

Aspect-based sentiment analysis on hair care product reviews

Malki Kothalawala*

Faculty of Graduate Studies and Research
Sri Lanka Institute of Information Technology, Sri Lanka
malkinimesha@gmail.com

Samantha Thelijagoda

SLIIT Business School
Sri Lanka Institute of Information Technology, Sri Lanka
samantha.t@slit.lk

Abstract: Nowadays, with almost everything being shared online, people are more verbal about their consumer experiences with products via reviews. Reviews can be vital for manufacturers to get insights into consumer opinions and consumers in their purchase decisions. Sentiment analysis, referring to the extraction of subjective opinions on a particular subject within a text, is a field within Natural Language Processing, that can convert this unstructured information hidden within reviews into structured information expressing public opinion. In regards to a specific product group like hair care products, certain brands are rising in the market due to their positive public opinion on particular aspects. While e-commerce websites facilitate users to view the reviews, they do not display which reviews contain which type of opinion on which aspect at a glance. This research aims to introduce an automated process that focuses on determining the polarity of online consumer reviews on different aspects of hair care products by using Aspect-based Sentiment Analysis. The system consists of processes like data gathering, pre-processing, aspect extraction and polarity detection and follows a sequential approach to achieve the intended goal. Consequently, by deciphering the aspect-wise polarity of the reviews, the implemented system demonstrates an accuracy of 85% from the test data for overall aspects, enabling consumers to get an at a glance idea about the public opinion and manufacturers to identify their strong and weak points.

Keywords: Aspect-based Sentiment Analysis, Natural Language Processing, Opinion Mining, Sentiment Analysis, Supervised Learning

I. INTRODUCTION

In a highly competitive business world, it is vital for brands and companies to get insights into consumer behaviors, their needs, likes and dislikes about a particular product, buying patterns and how consumers make decisions. Traditional businesses used to depend on practices like surveys in order to get insights into the consumer opinions. However, in the current modern world, businesses can utilize the power of the vast amount of both public and private information being generated on a daily basis due to the continuously growing number of users and devices connected to the internet. With an increasing amount of information being shared via the means of social channels, review sites, blogs and web forums, people tend to be more verbal about their consumer experiences online. Moreover, while traditional approaches like surveys often influence the responses of the consumers due to their use of target questions, reviews enable companies to get the natural opinions of the consumers. Furthermore, from the consumers' perspective, their purchase decisions are majorly influenced by the reviews and the experiences of the other consumers about the precise product which are available online. However,

consumers often have to read a long list of subjective reviews and opinions in order to decide whether the product fits their need since it is not visible at a glance and requirements are different from consumer to consumer.

Since general product reviews is a broad topic, it is specialized into a more specific type of reviews like the hair care product reviews. In regards to hair care product market, it has a wider scope with various types of hair care products being available. With different brands introducing various product lines intended for numerous hair goals for each category of hair care products, the market has become hypercompetitive. Due to factors like geographical limitations, ease of search offers and eliminating travel time and cost, consumers tend to opt for e-commerce websites like eBay, Amazon for buying their hair care products. When purchasing a hair care product via an e-commerce website, consumers come with different requirements in their mind and they often decide whether to purchase the product or not, by considering the reviews of the other consumers to validate against their hair goals. Moreover, once they use the product, they share their opinions so that those reviews can be helpful for the other consumers. For instance, a consumer might be interested in purchasing a certain shampoo from eBay and wishes to know the public opinion on how moisturizing the shampoo is, which is his or her hair goal. However, if most of the top reviews are talking about the scent of the shampoo, which therefore is not helpful for the consumer. As a result, the consumer often has to read a long list of reviews to find out whether the product matches his or her needs. Furthermore, in relation to the hair product manufacturers, they could be interested in knowing whether the consumers are satisfied with the ingredients of their new hair product line so that they can change their formulas according to consumer feedback. A mechanism of presenting the polarity of the reviews on the hair care product in terms of the aspects being talked about would be ideal in both of the aforementioned scenarios.

