# Opinion Mining and Sentiment Analysis Application for Opinion Classification from Education Questionnaire

Amir Hamzah[1]
[1]Departement of Informatics Engineering
Institut Sains & Teknologi AKPRIND
Yogyakarta, Indonesia

*Abstract*— Measurement of academic services using questionnaires with multiple choice answers generally provide comments and advice columns. In the data analysis results, comments and suggestions made by the thousands of students cannot be utilized due to the lack of analysis tools. Whereas comments and suggestions could be actually contain student opinions on various things, such as facilities, faculty, library and others. Opinion mining and sentiment analysis as a new tool in text mining can be applied to the data to utilize comments and suggestions. This research applied HMM-POS Tagger to give automatically POS TAG to the sentence based on training POS TAG data by using Hidden Markov Model. By implementing POS TAG pattern the comments can then be determined whether it was opinion or not. Furthermore if it were opinion it can be determined its target and also the orientation of the opinion whether it is positive or negative. The data used was 1,000 comments given POS-TAG manually and 1,000 comments as test data. Sentiment analysis is applied using four methods of classification, namely SVM, NBC, ME and KM-Clustering. The result showed that the accuracy of POS-Tagger was 0.95 and the average of accuracy of four classification method was 0.85.

*Keywords: HMM POS Tagger; opinion classification*

## I. INTRODUCTION

Opinion and opinion orientation are the most important part of decision making for a policy. The right decision is strongly influenced by opinion analysis from various sources related to decision making. For example in the business world, the addition of products by the production manager requires the analysis of the product review of goods on the market. Other examples such as the management of educational services in universities, the measurement of the level of service learning satisfaction can be measured from the opinions of students about the learning process. Opinions appear in a variety of situations, e.g. deliberately requested by an opinion assessment tool through a request for suggestions in a questionnaire activity, or appear naturally from an online forum provided by the college's official site. The volume of online opinion in the form of free text is getting more and more and generally not utilized because of its unstructured shape.

The existence of internet and other on-line information sources is growing very rapidly. Online data and information from companies and organizations are generally unstructured and generally in the form of text that reaches 80% [5]. The emergence of social media such as Facebook (2004) and Tweeter (2006) has encouraged activities such as reviews, discussion forums, blogs, micro-blogs, comments, and posts that multiply the existence of text documents on the internet. This is because the social media has been used both by individuals and organizations for various interests in conducting information sharing activities. This explosive condition of information further complicates the process of data mining as well as it has been predicted [10]. Therefore, the development of research in the field of opinion mining becomes a very important topic in addition to the previous topics, namely data mining and text mining.

One branch of research that evolved from the information explosion situation on the internet was sentiment analysis and opinion mining. Opinion mining is a challenging research because there is an accumulation of research challenges from *Information Retrieval* (IR) : *Information Extraction*, *Information Summarization*, *Document Classification* and from *Natural Language Processing* (NLP) fields such as *Named Entity Recognition* (NER) and *Document Subjectivity Analysis* [9]. This research branch examines how one extracts opinions from online media and performs an analysis of those opinions. Sentiment Analysis or opinion mining is a computational study of people's opinions, appraisals and emotions through entities, events and attributes [6].

The application of sentiment analysis and opinion mining to conduct policy evaluation and decision making promises a more practical and economical way than the classical method using the questionnaire approach. Critics of the method of questionnaire is that it is very time consuming and expensive method, while it is also providing results that sometimes cannot catch the real problem. Questionnaires and interviews are considered weak because in general people do not like to answer survey questions that are sometimes long-winded. In this position opinion mining answers the question of opinion extracting by listening rather than by asking like a questioner (*by asking*), thus more accurately reflecting the true reality [13]. Even further opinion mining allows capturing the emotions of opinion owners [7].

The Institute of Science and Technology AKPRIND as a higher education institution constantly wants to improve services in learning management. For that purpose at the end of each semester the academic administration should evaluate the learning service using the questionnaire instrument with the items provided. So far, there are data of questionnaires that cannot be used and analyzed that is student's suggestion data.

This data amounts to thousands of suggestions or more exactly the opinion of all participants from all courses. Suggestions or opinions can be about *academic atmosphere*, *lecturer*, *class room*, *air conditioners*, *Over Head Projector*, or other campus facilities. Within a few years this data is getting stacked up unusable.

