# Opinion Strength Identification in Customer Review Summarizing System Using Association Rule Technique

Parnicha Apisuwankun[1, a] and Janjao Mongkolnavin[1, b]

[1]Department of Statistics, Chulalongkorn Business School,
Chulalongkorn University, Bangkok 10330, Thailand
[a]parnicha.tfr@gmail.com, [b]janjao@cbs.chula.ac.th

*Abstract*— **The purpose of this research is to study opinion strength identification in customer review summarizing system in Thai language. We expect that opinion strength identification will help refining results from customer review summarizing system and improving quality of the results. In this paper, we present results from our preliminary study that was conducted with an aim at creating rules for identifying opinion strength by using an association rule technique. We used a paper survey to collect human opinion towards words used in 180 reviews of three groups of cosmetic products: blush on, facial cleanser, and foundation (60 reviews per group). The survey was done with 225 samples and resulted in 11,686 phases that had effects on samples' opinion towards the products. Association rule technique was applied to words in those phases to generate opinion strength identification rules. In our further study, those rules will be used to extend a capability of a customer review summarizing system. We plan to compare an effeciency of the system which is extended with the opinion strength identification rules to the one that is not, and also to a summary result obtained from human.**

*Keywords—Opinion strength; Opinion strength identification; Customer review summarizing; Association rules.*

## I. INTRODUCTION

Online customer reviews are useful for product owners as word-of-mouth marketing [4]. The reviews also can be used as feedbacks for improving their products or services to enhance their competitiveness [2], [6]. In aspects of potential customers, these reviews can help them make purchase decision [9], [16]. However, at present, number of product reviews on web sites has increased rapidly. As a result, it is very difficult for a potential customer to read all reviews. On the other hand, it is also difficult for product owners to track and aggregate every customer reviews [2], [6], [14]. Therefore, researchers are interested in finding approaches to summarize customer reviews that are expressed in form of text on various websites by using opinion mining methods [9], [6], [14]. Many customer review summarizing techniques have been proposed [2], [1], [8]. Most techniques can be divided into two steps: (1) feature extraction and (2) orientation identification (whether customers feel positive or negative towards each feature of the product) [2], [1], [8].

In addition to orientation or direction of opinions that appear in review text, there is also an opinion strength which indicates a degree of opinion that the reviewer has on each particular feature of the product [11], [10]. Reviews that have the same orientation can be very different in their degrees. For example, for a product such as eyeliner, a review on brand A can be "It can draw a very sharp line !!! that lasts all day long and the color is very black" and a review on brand B can be "Dark color and sharp line. It also doesn't mess up under the eyes." Both reviews are positive, but they show varying degrees of positive opinion. The review on Brand A sounds more positive than the one on Brand B.

The works that are related to customer review summarizing in Thai language include [12], [3], and [15]. However, their studies focused on extracting features of products or services and identifying orientation of the opinions towards those features. Up to our knowledge, a research on opinion strength identification in Thai language is not found.

In this article, we discuss related researches in Section II. In Section III, we explain our preliminary study including our approach to create opinion strength identification rules by using an association rules technique, the results, and an approach to incorporate them into a customer review summarizing system. Our conclusion and further experimental study is discussed in Section IV.

## II. RELATED WORKS

In this section, we discuss researches related to each part of a customer review summarizing system including Thai word segmentation and part-of-speech tagging, feature extraction, orientation identification, opinion strength

identification, and the overview of a customer review summarizing system using opinion strength identification.

## A. Thai Word Segmentation

There is no doubt that customer reviews that are available in various websites are treated by computers as text. Thus, they are needed to be segmented into words and tagged with proper part-of-speech before submitted to other processes.

At present, word segmentation techniques that are most commonly used for Thai text include Thai Lexeme Tokenizer (LexTo) [13] and Thai Lexeme Analyser (TLexs) [13]. LexTo uses a Longest Matching technique based on dictionary from Lexitron [13]. It also allows users to add specific words to the dictionary so that those words can be recognized by LexTo. TLexs uses Conditional Random Field model (CRF) [7] built from a five million word corpus. In our study, we used LexTo as our word segmentation tool because of its flexibility in recognizing ambiguous words that are not included in the dictionary. In context of customer review where spoken language is used more often than written one, LexTo appears to be a more appropriate tool.