Sentiment analysis, a field within Natural Language Processing (NLP), referring to the extraction of subjective opinions on a particular subject within the text, can be utilized in comprehending most of the unstructured information available in review sites, social media and web forums, etc. Sentiment analysis aids in extracting information from a text about the polarity, which denotes whether the speaker

expresses a positive or negative opinion, subject, which topic is being discussed and the opinion holder, the individual conveying the opinion. With the ever-increasing amounts of textual data in the world and most of them being unstructured, sentiment analysis is highly beneficial for commercial applications while enabling organizations to gain crucial

business insights by deciphering unseen meanings of the text in order to make better business decisions. While sentiment analysis can be assistive in deciding whether a person is expressing a negative, positive or neutral opinion, Aspect-based Sentiment Analysis (ABSA) is capable of uncovering which aspect or attribute of the product the person expresses the opinion on. It is significant to know precisely which feature of the hair care product is being talked about since consumers often look for reviews on different aspects.

Hence, the purpose of the system is to decipher the polarity of online consumer reviews on different aspects of hair care products using aspect-based sentiment analysis so that it can be beneficial for both the consumers and the product manufacturers.

II. RELATED WORK

A. Rule-based approaches

In a research conducted by Sonal Meenu Singh and Nidhi Mishra, it classified customer reviews on mobile phones as positive, negative or neutral based on manually extracted eight aspects including cost, size, battery, camera, OS, processor, storage and screen. The authors have conducted the research collecting 80 mobile phone reviews from commercial websites like www.amazon.com and www.cnet.com on Iphone 6, Moto G3 and blackberry Z10. The reviews collected have been pre-processed by removing stopwords, white spaces and special symbols and POS tagging has been performed on the pre-processed reviews with the assistance of Stanford parser. With the aid of Part-Of-Speech (POS) tagging, opinion words have been identified and stored from the reviews. SentiWordNet 3.0 has been used to identify scores for each aspect after the completion of pre-processing. In comparison with the online tools such as WordNet, SentiWordNet and MPQA, the system is depicted to have an accuracy of 75% and according to the results, the authors state that the system is capable of handling negation, intensifiers and synonyms more effectively than the compared tools. The drawback of the system is described as its inability to accurately identifying the polarity of the cost feature [1].

Paramita Ray and Amlan Chakrabarti have carried out research on sentiment analysis for product reviews using the Lexicon method. A framework for analyzing the sentiment of users on Twitter using the Twitter API has been proposed with the aid of R software. As the initial step, the authors have created a Twitter application in order to use Twitter API. After the completion of data collection, the tweets have been analyzed and pre-processed. As the following step, classification has been performed with the assistance of a lexical method or a dictionary-based approach. In order to test the results, tweets with the keyword 'Iphone' posted within a specific time period has been analyzed considering aspects like voice quality, battery quality, service, price, picture quality and size. Apart from being able to identify emoticons in tweets, the research was also done in both document level and aspect level analysis [2].

Moreover, I. K. C. U. Perera and H. A. Caldera have done a research on aspect-based opinion mining which specifically focuses on restaurant reviews. Reviews on 20 restaurants have been collected from the website www.zomato.com. For the purpose of extracting aspects, the pre-processed reviews have been tokenized with the aid of a POS tagger. Based on the

frequency of occurrence of words, aspects have been manually selected as place, food, service, time and staff. Using Stanford dependencies, which represent grammatical relations between words, opinion words have been identified and the filtered-out aspect words and opinion words have been sent to SentiWordNet. SentiWordNet assigns scores for the key value pairs and once all assigning scores for all opinion words against a certain aspect, the total of all scores has been calculated. Polarity has been decided based on the score value. The testing has been done both manually and systematically. In the systematic approach, the aspect polarity is checked by the system and in comparison to the manually checked results, the systems are depicted to show an accuracy of 70% [3].

