# Thai Social Media Alert System for Business

Supatta Viriyavisuthisakul[1], Parinya Sanguansat[1], Pisit Charnkeitkong[1]
and Choochart Haruechaiyasak[2]

[1]Faculty of Engineering and Technology, Panyapiwat Institute of Management
Nonthaburi, Thailand

[2]National Electronics and Computer Technology Center
Pathum Thani, Thailand

E-mail: [1]supattavir@pim.ac.th, [1]parinyasan@pim.ac.th, [1]pisitcha@pim.ac.th, [2]choochart.haruechaiyasak@nectec.or.th

**Abstract**:   Nowadays, social media has exponential growth of information generated by social interaction. In this way, social media has a huge impact on how businesses connect with customers. Pantip is the one of the most popular web board in Thailand in which the customers post their feedbacks about products and services. Among these feedbacks, some of them affects the business in negative way. To solve this problem, business should have a system that can alert their customer's negative feedbacks on social media. If the business can quickly deal with these negative feedbacks, the problems will be managed more easily. Therefore, many businesses can take advantage of this information for market analysis that can increase their opportunities. In this paper, the social media tools, that can utilized the data from social medias, are surveyed from 2005 to present. Each system has different features, but none of them cannot meet the business requirements, especially when the business cannot deal with the problem immediately because the negative sentiments are incorrectly classified and not alert. This paper proposes the social media alert system based on the business requirements. The machine learning technique is applied here to determine which information should be alerted to business. The system monitors the data in Pantip that relate to the business by keywords.  After the data were collected and preprocessed, feature vectors are extracted by Term Frequency-Inversed Document Frequency (TF-IDF) before feeding to Support Vector Machine (SVM). Experimental results show that it can achieved the good accuracy rate and also in terms of the sensitivity and specificity.

## 1. Introduction

Recently, Social media is widespread adoption of communication in general purposes for sharing experience opinion or attitude. User generated contents consist text, image, voice or video such as Facebook, Google+, Twitter, Myspace, YouTube and Microblog posts. Social media have broadcast their contents by using mobile devices, smartphones, tablets and computers. For this reason, social media is the convenient channel of retail business for customer relationship management. Customer feedbacks on social media about a product or service are extremely important to the business, if businesses are able to digest all of customer feedbacks that would cause a great impact on the business. But, in the real world, the business cannot manually manage them, because the data are very large. To solve this problem, business should have a system that can alert their customer's negative feedbacks on social media. If the business can quickly deal with these negative feedbacks, the problems will be managed more easily.

The first alert system on social media took place in 2005, as the result of social media became more popular than ever in 2004. Due to Facebook website was launched in that year [1], be the cause of the rapid increase of the number of the social media users [2]. Therefore, many businesses can take advantage of this information for market analysis that can increase their opportunities.

For sentiment analysis, there are many machine learning algorithms that are widely used to solve this classification problem, such as Naïve Bayes, Maximum Entropy, Support Vector Machine (SVM), etc. In [3], they presented a basic emotion model and classify them by SVM, that achieved a prediction accuracy at 96.43% on web data. In [4], they proposed an emotion classification of user in Twitter that related to sandy hurricane tag, which were real-time collected from 29 October to 1 November 2014. Among these data, Naïve Bayes and SVM were compared for classifying the data into 4 sentiments: positive, anger, fear and other. They concluded that SVM obtained the best accuracy in all cases. SVM get the best accuracy at 75.9% compared to 69.1% of Naïve Bayes. In [5], the business need to classify the customer opinion for Indonesian Message in Facebook. Maximum Entropy and SVM were used for sentiment classification into four classes: positive, negative, neutral and question. The result showed that SVM achieved the best accuracy at 83.5%.

In this paper, we focus on the social media in Thailand, especially Pantip [6]. It is the most popular social community web board in Thailand. Nowadays, Pantip have around 2 million members [7] and it has grown more quickly. In many times, Pantip have a significant influence on the decision-making process of consumers. Thus, many businesses in Thailand realize the importance of Pantip by supporting their customers directly via their official customer service accounts. Many products and services are reviewed and complained on this web board. This behavior is a great opportunity for business to be aware of the responding of their products, services or promotions directly from the customers. In case of customer's complaint, business should act immediately to address his complaint. However, sentiment classification from messages in social media has several challenges. Firstly, it uses Thai language which all of words need to be segmented before extracting the feature and many informal structures are always used in their posts. Second, the sentence often has many slangs and repeated letters. In [8], they collected Thai language data 1,638 documents from Pantip and then classify into 4 classes, including Pain, Gain, Need and Neutral. This research