## B. Feature Extraction

Feature extraction is a process to identify words that represent features of products or services in customer reviews. Researches that are related to extraction of product features in customer reviews that are written in Thai include Prombut (2007) [12], Haruechaiyasak et al. (2010) [3], and Thumrongluck (2010) [15].

In Prombut (2007) [12] and Haruechaiyasak et al. (2010) [3], words that represent features of products or services were determined by researchers. Thus, those words were subjective to both researchers and those particular domains of products or services. Thumrongluck (2010) [15] used the discrimination-based term extraction method [5] to extract features of products. In this method, Term Frequency Inverse Class Frequency (TFICF) of each term is calculated from frequency of words that appear in reviews of each product category. Words that have highest TFICF in each product category are chosen as feature words of that product category [15]. In our study, we will use Thumrongluck (2010)'s technique because it can be applied across different product categories, which makes it more general in adoption.

## C. Opinion Words and Orientation Identification

Orientation of opinions in customer reviews can be classified as positive or negative polarity through opinion words used. Researches related to opinion word identification and orientation of opinion in Thai text include Prombut (2007) [12], Haruechaiyasak et al. (2010) [3], and Thumrongluck (2010) [15].

Thumrongluck (2010) used a Wordnet which is a bipolar adjective structure to find synonym or antonym words to determine polarity of a found opinion word. The Wordnet is built from a seed list of opinion words which are commonly used to describe positive or negative opinion of each product category. A new word that is either a synonym or antonym of words in the seed list is added to an appropriate polar. A polarity of a found opinion word in the review is determined by the polar of the Wordnet that it belongs to. In Prombut (2007) [12] and Haruechaiyasak et al. (2010) [3], each opinion word was determined by researchers, which is its limitation in generalization to different product categories.

In order to identify positions of opinion words in review text, Prombut (2007) [12] and Haruechaiyasak et al. (2010) [3] created a possible list of sentence patterns, then matched the review text to a pattern in the list to identify feature and opinion words. Sentence patterns used in Prombut (2007) were created by using grammatical structure in Thai language, while Haruechaiyasak et al. (2010) [3] extracted sentence patterns in spoken language that are commonly used to express opinion in each product features. In Thumrongluck (2010), opinion words of each product feature were detected by Reverse-Distance-Weighting (RDW) method which was developed by Oelke et al. (2009) [5]. The concept of this method is using a distance (cutoff) between an opinion word and a feature word to determine whether that opinion word belongs to that particular feature word or not. An opinion word which is nearest to the feature word will be determined as the opinion that the reviewer had on that product feature, and the polarity of that opinion word determines its orientation towards that particular feature [15]. RWD method is appropriate for Thai text where there is no symbol to signal the end of each sentence [15]; therefore, RWD is used in our study.

## D. Opinion Strength Identification

Up to our knowledge, there is no research on identifying opinion strength in customer review in Thai language. However, there are related studies in sentiment strength in English text which include Thelwall et al. (2010) [11], Taboada et al. (2011) [10] and Meng (2012) [17].

The algorithm of sentiment strength detection in Thelwall et al. (2010) called "SentiStrength". In their work, the opinion word list contains words and strength of individual words which were determined by Human Coder Subjective Judgment. Scale ratings of both positive and negative comments are ranging from 1 to 5 (Scale 5-Point): (no positive emotion) 1, 2, 3, 4, 5 (very strong positive emotion), and (no negative emotion) -1, -2, -3, -4, -5 (very strong negative emotion) [11]. Thelwall et al. has proposed rules to increase degree of opinion strength, which are listed below [11].

- If negation word appears close to opinion word, invert opinion orientation. For example, if "very happy" has positive strength 4 then "not very happy" would have negative strength -4.
- If more than two repeated letters are used, boost 1 point of strength. For example, the word "haaaappy" has one more point of positive strength than the word "happy".
- If at least one exclamation mark is found, boost at least

2 points of strength for both positive and negative opinions. For example, the word "OMG!!!!!!!!!!" (Oh my god!!!!!!!!!!).