B. Automatic approaches

A research conducted by Zeenia Singla, Sukhchandan Randhawa, and Sushma Jain focused on sentiment analysis of consumer product reviews. The Statistical and Sentiment Analysis of Consumer Product Reviews (SACP) framework developed by the authors, is composed of two major components where one module handles Data Collection and Pre-processing and the other module works on Feature Selection and Analysis. The Data Collection module has collected more than 400,000 online reviews on around 4500 mobile phones from the e-commerce website Amazon.com. The data sets have been put to CSV file format and the data has then been pre-processed in order to remove white spaces, stop words, digits, punctuation and special symbols. The authors have used a package called *tm* for the purpose of text mining. As the initial process of the feature selection and analysis module, the extraction of related features from the data is performed. The module then carries out statistical analysis of the data in order to analyze the correlation between the features. As the final step of the module, the system performs sentiment analysis on the text in order to predict the sentiment as well as the polarity. The authors state that the classification proves to be efficient based on the 84.87% accuracy of Support Vector Machine (SVM) after cross validation [4].

Dhanalakshmi V., Dhivya Bino and Saravanan A. M. have conducted a research on opinion mining from student feedback data in 2016 using a supervised learning algorithm. The authors have used the tool Rapid Miner in mining and classifying the feedback provided by the students. The authors have used Support Vector Machine (SVM), Naïve Bayes (NB), K Nearest Neighbor (KNN) and Neural Networks (NN) as the supervised opinion mining algorithms. The authors have used the validation operator in Rapid Miner in order to concurrently train and test the classifiers. Using the Performance Operator in Rapid Miner, for each algorithm, accuracy, precision and recall values have been calculated for comparing the performance of the algorithms. From the results, the authors have identified that a 100% best result for precision is showed by K Nearest Neighbor (KNN) where a best of 97.07% and 99.11% for recall and accuracy has been shown by Naïve Bayes (NB) algorithm [5].

C. Hybrid approaches

Santhosh Kumar K L, Jayanti Desai, and Jharna Majumdar have conducted a research on opinion mining and sentiment analysis on online customer reviews. The system automatically extracts reviews from the Amazon website and classifies the reviews as positive or negative using algorithms such as Naïve Bayes classifier, Logistic Regression and

SentiWordNet. Based on the experimental results, the research concludes that Naïve Bayes classification happens to be the most efficient algorithm for opinion mining and that the system can be further extended with the aid of reviews from additional websites in order to decide the most efficient test classifier [6].

In research conducted by Francis F. Balahadia, Ma. Corazon G. Fernando and Irish C. Juanatas, they have proposed a performance evaluation tool for teachers with the aid of sentiment analysis. The system uses a dataset polarity created by the authors which contains a list of dictionaries of words. Using the Naïve Bayes algorithm, opinions have been extracted and analyzed in order to determine the polarity. As future work, the authors have intended to add a spelling and grammar checker to the system and also so enlarge the dataset dictionary they have created in order to improve the classification [7].

Through the literature review, it was identified that automatic or hybrid approaches prove to depict higher accuracy levels in comparison with the rule-based approach. Hence, a hybrid approach combining both machine learning algorithms and lexical resources was opted for the proposed system.

III. METHODOLOGY

The main goal of the system is to perform aspect-based opinion mining on the hair care product reviews. Since hair care products fall under various categories, the research specifically focused on hair products falling under the category of shampoo. Furthermore, due to shampoos being produced aiming different hair goals, there are numerous aspects to be considered. Therefore, four general aspects consisting of effectiveness, price, fragrance and formula have been opted manually after going through a range of reviews on hair shampoos. The system follows a sequential approach in order to achieve the objective. The entire system flow can be categorized into four major components including Data Gathering, Pre-Processing, Aspect Extraction and Polarity Detection. Fig. 1. illustrates a detailed process of the system architecture.