classified the data by k-NN, comparing with 10 distances. In experimental results, Bray-Curtis distance got the best performance at 58.62% with TF-IDF. In [9], they collected the Thai comments from YouTube and classify their sentiments. Two focused genres, 2771 Thai music video (MV) and 3077 Thai commercial advertisement video (AD) were classified into Ekman's six basic emotions: anger, disgust, fear, happiness, sadness and surprise. Multinomial Naïve Bayes (MNB), Decision Tree and SVM were compared. SVM obtained the best accuracy in the commercial advertisement (AD) at 76.14% and MNB achieved the best accuracy in the music video (MV) at 84.48%.

Out of the business, the social media data are also used in other areas such as prediction the political alignment of Twitter users [10]. They predicted the data of Twitter based on contents and structures of their political communication the United States midterm elections in 2010. They found that SVM can predict the political affiliation with 91% of accuracy.

In this paper, the study found that the tools currently in have limited language skills and data source. In some tools can be supported Thai language but cannot discriminated the message as should be, especially negative language which affect the business. Thus, we focus on the negative sentiment for business in Thai and then create the alert system. We collected data from Pantip from January 2013 to October 2015 with our methodology increase the accuracy rate of prediction. We use pre-processing noisy text and sentiment classification including 2 class separate by expert. After that, in process of classification the data split 60% for training data set, 40% for test data set after that use support vector machine (SVM) [11] because, the last review of the research from 2000 to 2015. All research has yielded in the same direction is SVM most efficient to group messages on social media and adjust on proper parameters have provide the most accurate then analyze the performance of the notification of negative information and non-negatives compared the accuracy with experts. To evaluate the predictive by confusion matrix and sensitivity and specificity by ROC (Receiver Operator Characteristic) curve.

## 2. Proposed Framework

We proposed overview of steps and techniques commonly used in sentiment classification approaches, as show in Figure [1]. The system monitoring the customer's opinions. If system find a message is relevant negative comments, business going to solve the problem and respond to the customer directly or via Social media or both the way. The responsibility of business through Social media influence customers directly and indirectly to maintain a good brand image and brand royalty of the customer which cannot estimate.

From technical analysis, we collected the data from social media this interest by crawler, which the Thai language data from Pantip. We extracted data from January 2013 to October 2015 by selected from the word filter. So, alert system must be retrieved by using the keyword Social media-related businesses and stored on the database.

After the data have stored in the database, the next step going to prepare the data into pre-processing include the word wrap of the sentence (word tokenization) and stopwords removal, in process of feature extraction is calculated by sum of weight of each term in the corresponding sentence. The weight of each term is calculated by multiplication of TF and IDF, this step for convert term into numerical format which the data will be ready in process of classification. Support Vector Machine is used as a classifier method on alert system that separate sentiment into 2 groups consist of negative and non-negative, especially negative sentiment must be measured performance are highly accuracy, sensitivity and specificity have evaluated by ROC curve. If some messages are classified in the negative group that system must be alert to business for solve the problem immediately which represented by alert line, but if in the non-negative group were sent to the same for acknowledge which have the meaning do not urgency, so represented by acknowledge line.
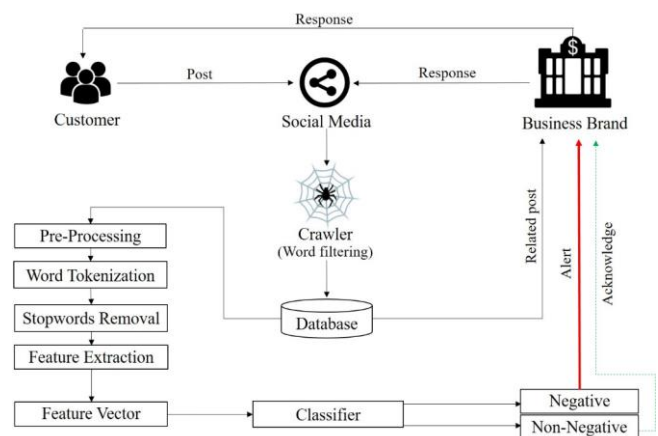


Figure 1. The proposed social media alert system.

## 3. Term Frequency and Inverse Document Frequency

In classification, feature extraction is an important process for data extraction to make the data suit to classifier. In this paper, we used TF- IDF ( Term frequency - Inversed document frequency) to transform Thai text data into a numerical statistic by using frequency. TF will provide the raw frequency of words in the document as

$$tf(t,d) = \frac{f(t,d)}{\max\{f(w,d) : w \in d\}} \tag{1}$$

When $t$ is the number of times the word appears in the document and $d$ is a corpus.