The problem in this research is how to build a software prototype that can extract opinions from a collection of commentary text documents, and then determine the target of opinion and opinion orientation.

The purpose of this study is to review the application of opinion mining techniques and sentiment analysis to analyze the data suggestions / opinions of students. This research is also designed to create prototype of opinion mining software and sentiment analysis that can extract opinion, analyze opinion, map target of opinion and set opinion orientation

## II.    RELATED WORKS

The study of opinion mining applications conducted by [4] on the *English National Health Service* website that captures 6,412 free comments from treated patients. Analysis of comments related to *hygiene*, *hospital services* and various aspects of hospital responsibilities has resulted in conform between 81% to 89% compared with quantitative rating method that was provided through the questionnaire.

In opinion mining the method to detect opinions is by applying HMM-POS Tagger. This method is applied to the collection of trainer data which is the opinion sentence and non-opinion sentence that has been given POS-TAG or Part of Speech's sign of the word on each word. The program is expected to extract opinions from the commentary text of the questionnaire, as well as search for the object of opinion.

There are several approaches to automated POS tagging, which are *rule-based*, *probabilistic*, and *transformational-based approaches*.  Rule-based tagger POS sets a tag to a word based on some of the manual linguistic rules created, for example a word is tagged NOUN if it follows AJECTIVE. The probabilistic approach determines the word tag of a token based on the context probability of manually specified token tags from a corpus. The transformational-based approach combines rule-based and probabilistic approaches to automatically derive symbolic rules from the corpus [11]. The use of the Hidden Markov Model for POS Tagger Bahasa Indonesia was examined by [14] resulted in accuracy of 96.2% and [15], which resulted in an accuracy of 92.2%

## III.    THEORITICAL BACKGROUND
### A.  HMM POS-TAGGER

HMM-POS Tagger is a method for performing POS-tagging in a sentence automatically based on an analysis and characteristic of POS-tag data from training data collection. Part-of-Speech (POS) tagging known as grammatical tagging is the process of providing a POS-tag to a word in a sentence text. Part-of-speech is a grammatical category of words in a sentence, e.g. verb (VERB), noun (NOUN), adjective (ADJECTIVE), and others. POS tagging is an important tool in many natural language processing applications such as disambiguation, parsing, question-answer, and machine translation systems. Because assigning manual part-of-speech tags to words in a sentence is costly, exhausting, and time-

consuming, there has been a wide interest in the automation of the POS tagging process [3].

If there is a sentence consisting of n words ($w_i$: i = 1, .., n), and will be assigned a POS-tag for each word that composes the sentence ($t_i$: i = 1, .., n), then this issue can be formulated as seeking the maximum value of:

$$\hat{t} = \arg\max_{t_1^n} P(t_1^n \mid w_1^n) \tag{1}$$

By applying Bayes  theorem in conditional probability, then (1) can be written to be:

$$\hat{t} = \arg\max \frac{P(w_1^n \mid t_1^n) P(t_1^n)}{P(w_1^n)} \tag{2}$$

Since the value of the denominator is always the same for every sentence, then (2) can be written to be:

$$\hat{t} = \arg\max P(w_1^n \mid t_1^n) P(t_1^n) \tag{3}$$

By making two assumptions, then equation (3) can be written:
(1) The probability of a word depends only on its POS-tag.

$$P(w_1^n \mid t_1^n) \approx \prod_{i=1}^{n} P(w_i \mid t_i) \tag{4}$$

(2) The probability of a POS-tag depends only on the previous POS-tag.

$$P(t_1^n) \approx \prod_{i=1}^{n} P(t_i \mid t_{i-1}) \tag{5}$$

By applying (4) and (5) to (3) to be obtained:

$$\hat{t} = \arg\max \prod_{i=1}^{n} P(w_i \mid t_i) \, P(t_i \mid t_{i-1}) \tag{6}$$

Furthermore from the results of HMM-POS Tagger will be generated sentences that have been given POS-tag. In the next step based on these POS-Tag patterns it will be determined whether a text is an opinion or not an opinion. With POS-Tag patterns can also be determined an object of opinion.

### B.  SENTIMENT ANALYSIS

Sentiment analysis for comments that have been detected as an opinion is done by applying an opinion classification to determine the opinion orientation. There are many methods available in text classification. In this research, we used four methods, namely *Naive Bayes Classifier* (NBC) method, *Support Vector Machine* (SVM), *Maximum Entropy* (ME) and K-Means Clustering (KMC).