A. Data gathering

The preliminary step was to gather data or to collect consumer reviews, which act as the initial inputs to the system. In order to collect the reviews, the system facilitates users to enter a URL of an eBay product page. Using the provided URL, the system performs HTML parsing on the webpage with the aid of JSOUP library [8]. The system is capable of traversing through different pages of reviews and extracting all reviews on the product. The collected review data is then written to a CSV file. The system only collects the name of the person who wrote the review, the subject of the review and the review body on which aspect-based sentiment analysis is performed.

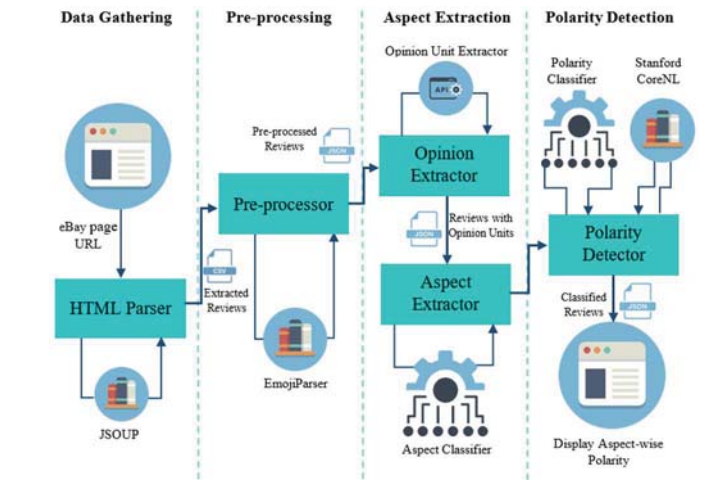


Fig. 1. High level diagram

B. Pre-processing

Once the data gathering stage is completed, the next stage is the pre-processing. The CSV file containing the reviews extracted from the data gathering stage is passed to the pre-processing module as the input. Since most of the consumer reviews tend to contain unnecessary noise such as smileys, emoticons, and capitalized words, which makes the reviews difficult to analyze, pre-processing is an essential step. The first step of the pre-processing module is to check whether the review is written in English. In case of it being written in any other language, that review is not taken for further processing. A regular expression is used to check whether the review is written in English language. Unnecessary special characters are eliminated and needlessly capitalized words are converted into simple letters. With the aid of EmojParser library, emojis included within the review are removed [9]. Moreover, emoticons within the review are also eliminated by using a regular expression pattern.

C. Opinion extraction

Since a single sentence of a review can contain information on multiple aspects, opinion units within a sentence of a review are broken down with the aid of opinion unit extracting API provided by MonkeyLearn [10]. The API takes the text as an input and returns a JSON object containing the opinion units within the given text. The data is then written into a JSON file.

D. Aspect extraction

At the completion of pre-processing and opinion extraction, the extraction of aspects is started. In order to extract the aspects out of the opinion units, as the initial step, an aspect classification model has been created using Support Vector Machine (SVM). A csv file consisting of a list of 200 pre-processed reviews on shampoo has been used to train the model. The training data has been tagged with four manually defined aspects consisting of effectiveness, price, fragrance and formula. The aspect classification model has been tested using both Support Vector Machine and Multinomial Naïve Bayes algorithms since most of the related work has been done utilizing the two algorithms. SVM showcased a higher accuracy level, thereby SVM being finalized as the classification algorithm for the model.

E. Polarity detection

Following the extraction of aspects, the opinion units are passed to the Polarity Detection module. The modules utilize two mechanisms in identifying the polarity of an opinion unit. Initially, the system uses Stanford CoreNLP library in order to identify the polarity of an opinion unit [11]. By taking the opinion unit as an input, Stanford CoreNLP predicts the polarity according to fine-grained sentiment analysis. In other words, it predicts the polarity as very positive, positive, neutral, negative and very negative.