Taboada et al. (2011) [10] and Meng (2012) [17] use a similar method called "Semantic Orientation CALculator (SO-CAL)". Their intensifier list contains words and modifier which is a percentage of increasing or decreasing opinion strength. To calculate strength of opinions, they start by detecting an intensifier that appears close to an opinion word and then using its modifier to assign a degree of strength from score of opinion word [10], [17]. The limitation of this method is that their word list seems to be more appropriate for formal language which is not a context of customer reviews in websites where spoken language and slang are used rather often. Thus, we chose to partly adapt the human coder judgment approach of Thelwall et al. (2010) [11] to determine strength of Thai opinion words.

*E. The Review Summarizing System.*

Putting the relevant researches together, we can draw the whole picture of the review summarizing system as showed in Figure 1. The system can be divided into three main parts. Our preliminary study is related to the third part of the system where strength of the opinion found is adjusted by intensifier words (intensifiers) in the review. In the next section, we will discuss the approach that we use to generate opinion strength identification rules and the results.
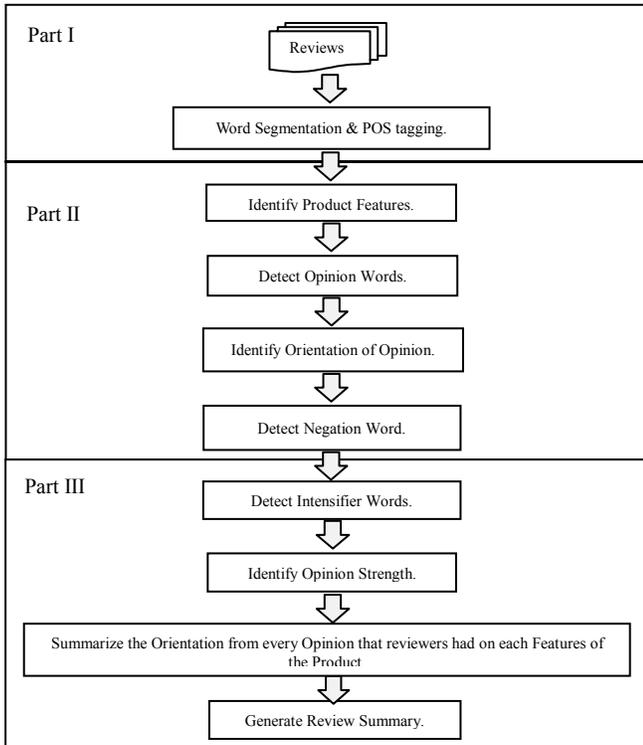


Figure 1. Three parts of the customer review summarizing system.

## III. CREATING OPINION STRENGTH IDENTIFICATION RULES BY USING ASSOCIATION RULE TECHNIQUE

*A. Data Collection*

The data used for calculating opinion strength include the opinion word list and the associate strength degree of each word which indicates how much the opinion word goes on positive or negative orientation. From the literature review, we found opinion word lists and strength degrees in English only. For our study which focuses on Thai text, an opinion word list has been gathered by using a paper questionnaire. In the questionnaire, each sample is asked to read each product review and identify words that affect his/her opinion towards that product or service, and rate each word with score -3 to 3 : -3 (very strong negative opinion), -2, -1 (slightly negative opinion), +1 (slightly positive opinion), +2, +3 (very strong positive opinion) as showed in Table I.

TABLE I.    THE SYMBOLS USED IN OPIONION STRENGTH RATING AND THEIR MEANINGS.

| Symbols | Meaning |
|---------|---------|
| + | "slightly positive opinion" |
| + + | "positive opinion" |
| + + + | "very strong positive opinion" |
| - | "slightly negative opinion" |
| - - | "negative opinion" |
| - - - | "very strong Negative opinion" |

180 reviews of three product categories: blush on, facial cleanser, and foundation were used to build questionnaires. They were latest 180 reviews of three product categories from jeban website (www.jeban.com : the most popular cosmetic product review website in Thailand). In each product category, reviews of three brands were selected and in each brand, twenty reviews were included. Thus, there were 20 reviews for each of the 3 x 3 brands from three product categories. Nine groups of twenty reviews were used to build 9 questionnaires.