### NBC Method

The NBC method assumes the collection of opinion documents as D ={d1, d2, ..., d | D |} and collection of categories C = {c1, c2, ..., c | C |}. The NBC classification is done by finding the probability P (C = $c_j$ | D = $d_i$), i.e. the probability of category $c_j$ if the document is found in. The document is seen as a tuple of words, i.e. <$w_1$, $w_2$, ..., $w_n$>, whose frequency of occurrence is assumed to be a random variable with Bernoulli probability distribution [8]. The document classification is to find the maximum value of:

$$V_{MAP} = \arg\max_{c_j \in C} P(C_j \mid w_1, w_2, ..., w_n) \tag{7}$$

By providing the Bayes theorem obtained

$$V_{MAP}= \arg\max_{c_j \in C} \frac{P(w_1, w_2, ..., w_n \mid c_j)P(c_j)}{P(w_1, w_2, ..., w_n)} \quad (8)$$

Since the denominator value is constant for a document, and assuming that each word is independent of each other the equation (8) can be written as:

$$V_{MAP} = \arg\max_{c_j \in C} \prod_{i=1}^{n} P(w_i \mid c_j)P(c_j) \quad (9)$$

Practically the calculation of P ($c_j$) is approximated by:

$$P(c_j) = \frac{\mid doc_j \mid}{\mid contoh \mid} \quad (10)$$

| $doc_j$ | is the number of document categories j and | sample | is the number of sample documents (training).

In computation P ($w_i \mid c_j$) is approached by:

$$P(w_i \mid c_j) = \frac{\mid n_i + 1 \mid}{n + \mid vocabulary \mid} \quad (11)$$

$n_i$ is the frequency of occurrences of the word $w_i$ in category $c_j$, and n is the word frequency in the category document $c_j$ and | vocabulary | is the number of occurrences of all words in the sample document collection.

*SVM Method*

*Support Vector Machine* (SVM) was first developed by [1] and continued with a more detailed description by [2]. The SVM concept can be explained as the search for the best hyperplane that serves as a separator of two classes in the input space. For a collection of documents in the form of:

$$\mathcal{D} = \{(x_i, y_i) \mid xi \; \varepsilon \; Rp, \; yi \; \varepsilon \; \{-1,1\}\} \quad (12)$$

where $y_i$ is 1 or -1, indicates which class the $x_i$ is located. Each $x_i$ is a real p-dimensional vector. There will be a maximum margin hyperplane that divides the points for points that have $y_i = 1$ of which has $y_i = -1$.

Any hyperplane can be written as set of points x satisfying w. x - b = 0 where (.) indicates dot product. The vector w is a normal vector that is perpendicular to the hyperplane. Parameter ‖ w ‖ determine the offset of hyperplane from the origin along the normal vector w. It will be attempted to choose w and b to maximize margins, or the distance between parallel hyperplanes apart as far as possible to separate data.

These hyperplanes can be described by the equation:

$$w. x - b = 1 \quad (13)$$

and

$$w. x - b = -1 \quad (14)$$

The distance between the two hyperplanes is 2 / ‖ w ‖, so we want to minimize ‖w‖. To prevent data points from falling into the margins, the following limits must be added:
For each value i

$$w. x_i - b \geq 1: \text{first class x} \quad (15)$$

or

$$w. xi - b \leq -1: xi \text{ second class} \quad (16)$$

It can be rewritten as:

$$y_i (w \; x_i - b) \geq 1,$$
$$\text{for all } 1 \leq i \leq n \quad (17)$$

So the problem of finding the maximum hyperplane is the optimization problem:
Minimize ‖ w ‖ with constraints for every i = 1, ..., n

$$y_i (w \; x_i - b) \geq 1 \quad (18)$$

*Maximum Entrophy (ME) Method*

Classification with Maximum Entropy (ME) applies information theory. Entropy is the average set of information contained in a set of events X = {$x_1$, $x_2$, ..., $x_3$} which can be expressed in:

$$H(p) = \sum_{x \in X} p(x) \log_e \left( \frac{1}{p(x)} \right) \quad (19)$$

With the value of H (p) is the set of information from the set of events X, and p (x) is the probability of the occurrence of x in the set X. The Maximum Entropy (ME) method is a method to maximize the value of H (p). The maximum value of H (p) will be obtained if the value of X is uniform so that p (x) = 1 / | X | with | X | is the cardinality of the set X.