The system then uses the output given by a custom sentiment analysis model created using Support Vector Machine (SVM). The sentiment analysis polarity classification model has been created beforehand, and it has been trained using the same review dataset which was used to train the aspect classification model. In order to train the model, each review sentence has been manually tagged as positive, neutral or negative. Once again, the model has been created using both SVM and Support Vector Machine and Multinomial Naïve Bayes algorithms and SVM has been opted as the finalized algorithm due to it having a higher accuracy level for the test data used.

The two predictions from Stanford CoreNLP and polarity classification model are compared and a final polarity value is decided. The reason for using two predictions is due to Stanford CoreNLP being a general sentiment analysis library, at times providing incorrect outputs when certain hair care domain-specific words are used. Therefore, The JSON file containing the reviews along with the opinion units is then updated to include the aspect as well as the polarity for each opinion unit. With the assistance of the updated JSON file, the results are graphically represented through a JSP page by using web technologies like Ajax, jQuery and Bootstrap.

IV. RESULTS

The system was evaluated by providing multiple eBay product pages having single page reviews as well as multiple page reviews as inputs. The system proves to be able to clearly display aspect-wise polarity for the selected aspects for different opinion units of the reviews on shampoos.

The intended outcome of the initial data gathering process was a CSV file containing the author, subject and content of each review written for a shampoo listed on eBay. The results demonstrate that the system is capable of extracting the aforementioned information not only for single-paged reviews but also for multiple-paged reviews by using the product page URL. Fig. 2. demonstrates a sample CSV file output generated by the data gathering module.

samtot1997	Nice	It's a nice cowash. My curls/coils are left soft, clean and defined. I can finger de
rose.nisbet	Brilliant!	I have very dry, very thick curly hair and I am always having trouble finding the
jomary-austin	Yummy	Love using this cleansing conditioner...it smells lush & although it doesn't lather
robysba40	Hair felt clean	I was worried that co-washing for the first time would leave my hair feeling un
barbados41	Perfect product to improve con	My hair has improved already! I am cutting down on conventional shampoo: t
super_12_vz	Perfect	Delivery was quick. I was surprised how much product you get, so great value f
clariclarinha	Washes well hair and doesn't d	It's perfect for co-wash. It really cleans the hair well, the fragrance is awesome
dorota8074	Just great!	I love it for my daughter's African hair
kisdr-25	very good cowash	Perfect for my wavy hair! My hair is clean and soft and shiny after using it.
mariaold3488	Love this sooooo much x	
laura.flawn	Consitions my hair beautifully!	
artist_painter	Good product	Good cowash - protein free for curly girls
looray75	Feels like a much more natural	Leaves my hair feeling clean, healthier and thicker - I have fine dry hair.
lidihiidr	Good for 4b/4c	Excellent product for cleansing my Afro hair. Cleanses well without drying, sm
starseed-uk	Perfect cleanser	Made my hair silky and shiny.
bim_zy	okay	okay

Fig. 2. CSV containing the reviews

Using the review data written to the CSV file, the pre-processing stage is capable of successfully identifying whether the review is written in English language, eliminating unnecessary special characters, converting needlessly capitalized words to simple and removing emojis and emoticons within the review. Fig. 3. shows how a sentence is pre-processed, completing the aforementioned tasks.

Before pre-processing : LOVE it! :) 😊
After pre-processing : Love it.

Fig. 3. A pre-processed review

The pre-processed data is then passed on for the extraction of opinions. As illustrated by Fig. 4, opinion units are separated within the sentence and data is written into a JSON file.