After a pilot test, we distributed each type of questionnaire to 25 samples. We, thus, have 225 (9 x 25) samples in total. Words and their associate strength were identified by human coders (samples) using symbols as showed in Table I. Figure 2 shows an example of the results from human coders.



Figure 2.  An example of the results from samples (human coders).

18

## B. Data Preparation and Data Mining Process

An example of data collected with the questionnaires is showed in Figure 2. The data are then rearranged into a table of phases and opinion strength as showed in Table II. The table consists of 11,686 rows. After that LexTo was used to segment each phase in words. The results from this step are showed in Table III. After data preparation, an association rule technique was employed to find associations between words and opinion strength by setting a minimum support percentage and a minimum confidence level to 3 and 10 respectively. SAS Enterprise Miner 6.1 was used to conduct this step.

TABLE II. AN EXAMPLE OF THE DATA TABLE WHICH IS SUMMARIZED FROM THE QUESTIONNAIRES.

| ID | Word | Strength |
|----|------|----------|
| 1 | ชอบกลิ่นมาก | 2 |
| 2 | หอม | 1 |
| 3 | ถูกใจมาก | 3 |
| 4 | ไม่ขาวเท่าไหร่ | -1 |
| 5 | ไม่ถึงกะหมอง | -1 |
| 6 | ธรรมชาติ | 2 |
| 7 | เนียนมาก | 3 |
| 8 | ไม่ปกปิด | -2 |

TABLE III. THE DATA TABLE AFTER WORD SEGMENTATION

| ID | Word | Strength |
|----|------|----------|
| 1 | ชอบ ∣ กลิ่น ∣ มาก ∣ | 2 |
| 2 | หอม ∣ | 1 |
| 3 | ถูกใจ ∣ มาก ∣ | 3 |
| 4 | ไม่ ∣ ขาว ∣ เท่าไหร่ ∣ | -1 |
| 5 | ไม่ ∣ ถึง ∣ กะ ∣ หมอง ∣ | -1 |
| 6 | ธรรมชาติ ∣ | 2 |
| 7 | เนียน ∣ มาก ∣ | 3 |
| 8 | ไม่ ∣ ปกปิด ∣ | -2 |

## C. Results

From all results that were generated from SAS Enterprise Miner 6.1, we selected only rules that showed significant associations between words and strength by considering their confidence, support, and lift. The selected rules are showed in Table IV where the left and the right side of symbol ==> represent word and its strength respectively.

Since the reviews used in this study are from the same product categories as Thumrongluck (2010) [15], we found words that belong to the opinion word list of Thumrongluck (2010) (Showed in highlighted rows in Table IV). For example, word "สวย (beautiful)" is associated with +1 point strength and word "มัน (oily)" is associated with -1 point strength, which is consistent to Thumrongluck (2010) [15]. That is, the results confirm that those words are opinion words

not intensifier words (or words used to identify opinion strength) because they are either associated to +1 or -1 point strength. So they can only be used to identify orientation of opinion whether it is positive or negative.

TABLE IV. RULES THAT SHOW ASSOCIATIONS BETWEEN WORDS AND OPINION STRENGTH.