The application of the ME Method for document classification is performed by the conditional probability approach of a class of documents when a document is present. Suppose that the set of class of documents is A = {$a_1$, $a_2$, ..., $a_c$} and the collection of documents is D = {$d_1$, $d_2$, ..., $d_n$}. The determination of class a of document d will be seen by determining the conditional probability value p (a | d) of maximum value of the probability distribution with maximum entropy.

*KM-C Method*

K-means clustering approach is clustering by using a cluster center as grouping criterion. The cluster center is the average value of all cluster members' objects. Suppose we have a collection of documents D = {$d_i$ | i = 1,2, ... | D |} = {$d_1$, $d_2$, ..., $d_{\mid D \mid}$} to be clustered into K clusters. In this case $d_i$ is a real valued vector that represents the document. The vector has a dimension n, which is the number of unique words in the document collection. Document collections can be represented by the n x |D| size matrix, denoted by [$X_{ij}$], with the $x_{ij}$ element representing Term Frequency (TF) i.e. the occurrence of the term (word) i in the document j. To obtain better accuracy in computation, the matrix containing the term frequency value is converted into a real element matrix that takes into account the frequency of occurrence of documents containing the i-th word by weighting *Invert Document Frequency* (IDF). Furthermore, it is also strived that the length of the document vector is always 1 by normalizing the document vector. Weighting that combines TF with IDF then known as TF-IDF weighting that can be formulated as:

$$w_{ij} = \frac{(\ln(f_{ij})+1).\log\left(\dfrac{N}{n_i}\right)}{\sqrt{\left((\ln(f_{ij})+1).\log\left(\dfrac{N}{n_i}\right)\right)^2}} \qquad (20)$$

K-Means clustering algorithm is done by taking K initial vector as the seeds of the centers of the cluster. Furthermore the entire vector document is calculated distance against each cluster center. The document vector that closest to a cluster center then is defined to be a new member of the cluster. The cluster center is the average vector of all vectors in a given cluster.

The algorithm of K-means can be written as follows:
1) Take K object as seed from K center cluster
2) For all objects: locate the cluster with the closest distance, and assign the object into the cluster.
3) Recalculate the center of the cluster with the average of all objects in the cluster
4) Calculate the criteria function and evaluate. If the criterion function does not change the algorithm stops.

## IV. PROPOSED METHOD

### A. Research Materials

The materials used in this research are POS TAG for HMM-POS Tagger and pairs of POS TAG as rules for opinion detection and target detection which were taken from [14] and [12] with modification. These rules are in the form of POS TAG collections. The POS TAG related to the HMM POS Tagger is shown as in TABLE I, while the POS Tagger data training that use this table is shown in Fig.1. The rule of opinion detection was based on POS TAG pair as shown in TABLE II. The rule of target detection was based on POS TAG pairs as shown in TABLE III.

The source of comments data were text documents consisting of comment collection from suggestion on student questionnaire from Institute of Science and Technology AKPRIND for 6 semesters :
1) Semester 1 2014/2015: 4.128 comments
2) Semester 2 2013/2014: 3.152 comments
3) Semester 1 2013/2014: 3.801 comments
4) Semester 2 2013/2014: 2.551 comments
5) Semester 1 2012/2013: 3.663 comments
6) Semester 2 2012/2013: 1883 comments

From these collections 1000 documents were selected as training data and 1000 documents were selected as test data.

### B. Data Analysis

The research steps taken in solving the problem are divided into two steps: the *pre-processing* step of the document and the step of *detection* and *extraction of opinion*, the target detection and finally the classification of opinion

### Pre-processing Document and training

Documents in the form of student comments with irregular form are converted into standard sentence. The next step, the parsing process in the collection of training data so that it can be determined the frequency of occurrence of each POS-TAG word in the collection of train data with POS-TAG

code as in Table I. Collection of data training for POS-Tagger has format as shown in Fig 1.

```
<DOC-0001>tolong/VBI gedungnya/NNG
dibersihkan/VBP <*DOC>
<DOC-0002>birokrasi/NNU kampus/NNC yang/PR
jelek/JJ membuat/VBT
mahasiswa/NNC merasa/VBT dikhianati/VBP
<*DOC>
<DOC-0003>kursi/NNC di/IN dalam/IN
ruangan/NNU kurang/RB
mengenakkan/JJ <*DOC>
<DOC-0004>lebih/RB ditingkatkan/VBP lagi/RB
<*DOC>
<DOC-0005>acnya/NNG kurang/RB dingin/JJ
<*DOC>
...
<DOC-0999>pasang/VBT ac/NNC pak/UH
soalnya/NNG panas/JJ <*DOC>
```