IDF measures the importance of words in the document by dividing the total number of document and taking logarithm as

$$idf(t,d) = \log \frac{N}{|\{d \in D : t \in d\}|} \tag{2}$$

When N is the number of all documents in corpus.

$|\{d \in D : t \in d\}|$ is the common adjust the word role this format suit to the term appears more than one, if the term does not appear in the document, this will change denominator to $1 + |\{d \in D : t \in d\}|$. Then TF-IDF can be calculated as

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D) \qquad (3)$$

## 4. Support Vector Machine

Support Vector Machine is a supervised learning model that can be used as a binary classifier. It finds the optimal hyperplane that can split data into two classes. As shown in Figure. 2 many hyperplanes can correctly classify data into 2 classes but SVM choose the optimal hyperplane that maximizes the margin between two classes as in Figure. 3. With this optimal hyperplane, the samples on this boundary are called Support Vectors.
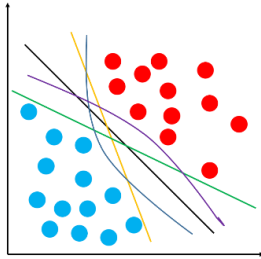


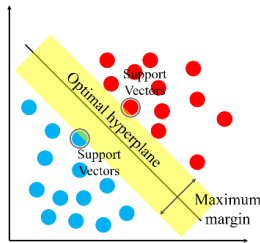Figure 2. The hyperplane can be split in 2 class.



Figure 3. The margin between the vector of the hyperplane.

Moreover, SVM can efficiently deal with the non-linear classification problem using the kernel method, that mapping input vectors into high-dimensional space in which the samples can be clearly classified. The common kernel functions include:

- Linear: $\langle x, x' \rangle$
- Polynomial: $\left( \gamma \langle x, x' \rangle + r \right)^d$
- Radial basis function (rbf): $-\gamma \left| x - x' \right|^2$
- Sigmoid: $\tanh \left( \gamma \langle x, x' \rangle + r \right)$

## 5. Experimental Results

We collect the data from PANTIP.com crawl by retail business keywords from January 2013 to October 2015, there are 67,496 documents which are labeled by the specialists into 2 classes consist of Non-negative (62,496) and negative (5,000). Non-negative appeared to outnumber the negative twelve to one here because the nature of social media is the neutral sentiment is the most part of them. In preprocessing step, word tokenization was applied by KUCUT [23] and then remove our custom stop words before feature extraction which this process going to change word (term) into feature vector by TF-IDF. After that, SVM was applied to classify these feature vectors into 2 classes

Alert system classifier experiments were performed by supervised learning using SVM in Scikit-learn [6]. The data were divided into 60% for training set and 40% for test set, respectively. For the SVM parameters, which obtained the best accuracy at 95%, are set as follows: C = 50000, degree = 1, kernel = rbf, gamma = 0.1. The confusion matrix shows the misclassifying rate in each class is around 5%, as shown in Table 1.

Table 1. The confusion matrix of test data.

|  |  | Predicted | |
| --- | --- | --- | --- |
|  |  | Non-negative | Negative |
| Actual | Non-negative | 18045 | 891 |
|  | Negative | 397 | 7666 |

For measuring the sensitivity and specificity of our system, the Receiver Operator Characteristic (ROC) curve was calculated, as shown in Figure. 4. This curve is to find the relationship between true positive rate (TPR) and false positive rate (FPT) at various threshold settings. The true positive rate is known as sensitivity and false positive rate is also known as specificity can be calculated as

$$sensitivity = \frac{true\ positive}{true\ positive\ +\ false\ negative} \qquad (4)$$

and

$$specificity = \frac{true\ negative}{true\ negative\ +\ false\ positive} \qquad (5)$$

where
- True positive (TP) is the negative comment correctly alerted as negative.
- False positive (FP) is the non-negative comment incorrectly alerted as negative.
- True negative (TP) is the non-negative comment correctly classified as non-negative.
- False negative (TP) is the negative comment incorrectly classified as non-negative.

