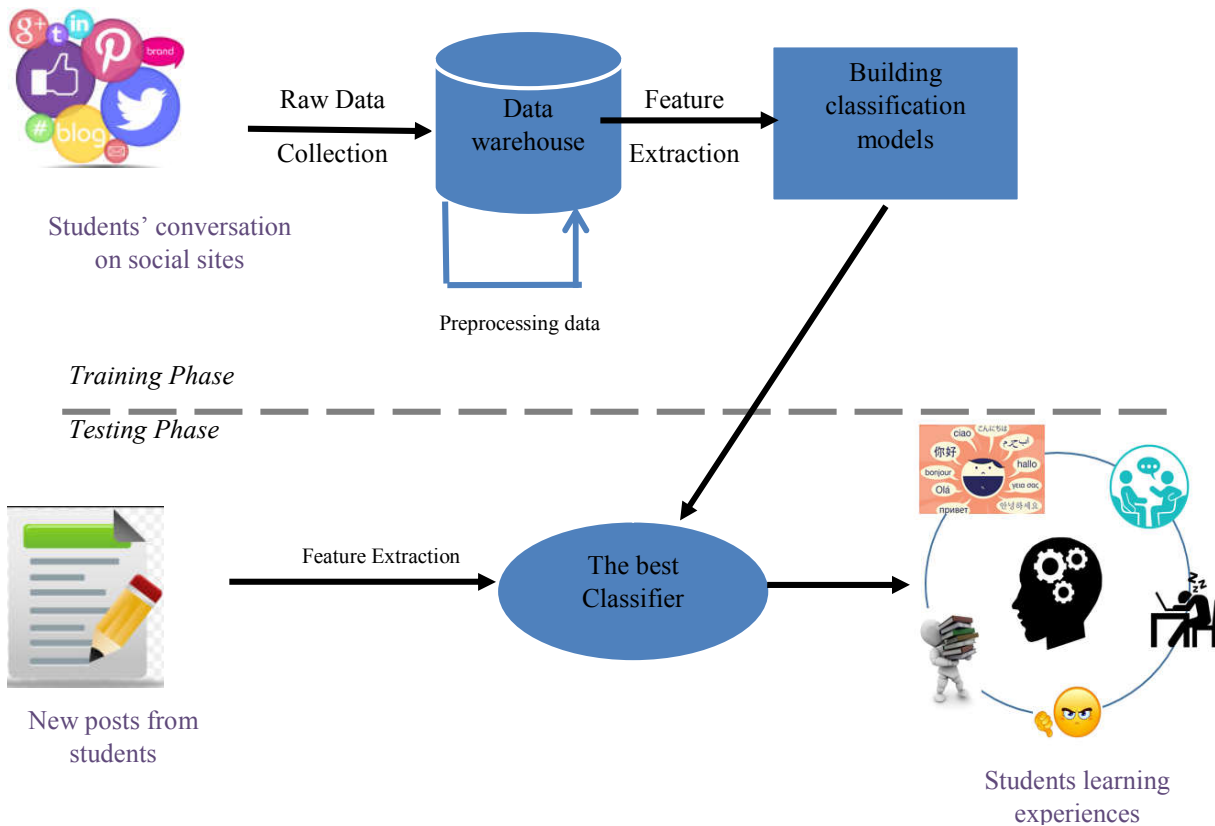
Figure 2 . A framework for mining social media data using data mining techniques.

As discussed above, this task can be recast as a multi-label classification problem, a variant of the classification problem where multiple target labels must be assigned to each post. Formally, multi-label learning can be phrased as the problem of finding a model that maps inputs $\mathbf{x}$ to binary vectors $\mathbf{y}$, rather than scalar outputs as in the ordinary classification problem. The task of learning from multi-label classification problem can be addressed by transformation techniques. This technique turns the problem into several single-label classification problems. There are two main methods of this techniques called "*binary relevance*" and "*label combination*".

● Binary relevance (BR): If there's q labels, the binary relevance method create q new data sets, one for each label and train single-label classifiers on each new data set. One classifier only answer yes/no to the question "does it belong to label i?". The final multi-label prediction for a new instance is determined by aggregating the classification results from all independent binary classifiers

● *Label combination (LC)*: BR is simple but does not work well when there's dependencies between the labels. This method tries to solve that drawback by taking into account label correlations. Each different combination of labels is considered to be a single label. After transformation, a single-label classifier $H: X \rightarrow \mathcal{P}(L)$ is trained on $\mathcal{P}(L)$ the power set of all labels. The main drawback of this approach is that the number of label combinations grows exponentially with the number of labels. This increases the run-time of classification.