Fig. 1. POS-Tagger data training format

The collection of training data for opinion classification has format as shown in Fig. 2. In the training data the orientation of opinion document has been manually set, which are positive opinions and negative opinions. The orientation of opinion on the training data is labeled on the document in the form of code (+) and (-).

```
<DOC-0001(+)>mengajarnya sudah oke, perlu
ditingkatkan </DOC>
<DOC-0002(-)> kamar mandi kotor sekali,
tolong dibersihkan yang rutin </DOC>
<DOC-0003(+)>pak joko ngajarnya memang
joss..</DOC>
<DOC-0004 (-)>kipas angin kurang baik,
fasilitas wivi tolong lebih ditingkatkan
lagi biar baik </DOC>
<DOC-0005(-)>lebih ditingkatkan lagi
fasilitas perpustakannya </DOC>
<DOC-0006(+)>birokrasi kampus sangat jelek
urusan terlalu bertele-tele </DOC>
```

Fig. 2. Classification data training format

After POS Tagger training data has been processed, the information of POS-TAG to each word in the sentence collection can be determined. The next step, i.e. opinion extraction and opinion classification was conducted follow the flowchart as shown in Fig.3.

TABLE I. POS TAGGER DATA

| No | POS | POS Name | Contoh |
|----|-----|----------|--------|
| 1 | OP | Open Parenthesis | ({[ |
| 2 | CP | Close Parenthesis | )}] |
| 3 | GM | Slash | / |
| 4 | ; | Semicolon | ; |
| 5 | : | Colon | : |
| 6 | " | Quotation | " and ' |
| 7 | . | Sentence terminator | . |
| 8 | , | Comma | , |
| 9 | - | Dash | - |
| 10 | … | Ellipses | … |
| 11 | JJ | Adjective | Baik, Bagus |
| 12 | RB | Adverb | Sekali, sangat |

| 13 | NNC | Countable Noun | Kursi, Kulkas |
|----|-----|----------------|---------------|
| 14 | NNU | Uncountable Noun | Gula, hujan |
| 15 | NNP | Proper Noun | Toyota, Sony |
| 16 | NNG | Genetive Noun | Motornya |
| 17 | VBI | Intransitive Verb | Pergi |
| 18 | VBT | Transitive Verb | Membeli |
| 19 | VBP | Passive Verb | ditingkatkan, diperbaiki |
| 20 | IN | Preposition | Di, Dari, Ke |
| 21 | MD | Modal | Bisa |
| 22 | CC | Coor-Conjunction | Dan, Atau, tetapi |
| 23 | SC | Subor-Conjunction | Jika, Ketika |
| 24 | DT | Determiner | Para, Ini, Itu |
| 25 | UH | Interjection | Wah, Aduh, Oi |
| 26 | CDO | Ordinal Numerals | Pertama, Kedua, Ketiga |
| 27 | CDC | Collective Numerals | Berdua |
| 28 | CDP | Primary Numerals | Satu, Dua, Tiga |
| 29 | CDI | Irregular Numerals | Beberapa |
| 30 | PRP | Personal Pronouns | Saya, Mereka |
| 31 | WP | WH-Pronouns | Apa, Siapa, Dimana |
| 32 | PRN | Number Pronouns | Kedua-duanya |
| 33 | PRL | Locative Pronouns | Sini, Situ |
| 34 | NEG | Negation | Bukan, Tidak |
| 35 | SYM | Symbols | #,%,^,&,* |
| 36 | RP | Particles | Pun, Kah |
| 37 | FW | Foreigns Words | Word |