| No. | CONFI-DENT | SUP-PORT | LIFT | COUNT | RULE |
|-----|-----------|----------|------|-------|------|
| 1 | 30.36 | 4.92 | 2.09 | 575 | ไม่ ==> -1 |
| 2 | 35.14 | 4.38 | 1.29 | 512 | มาก ==> 2 |
| 3 | 50.26 | 4.18 | 1.61 | 488 | ดี ==> 1 |
| 4 | 34.21 | 3.12 | 1.25 | 364 | ชอบ ==> 2 |
| 5 | 50.74 | 2.34 | 3.50 | 273 | ไม่ & ค่อย ==> -1 |
| 6 | 44.93 | 2.20 | 1.44 | 257 | สวย ==> 1 |
| 7 | 46.64 | 2.02 | 1.49 | 236 | สะอาด ==> 1 |
| 8 | 20.18 | 1.69 | 1.63 | 198 | ๆ ==> 3 |
| 9 | 40.31 | 1.57 | 3.26 | 183 | กก ==> 3 |
| 10 | 40.00 | 1.54 | 1.28 | 180 | แห้ง ==> 1 |
| 11 | 61.92 | 1.49 | 1.98 | 174 | โอเค ==> 1 |
| 12 | 60.89 | 1.41 | 1.95 | 165 | ไม่ & แห้ง ==> 1 |
| 13 | 39.85 | 1.34 | 1.27 | 157 | เนียน ==> 1 |
| 14 | 27.62 | 1.28 | 1.90 | 150 | ทน ==> -1 |
| 15 | 59.67 | 1.24 | 4.11 | 145 | ไม่ & ทน ==> -1 |
| 16 | 45.03 | 1.24 | 1.44 | 145 | น่ารัก ==> 1 |
| 17 | 42.42 | 1.20 | 3.43 | 140 | มากก ==> 3 |
| 18 | 42.42 | 1.20 | 3.43 | 140 | มากก & กก ==> 3 |
| 19 | 44.48 | 1.17 | 1.42 | 137 | ง่าย ==> 1 |
| 20 | 39.03 | 1.17 | 1.25 | 137 | ปกปิด ==> 1 |
| 21 | 45.02 | 1.12 | 1.44 | 131 | ตึง ==> 1 |
| 22 | 60.29 | 1.05 | 1.93 | 123 | ไม่ & ตึง ==> 1 |
| 23 | 20.55 | 1.02 | 1.42 | 119 | มัน ==> -1 |
| 24 | 46.06 | 0.95 | 3.17 | 111 | ไป ==> -1 |

However, Thumrongluck (2010) [15] focused only on identifying orientation of opinion and not on its strength which is the emphasis of our work. For this reason, we filter rules that are related to Thumrongluck (2010)'s opinion words from Table IV. The rest of the rules are showed in Table V.

## D. Applying opinion strength identification rules in customer review summarizing system

We propose an approach to apply opinion strength identification rules in customer review summarizing system as Step 5 in the following 6 steps to summarize customer reviews (showed in Figure 3):

1) Identify feature(s) of the product in a review.

2) Identify opinion(s) in the review that is related to that feature.

3) Identify orientation of the opinion by score 1 or -1 (positive or negative respectively).

4) Detect a negation word (if there is any) to invert

orientation of the opinion.

5) Apply the opinion strength identification rules if it exists.

6) Summarize the orientation and strength of opinion that reviewers had on each feature of the product from all reviews.

TABLE V.　THE RULE FOR IDENTIFING OPINION STRENGTH

| No. | Words | Strength | Example |
|---|---|---|---|
| 1 | "มาก (very)" | 2 | "สวยมาก (very beautiful)", "ดีมาก (very good)" |
| 2 | "ชอบ (like)" | 2 | "ชอบเม็ดสี (like the pigment)", "ชอบค่ะ (like)" |
| 3 | "ๆ" (repeated word symbol) | 3 | "สวยมากๆๆ (very very very beautiful)", "เนียนสุดๆ (super super smooth)" |
| 4 | "มากก&กก" (very&yy)" (more than 2 repeated letters) | 3 | "ชอบมากกก (reallyyyy like)", "สวยมากกก (veryyyy beautiful)" |
| 5 | "โอเค (ok)" | 1 | "ราคาโอเค (the price is ok.)", "คุมมันโอเค (oil control is ok.)" |
| 6 | "ไป (too)" | -1 | "เข้มไป (too dark)", "หนาไป (too thick)" |

In Step 5 which is the process to identify degree of opinion strength, we will calculate opinion strength by using rule showed in Table V. The polarity of the opinion word from step 4 will be multiplied by the opinion strength of the rule that matches words in the review. Examples of score calculation can be showed as follow.