## 6. Experiments

### 6.1. Evaluation metrics for multi-label classifiers

In the single-label classification, metrics such as accuracy, precision, recall, and the F1 score were commonly used to evaluate the performance. However, in the multi-label classification the evaluation metrics are more complicated because of some reasons: one post can be assigned more than one label; and some labels can be correct while some are incorrect. In this situation, researchers proposed two types of metrics which are example-based measures and label-based measures.

**Example-based measures**

These measures are calculated based on examples (in this case each post is considered as an example) and then averaged over all posts in the dataset.

Suppose that we are classifying a certain post *p*, the gold (true) set of labels that p falls into is G, and the predicted set of labeled by the classifier is P, the example-based evaluation metrics are calculated as follows:

$$Accuracy = \frac{1}{N} \sum_{i=1}^{N} \frac{G_i \cap P_i}{G_i \cup P_i}$$

$$Precision = \frac{1}{N} \sum_{i=1}^{N} \frac{G_i \cap P_i}{P_i} \quad \text{and}$$

$$Recall = \frac{1}{N} \sum_{i=1}^{N} \frac{G_i \cap P_i}{G_i}$$

$$F1 = \frac{1}{N} \sum_{i=1}^{N} \frac{2.Precision_i.Recall_i}{Precision_i + Recall_i}$$

where N is the number of posts in the dataset.

**Label-based measures**

These measures are calculated based on label and then averaged over all labels in the dataset. For each classifier for a label *l*, we create a matrix of contingency for that particular label *l*. Table 2 shows that matrix.

Table 2. Contingency Table per label. (note that the sum of tp, tn, fn, and fp equal to the number of posts).

| | | Gold Standard | |
| --- | --- | --- | --- |
| | | True l | True not l |
| Classification Outcome | Predicted as *l* | True postive (tp) | False positive (fp) |
| | Predicted as not *l* | False negative (fn) | True negative (tn) |

Based on that matrix, we calculate the measures as follows:

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn}$$

$$Precision = \frac{tp}{tp + fp} \quad \text{and} \quad Recall = \frac{tp}{tp + fn}$$

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$

There are two more commonly used measures to estimate the performance of multi-labeled classification which are micro-average F1 and macro-average F1. The former gives equal weight to each per-post classification decision, while the latter gives equal weight to each label. They are variants of F1 used in different situation. In the case there is no label whose probability is greater than a threshold T, we assign the post to the label with the largest probability.

### 6.2. Experimental setups

To train and test the model, we performed 10-fold cross validation test. In building and testing models, we exploited the following tools:

Classifiers: WEKA (http://www.cs.waikato.ac.nz/ml/weka/)

- Word segmenter: vnTokenizer (http://mim.hus.vnu.edu.vn/phuonglh/softwares/vnTokenizer)
- Stop-word list: containing about 200 common words

*6.3. Experimental results*

6.3.1. Estimating the effect of using different machine learning techniques

With 7 labels, we have $2^6$=64 possible label sets for each post. The thresholds in the Decision Tree classifier are determined by the one which yields the best performance on evaluation metrics. By experiments, we set the thresholds for J48 to 0.8.

Table 3 shows experimental results. From experiments, we can see that machine learning-based classifiers achieved significant improvement in comparison to the random guessing baseline, Zero Rule - a baseline classification uses a naive classification rule in both settings of multi-label classification, binary relevance and label combination.

|  | Accuracy | Recall | Precision | F1 micro | F1 macro |
|---|---|---|---|---|---|
| Binary Relevance |  |  |  |  |  |
| Zero Rule |  |  | Very low |  |  |
| J48 (threshold = 0.8) | 0.443 | 0.504 | 0.633 | 0.559 | 0.56 |
| Label Combination |  |  |  |  |  |
| Zero Rule | 0.251 | 0.143 | 0.036 | 0.24 | 0.058 |
| J48 | 0.565 | 0.548 | 0.571 | 0.583 | 0.558 |

6.3.2. Performance of classifying each category

Table 4 shows experimental results measuring label-based accuracy and F1 score for each category using Decision Tree. These results are quite promising in detecting students' learning experiences from online posts. This suggests that it is appropriate to use the best classifiers to apply for detecting students' learning experiences when having new posts from students.