TABLE II.    RUL E FOR OPINION DETECTTION

| No | Rule | Examples |
|----|------|----------|
| 1 | RB JJ | sangat buruk,  dengan bagus , memang jelek |
| 2 | RB VB | semoga berjalan,  jika memilih |
| 3 | NN JJ | LCDnya jelek, alatnya bagus |
| 4 | NN VB | Ngajarnya membosankan, perkataannya menjengkelkan |
| 5 | JJ VB | mudah difahami, cepat memahami |
| 6 | CK JJ | bagus atau baik, tetapi malas |
| 7 | JJ BB | sama bagus |
|  |  |  |
| 8 | VB VB | membuat pusing, membikin bosan |
| 9 | JJ RB | indah sekali, bagus sekali |
| 10 | VB JJ | membikin bingung |
| 11 | NEG JJ | tidak seindah, tidak semudah |
| 12 | NEG VB | tidak mengerti, tidak memahami, bukan mengajar |
| 13 | PRP VBI | saya menyukai, kita suka |
| 14 | PRP VBT | kita suka |
| 15 | VBT NN | memiliki kedekatan, memiliki kepekaan |
| 16 | MD VBT | Perlu mengambil referensi |
| 17 | MD VBI | Perlu dikembangkan |
| 18 | UH  VBP | Tolong dicat, tolong diperbaiki |
| 19 | JJ  VBP | Mudah diterima, sulit dipahami |

TABLE III.   RULE FOR TARGET DETECTTION

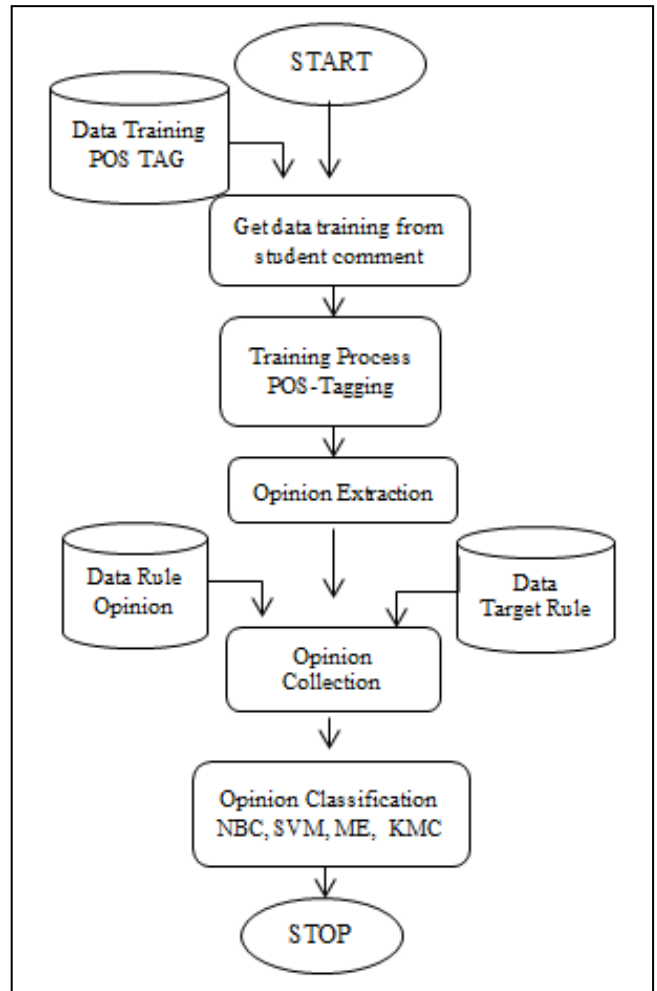| No | Rule | Examples |
|----|------|----------|
| 1 | NN | Ac, lcd, internet |
| 2 | NNG | Laboratoriumnya, lcdnya |
| 3 | NNP | Pak Joko, Bu Yuli, Pengok |
| 4 | NN NN | kantin kampus, ac pengok |
| 5 | NN CC NN | Kampus dan lab |
| 6 | NN IN NN | AC di Klas |



Fig. 3.  Flowchart of Opinion extraxtion and classification

*Opinion Extracttion  and Target detection*

The step of opinion extraction was done after pre-processing and POS tagging had been done. By applying the rule fo opinion detection from TABLE II the document  using opinion and opinion targeting rules using target rule is determined as the scheme in Fig.3. After opinions and opinions targets was obtained the last step was to make an opinion classification to determine whether an opinion be positive or negative.