Example 1. "ราคาถูกมากกกก (The price is veryyyy cheap)"

- Product feature is showed by word "ราคา (price)".
- Opinion word is showed by word "ถูก (cheap)". (Polarity of word "ถูก" is +1.)
- Negation word does not exist.
- Intensifier is showed by the word "มากกกก (veryyyy)" (Opinion strength of word"มากกกก" is +3 points.)

So the score for the above opinion is 3 (1 × 3) points. The opinion can be summarized that product feature is price and its opinion score is +3, which means that the reviewer has a very strong positive opinion towards the product's price.

Example 2. "สีเข้มไป (The color is too dark.)"

- Product feature is showed by word "สี (color)".
- Opinion word is showed by word "เข้ม (dark)". (Polarity of word "เข้ม" is +1.)
- Negation word does not exist.

- Intensifier is showed by word "ไป (overly / too)". (Opinion strength of word "ไป" is -1 points.)

Thus, the score for the opinion in Example 2 is -1 (1 × -1) point. The opinion can be summarized that product feature is color and its opinion score is -1, which means that the reviewer has a slightly negative opinion towards the product's color.
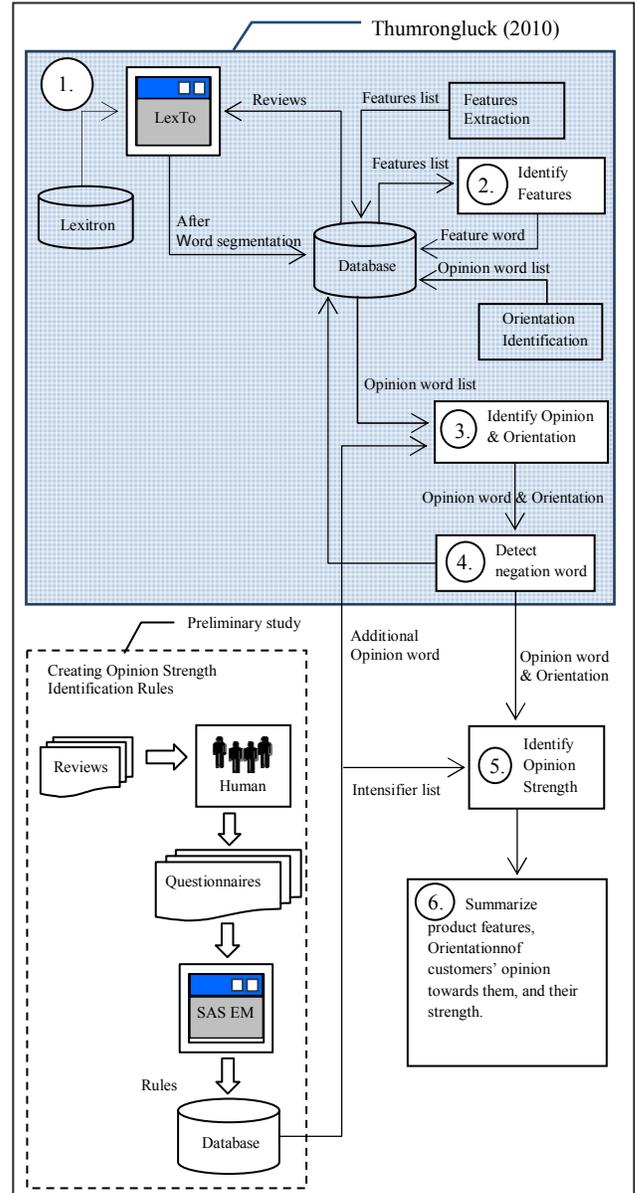


Figure 3. Overview of the customer review summarizing system using opinion strength identification rules.

IV.　CONCLUSION AND FURTHER EXPERIMENTAL STUDY

In our preliminary study, phases in customer reviews that have effect on human opinion were collected from human coders together with their associate opinion strength using

nine paper questionnaires. After appropriate data preparation, the association rule technique is employed to find significant associations between words and opinion strengths. With our approach, we expect that an efficiency of an automated product review summarization system can be improved. This is because in spoken language where proper grammars are not strictly used and the language itself may continuously evolved over time, a data mining technique such as association rules may be used to capture patterns if they may exist in an informal communication such as web blogs. In the result of our preliminary study, we found opinion strength identification rules especially intensifiers that do not exist in Thumrongluck (2010) [15] and we expect that those rules would help increasing precision and recall of the result obtained from the automated customer review summarizing system when compared to results from human.