Table 3. Label-based accuracy and F1 scores for each category using Decision Tree

|  | Heavy Study Load | Negative Emotion | Career targets | English barriers | Others | Material Resources | Diversity Issues |
|---|---|---|---|---|---|---|---|
| Accuracy | 0.81 | 0.845 | 0.839 | 0.85 | 0.697 | 0.814 | 0.981 |
| F1 | 0.530 | 0.494 | 0.502 | 0.698 | 0.487 | 0.608 | 0.609 |

## 7. Conclusion and future work

This study explores social media data in order to understand students' learning experiences in Vietnamese by integrating both qualitative analysis and data mining techniques. By the qualitative method, we found that students are struggling with heavy study load, sleep problems, language barriers, negative emotion, career targets, and diversity problems. Building on top of the qualitative analysis, we implemented and evaluated a multi-classifiers to automatically detect students' learning experiences on a dataset collected from a forum of a university in Vietnam. By applying data mining techniques, the proposed framework can overcome the limitation of analyzing large-scale data manually. The experimental results are promising, and can able to classify new posts with high accuracy. This will help administrators, educational managers to catch up immediately students' learning experiences in order to make relevant decisions to support

students and therefore enhance education quality of universities in Vietnam.

Our work is the first step toward revealing insights from informal social data in order to improve quality of education. The limitation of this work will also lead to many possible direction for future work. For examples, we did find a small number of posts refering to good things at schools. However, in this work, we only chose to focus on issues/problems because these could be the most informative for improving universities' quality. Therefore, in the future we will compare both good and bad things in students' posts. In addition, we will also investigate other texts in social media such as Facebook, Twitter, etc.

## References

[1] Z. Zerihun, J. Beishuizen, W. V. Os.: Student learning experience as indicator of teaching quality. In Educational Assessment, Evaluation and Accountability., Volume 24, Issue 2, pp 99–111. DOI: 10.1007/s11092-011-9140-4 (May 2012).

[2] J. Gordon, J. Ludlum, J.J. Hoey.: Validating the NSSE against student outcomes: Are they related? Research in Higher Education, 2008(49), 19-39 (2008).

[3] B., Liu.: Sentiment analysis and subjectivity. Handbook of natural language processing, 2, 627-666 (2010).

[4] J.P. Sue, C. Linehan, L. Daley, A. Garbett, S. Lawson: "I can't get no sleep": Discussing #insomnia on Twitter. Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, Austin, Texas, USA [doi>10.1145/2207676.2208612] (May 2012).

[5] B. Yu.: The emotional world of health online communities. Proc. of iConference 2011, February 8-11, pp. 806-807 (2011).

[6] H. Jafarkarimi; A.T.H. Sim and R. Saadatdoost: A Naïve Recommendation Model for Large Databases. International Journal of Information and Education Technology, 2 (3). pp. 216-219. ISSN 2010-3689 (June 2012)

[7] C. Romero, S. Ventura.: Educational Data Mining: A review of the state of the art. IEEE transactions on Systems, Man and Cybernetics, 40(6), 601–618(2010).

[8] N. Thai-Nghe, T. Horvath.: Personalized forecasting student performance. In: Proceedings of 11th IEEE International Conference on Advanced Learning Technologies (ICALT2011), 412–414 (2011).

[9] X. Chen, M. Vorvoreanu, and K. Madhavan.: Mining Social Media Data for Understanding Students' Learning Experiences. IEEE TRANSACTIONS ON LEARNING TECHNOLOGIES, 7(3), pp. 246-259 (2014).

[10] D. Andrzejewski, X. Zhu.: "Latent dirichlet allocation with topic-in-set knowledge". In: Proceedings of the NAACL HLT 2009 Workshop on Semi-Supervised Learning for Natural Language Processing. Association for Computational Linguistics. pp. 43–48 (2009).

[11] C. Cortes, V. Vapnik.: Support-vector networks. Machine Learning, 20(3), 273–297(1995).

[12] D.J.C. Mackay.: Information Theory, Inference, and Learning Algorithms. Cambridge University Press, 640 pages (2012).

[13] J.R. Quinlan.: Simplifying decision trees. International Journal of Human-Computer Studies, 51(2), 497–510(1999).

[14] S.R. Porter.: R. Self-Reported Learning Gains: A Theory and Test of College Student Survey Response. Research in Higher Education, 2013(54), 201-226 (2013).

[15] G. Tsoumakas, I. Katakis, I. Vlahavas.: Mining Multi-label Data. Chapter Data Mining and Knowledge Discovery Handbook, pp 667-685 (2010).