## V.    RESULT AND DISCUSSION

The prototype for *Pre-processing* document collection for HMM POS-Tagger has been developed in the interface as shown in in Fig.4.
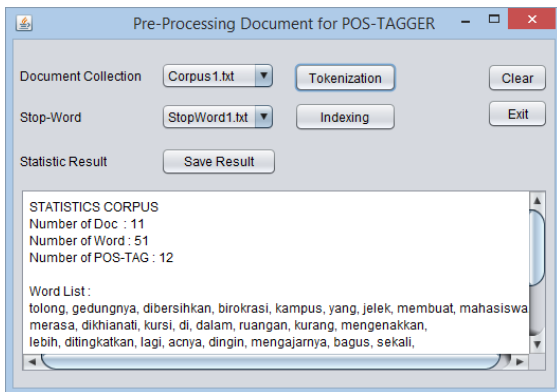


Fig. 4.    The result of document pre-processing

For HMM POS-Tagger demonstration to assign POS TAG to the document that is input to the system as shown in Fig.5. The sentence "*tolong gedungnya dibersihkan*" is tagged by system as depicted in Fig.5 :  *tolong/VBI  gedungnya/NNG dibersihkan/BP* with accuracy 100%.
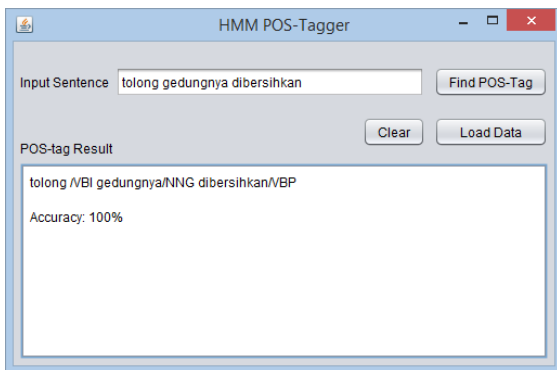


Fig. 5.    Result of  HMM POS Tagger

Opinion extraction from various target opinions can be shown in Fig.6., where target can be chosen selected from combo box component.
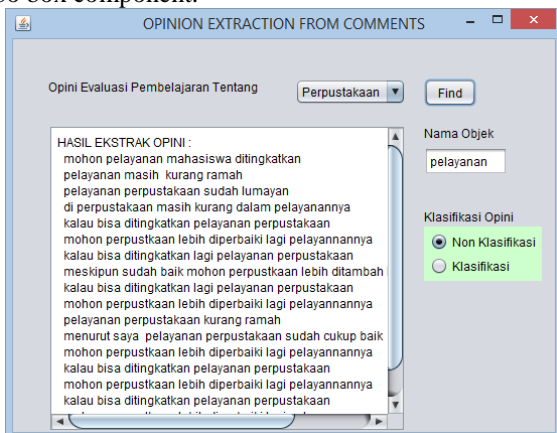


Fig. 6.    Opinion Ekstraction without classification

The classification of opinions into two class i.e. positive or negative  can be demonstrated in Fig.7. The classification method can be selected from four options (NBC, SVM, MAxEnt and KM-Klus). Those figure also shown that NBC method classification for target "Perpustaaan" has accuracy of 81.54%.
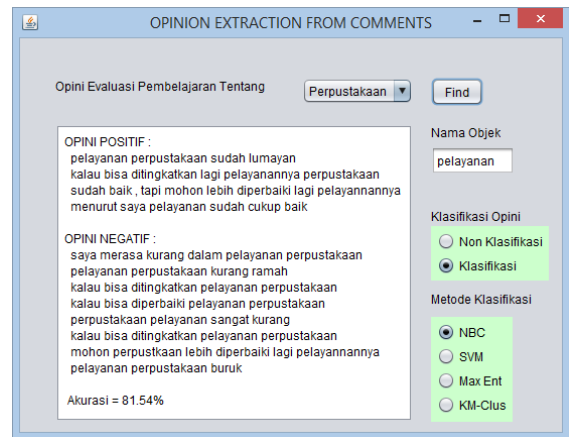


Fig. 7.    Opinion Ekstraction with classification

Performance Testing Program

Program performance is measured through precision and recall parameters for POS Tagger HMM performance in opinion detection and target detection. Test results using 1000 data training and 1000 test data is shown in Table IV.

TABLE IV.    PERFORMANCE OF HMM POS TAGGER

| PerformanceMeasurement | Application of HMM-POS Tagger | |
|---|---|---|
| | Opinion Detection | Target Detection |
| Precision | 0.95 | 0.91 |
| Recall | 0.92 | 0.89 |

The performance of the classification program is measured through accuracy parameters. The comparison of four methods test performance is shown in Table 5.  The result showed that KMC-method has best performance compared to the three other methods.