In our further experimental study, we plan to conduct three experiments on the same three product categories as we used in our preliminary study. However, different sets of reviews will be used in the experiments. The three experiments include the automated customer review summarizing system that uses opinion strength identification rules, the one that does not use the rules, and human.

## REFERENCES

[1] A.-M. Popescu and O. Etzioni, "Extracting product features and opinions from reviews," in Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing. Vancouver, British Columbia, Canada: Association for Computational Linguistics, October 2005.

[2] B. Liu, M. Hu, and J. Cheng, "Opinion observer: analyzing and comparing opinions on the Web," Proceedings of the 14th International Conference on World Wide Web, 2005, pp. 342-351.

[3] C. Haruechaiyasak, A. Kongthon, P. Palingoon, and C. Sangkeettrakarn, "Constructing Thai Opinion Mining Resource: A Case Study on Hotel Reviews," Proceedings of the 8th Workshop on Asian Language Resources, 2010, 64–71.

[4] D. H. Shin, "User experience in social commerce: in friends we trust," Behaviour & Information Technology, 2012, pp. 1-16.

[5] D. Oelke, M. Hao, C. Rohrdantz, D. Keim, U. Dayal, L. Haug, and H. Janetzko, "Visual Opinion Analysis of Customer Feedback Data," IEEE Symposium on Visual Analytics Science and Technology, Atlantic City, New Jersey, USA, 2009.

[6] H. Chen, "AI and Opinion Mining," IEEE Intelligent Systems, 2010, vol. 25, no. 3, pp. 74-80.

[7] J. Lafferty, A. McCallum, and F. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data", Proceedings of the 18th International Conference on Machine Learning. Morgan Kaufmann, 2001, pp. 282–289.

[8] L. Zhang, B. Liu, S. H. Lim, and E. O'Brien-Strain, "Extracting and ranking product features in opinion documents," Proceedings of the 23rd International Conference on Computational Linguistics: Posters, 2010, pp. 1462-1470.

[9] M. Hu and B. Liu, "Mining and summarizing customer reviews," in KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining. New York, NY, USA: ACM, 2004, pp. 168–177.

[10] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-Based Methods for Sentiment Analysis," Association for Computational Linguistics, 2011, vol. 37, no. 2, pp. 1-42.

[11] M. Thelwall, K. Buckley, G. Paltoglou, and D. Cai, "Sentiment Strength Detection in Short Informal Text," Journal of the American society for information science and technology, 2010, vol. 61, no.12, pp. 2544-2558.

[12] N. Prombut, "Mining Opinion in Product Reviews : A Case Study of Mobile Phone Reviews," Master's Project, Department of Computer Science Faculty of Applied Sciences King Mongkut's University of Technology North Bangkok, 2007.

[13] National Electronics and Computer Technology Center (NECTEC), "About LEXiTRON" [online]. 2008. Available from: http://lexitron.nectec.or.th/2009_1/index.php?q=common_manager/aboutlex#about [2013, January].

[14] S. M. Mudambi, and D. Schuff, "What makes a helpful online review? A study of customer reviews on amazon.com," MIS Quarterly, 2010, vol. 34, no. 1, pp. 185-200.

[15] T. Thumrongluck, "An Automated System for Summarizing Structured Product Reviews," Master's Thesis, Department of Statistics Faculty of Commerce and Accountancy Chulalongkorn University, 2010.

[16] Y. Chen, S. Fay, and Q. Wang, "The Role of Marketing in Social Media: How Online Consumer Reviews Evolve," Journal of Interactive Marketing, 2011, pp. 1-32.

[17] Y. Meng, "Sentiment Analysis: Study on Product Features," Dissertations and Theses from the College of Business Administration, 2012, pp. 1-100.