The possibility is extremely low the value of TP and FP will peak simultaneously, because when the high sensitivity will result in low specificity and the low sensitivity will result in high specificity. If desired, the high value of specificity, false positive must drop. The area under ROC curve (AUC) can indicate the overall performance of our alert system. This area measures the discrimination, that is the ability of the test to correctly classify user's post with and without the need of notification. Consider the situation

in which the posts are already correctly classified into two groups. The post with the more negative should be the one from the need to alert group. The area under the curve is the percentage of randomly drawn pairs for which this is true. In our system, we got AUC at 95% that is quite good result, 85% specificity limit means that at most 3 in 20 normal cases analyzed may yield a false alarm.
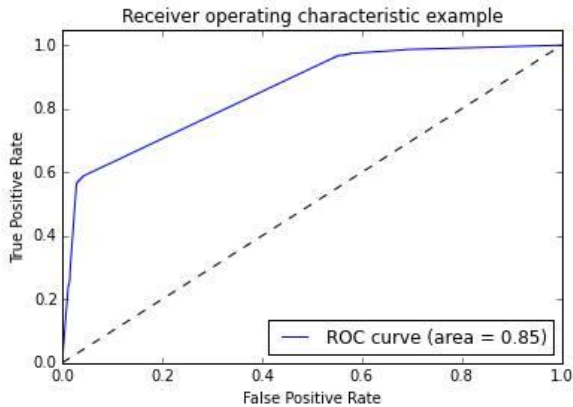


Figure 4. The ROC curve show the performance of SVM classifier.

## 4. Conclusion

This paper proposes a system framework for social media monitoring which collects the business related data from the most Thai famous web board, Pantip. The crawler uses business related keywords to acquire the data and store them in the database in which the classifier performs the sentiment classification on each post into negative and non-negative. The alert system will notify to business when it found the negative sentiment while non-negative class will be acknowledged to business for market analysis. With the fine-tuning parameters of SVM, as reported in our experiments, the results show that it obtained accuracy at 95%, the confusion matrix indicates the error rate of both classes is around 5% and the sensitivity and specificity can cover 85% of the AUC.

In future work, we intend to use the machine learning to automatic response for the negative post in many social medias. The data will be crawled from many social medias, such as Facebook and Twitter. In this way, business can extensively handle their customers. To improve the accuracy of the system, Thai word tokenization need to be improved and more specific to the words in our focused social medias. In additional, we will develop the system that can notify to business in many ways that can make business can deal with the problem immediately.

## Acknowledgement

## References

[1]    S. Phillips, "A brief history of Facebook," *the Guardian*, 2007. [Online]. Available: http://www.theguardian.com/technology/2007/jul/25/media.newmedia. [Accessed: 11-Dec-2015].

[2]    "Facebook Reports First Quarter 2015 Results." 2015.

[3]    H. Binali, C. Wu, and V. Potdar, "Computational Approaches for Emotion Detection in Text," in *4th IEEE International Conference on Digital Ecosystems and Technologies - Conference Proceedings of IEEE-DEST 2010, DEST 2010*, 2013, pp. 172 – 177.

[4]    J. Brynielsson, F. Johansson, C. Jonsson, and A. Westling, "Emotion classification of social media posts for estimating people's reactions to communicated alert messages during crises," *Secur. Inform.*, vol. 3, no. 1, p. 7, Aug. 2014.

[5]    A. Naradhipa, A.R.; Purwarianti, "Sentiment classification for Indonesian message in social media," in *Proceedings of the 2011 International Conference on Electrical Engineering and Informatics*, 2011, pp. 1–4.

[6]    "Pantip - Learn, Share & Fun," 1997. [Online]. Available: http://pantip.com/. [Accessed: 20-Dec-2015].

[7]    "Pantip.com Site Info." Alexa Internet, 2014.

[8]    S. Viriyavisuthisakul, P. Sanguansat, P. Charnkeitkong, and C. Haruechaiyasak, "A comparison of similarity measures for online social media Thai text classification," in *2015 12th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology (ECTI-CON)*, 2015, pp. 1–6.

[9]    P. Sarakit, T. Theeramunkong, C. Haruechaiyasak, and M. Okumura, "Classifying emotion in Thai youtube comments," in *2015 6th International Conference of Information and Communication Technology for Embedded Systems (IC-ICTES)*, 2015, pp. 1–5.

[10]   M. D. Conover, B. Goncalves, J. Ratkiewicz, A. Flammini, and F. Menczer, "Predicting the Political Alignment of Twitter Users," in *2011 IEEE Third Int'l Conference on Privacy, Security, Risk and Trust and 2011 IEEE Third Int'l Conference on Social Computing*, 2011, pp. 192–199.

[11]   Scikit-learn, "sklearn.svm.SVC — scikit-learn 0.17 documentation," *BSD License*, 2010. [Online]. Available: http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html. [Accessed: 27-Dec-2015].