TABLE V.    PERFORMANCE OF CLASSIFICATION

| PerformanceMeasurement | Classification Methods | | | |
|---|---|---|---|---|
| | *NBC* | *SVM* | *ME* | *KMC* |
| Accuracy | 0.84 | 0.83 | 0.84 | 0.88 |

## VI.    CONCLUSION

In this study the prototype of preprocessing the collection of comments has been developed. The system prototype of opinion extraction and classification has also been developed. The program has been demonstrated how to extract the opinion using HMM POS-Tagger which is equipped with sentiment analysis using four opinion classification methods, namely NBC, SVM, ME and KM-Clustering.  Accuracy in opinion detection and object opinion detection averages over 90% of class of opinion opinions with four methods of accuracy reaching an average above 80%. The result also showed that the KMC-method has the best performance in accuracy compare to the three other method.

## REFERENCES

[1] Boser,B.E.,Guyon, I.M. and Vapnik,V.N.,1992, "A Training Algorithm for Optimal Margin Classifiers",Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory,1992,pp. 1171-1183.

[2] Cortes,C. and Vapnik,V., 1995,"Support-Vector Networks", Machine Learning, 20, pp.273-297

[3] Cutting, D., Kupiec,J., Pesderson,J. and Sibun,P., "A Practical Part-ofspeech Tagger, Xerox Palo Alto Research Center", in Proceeding of the third Conference on applied Natural Language Processing , 1992, pp.133-140.

[4] Greaves, F., D.R. Cano, C. Millet, A.Darzi, and L. Donaldson, "Use of Sentiment Analysis for Capturing Patient Experience From Free-Text Comments", Journal of Medical Internet Research, 2013, 15:11, e239. Online publication date: 1-Jan-2013

[5] Grimes, S., 2013, *Unstructured Data and the 80 Precent Rule*, Clarabridge Bridgepoints.

[6] Liu,B.,2010, "Sentiment Analysis: Multi Facet Problem", IEEE Intelligence System, 25 (3),pp:76-80

[7] Loia, L. and Senatore, S.,"A fuzzy-oriented sentic analysis to capture the human emotion in Web-based content*",* Knowledge-Based Systems 58, 2014, pp. 75-85 Online publication date: 1-Mar-2014

[8] McCallum, A. and Nigam, K., 1998, "A Comparison of Event Models for Naive Bayes Text Classification", AAAI/ICML-98 Workshop on Learning for Text Categorization, pp. 41-48

[9] Pang, B., Lee, L. and Vaithyanathan, S., 2002, Thumbs up?: sentiment classification using machine learning techniques", Proceeding EMNLP '02 Proceedings of the ACL-02 conference on Empirical methods in natural language processing - Volume 10, pp: 79-86

[10] Putten,P.V.D., Kok, J. and Gupta,A., 2002, *Why the Information Explosion can be Bad for Data Mining, and How Data Fusion Provides a Way Out*, Proc.of the 2nd SIAM International Conference on Data Mining,pp:11-13

[11] Pisceldo. F., Manurung, R. and Adriani,M. "Probabilistic Part-of-Speech Tagging for Bahasa Indonesia", Third International MALINDO Workshop, colocated event ACL-IJCNLP, Singapore, 2009

[12] Rozi,I.F., Pramono,S.H. dan Dahlan, E.A.,2012, Implementasi Opinion Mining (Analisis Sentimen) untuk Ekstraksi Data Opini Publik pada Perguruan Tinggi, Jurnal EECCIS, Vol.6,No.1,Juni 2012

[13] Shelke,N.M., Deshpande,S. and Thakre,V., "Survey of Techniques for Opinion Mining," International Journal of Computer Applications (0975 – 8887) Volume 57– No.13, November 2012

[14] Wicaksono, A.F. and Purwarianti,A., 2010, "HMM Based Part-of-Speech Tagger for Bahasa Indonesia", Proceedings of the 4th International MALINDO Workshop, Jakarta.

[15] Widhiyanti, K and A. Harjoko, "POS Tagging for Bahasa Indonesia dengan HMM dan Rule Based", INFORMATIKA Vol.8.,No.2., November 2012, pp.151-167