Reading the JSON file containing the opinion units for each review, aspects are identified using the aspect classification model. Then, the polarity for each opinion word is identified by combining the outputs of Stanford CoreNLP and polarity classification model created. Fig. 5. demonstrates the data written to the JSON file at the completion of aspect extraction and polarity detection.

```

"opinions": [
  {
    "opinion_unit": "I have very dry, very thick curly hair and i am always havin",
  },
  {
    "opinion_unit": "although it doesn't soap up in your hair like normal shampoo",
  },
  {
    "opinion_unit": "my curls are more bouncy and i don't need to put so much ser",
  },
  {
    "opinion_unit": "i will carry on using cleansing conditioners from now on.",
  },
  {
    "opinion_unit": "game changer.",
  },
  {
    "opinion_unit": "want to now try their conditioners and other products."
  }
]

```

Fig. 4. Opinion units

```

{
  "opinion_unit": {
    "classifications": [
      {
        "category_id": 122849811,
        "probability": 0.744,
        "confidence": 0.867,
        "label": "price"
      }
    ]
  },
  "text": "i was surprised how much product you get, so great value for money.",
  "polarity": "POSITIVE"
},
]

```

Fig. 5. Classified opinion units

Using the updated JSON file, the system is capable of graphically representing the results clearly as shown in Fig. 6. Tagged aspects are displayed for each review, denoted by the colours of green, yellow and red for positive, neutral and negative polarities respectively. The users have the ability to filter the reviews either by aspects or by aspect and the polarity.

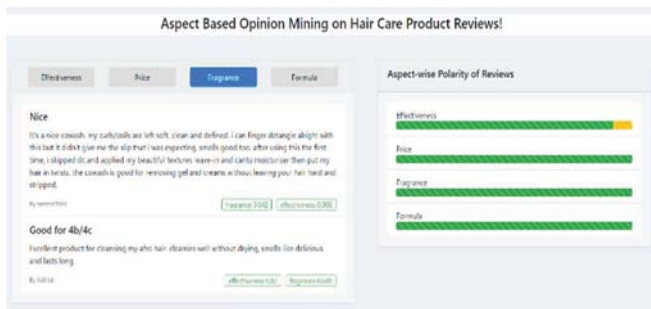


Fig. 6. Graphical representation of the output

V. DISCUSSION

The focus of the research was to decipher different public opinions on various aspects of hair shampoos expressed via reviews. Displaying the aspect-wise polarity was aimed since it is beneficial for both the manufacturers and consumers in gaining insights into consumer opinions and making purchase decisions respectively. Thus, the results of the system demonstrate that it is capable of dynamically displaying the aspect-wise polarity of reviews written on a given eBay product page.

The results received through test runs demonstrate that the system is capable of producing fairly accurate outputs for the given aspects. However, it takes a considerable amount of several seconds to produce the output due to the processes of web scraping, writing to and reading from files and the API calls involved as the number of reviews increases. Moreover, since it is a general sentiment analysis library, the Stanford CoreNLP library returns 'NEUTRAL' for certain opinion units like 'It makes my hair feel moisturized' whereas in an actual case, it contains a positive opinion. Therefore, the problem has been solved by integrating the polarity returned from the custom polarity classification model into the Stanford CoreNLP output. After the integration, the system proves to be able to accurately predict such a sentence as 'POSITIVE' as illustrated by Fig. 7.

Sentence : It makes my hair feel moisturized
Polarity from Stanford CoreNLP : NEUTRAL
Polarity from Classification Model : Positive

Fig. 7. Polarity prediction by CoreNLP and custom polarity classifier

However, the aspect classification model demonstrates to be incapable of correctly identifying sentences related to the aspect formula. As shown in Fig. 8., the phrase 'protein-free' refers to the aspect of the formula, but the classifier has only been able to identify the aspect of effectiveness denoted through the phrase 'Good cowash'.

```
"opinion_unit":{
  "classifications":{
    {
      "category_id":122849862,
      "probability":0.614,
      "confidence":0.958,
      "label":"effectiveness"
    }
  ],
  "text":"Good cowash - protein free for curly girls",
  "polarity":"POSITIVE"
}
```

Fig. 8. Aspect prediction by custom aspect classifier

VI. CONCLUSION

The main objective of the research was to decipher the aspect-wise polarity of the hair care product reviews so that, it can be useful for both hair product manufacturers and consumers. Manufacturers can identify their strong and weak points in order to make future decisions and change products to fit the requirement of the majority of the consumers.

Consumers can get an at a glance idea about the public opinion on a certain hair goal they intend to achieve by purchasing the product. In order to achieve the aforementioned objective, the initial step was to study and analyze related work and different mechanisms used in sentiment analysis. Using the knowledge gained through the literature review, methods for data collection, pre-processing, aspect extraction and polarity detection were derived. Using the HTML parser library JSOUP, reviews were dynamically collected by taking an eBay product page URL as the input. The reviews were pre-processed by removing unnecessary special characters, emoji and emoticons. Opinion units were broken down from review sentences and two classification models were created to classify the aspect and polarity respectively using SVM algorithm. After training the models, the aspect-wise polarity of the hair care product reviews was determined and graphically represented using a web page as the final output. The system can be further improved to support more aspects of a hair care product. The research is based on manually selected aspects for products falling under the category of shampoo. Hence, a mechanism of dynamically detecting the aspects would be ideal a future improvement for the system. Moreover, the accuracy of the system in terms of both aspect and polarity can be enhanced with the aid of more training data.

REFERENCES

- [1] S. M. Singh and N. Mishra, "Aspect based Opinion Mining for Mobile," *2nd International Conference on Next Generation Computing Technologies (NGCT-2016)*, Dehradun, India, 2016.
- [2] P. Ray and A. Chakrabarti, "Twitter Sentiment Analysis for Product Review Using Lexicon Method," *International Conference on Data Management, Analytics and Innovation (ICDMAI)* Zeal Education Society, Pune, India, Feb 24-26, 2017.
- [3] I. K. C. U. Perera and H. Caldera, "Aspect Based Opinion Mining on Restaurant Reviews," *2nd IEEE International Conference on Computational Intelligence and Applications*, 2017.
- [4] Z. Singla, S. Randhawa and S. Jain, "Statistical and Sentiment Analysis of Consumer Product Reviews," *8th ICCNT 2017*, IEEE - 40222, July 3-5, 2017, IIT Delhi, Delhi, India, Delhi, 2017.

- [5] V. Dhanalakshmi, B. Dhivya and A. Saravanan, "Opinion mining from student feedback data using supervised learning algorithms", *3rd MEC International Conference on Big Data and Smart City (ICBDSC)*, Muscat, Oman, 2016, pp. 1 - 5.
- [6] K. L. S. Kumar, J. Desai and J. Majumdar, "Opinion Mining and Sentiment Analysis on Online" *IEEE International Conference on Computational Intelligence and Computing Research (ICCIC)*, 2016.
- [7] F. Balahadia, M. Fernando and I. Juanatas, "Teacher's Performance Evaluation Tool Using Opinion Mining with Sentiment Analysis", *IEEE Region 10 Symposium (TENSYP)*, Bali, Indonesia, 2016, pp. 95 - 98.
- [8] J. Hedley, "jsoup Java HTML Parser, with best of DOM, CSS, and jquery", *Jsoup.org*, 2019. [Online]. Available: <https://jsoup.org/>. [Accessed: 12- Dec- 2019].
- [9] "vdmont/emoji-java", *GitHub*, 2019. [Online]. Available: <https://github.com/vdmont/emoji-java>. [Accessed: 12- Dec- 2019].
- [10] "MonkeyLearn", *App.monkeylearn.com*, 2019. [Online]. Available: https://app.monkeylearn.com/main/extractors/ex_N4aFcea3/. [Accessed: 12- Dec- 2019].
- [11] Manning, D. Christopher, M. Surdeanu, J. Bauer, J. Finkel, S. J. Bethard, and D. McClosky, "The Stanford CoreNLP Natural Language Processing Toolkit" in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2014, pp. 55-60