

# EMOTIONAL ANALYSIS AND PRODUCT QUALITY ASSESSMENT THROUGH SOCIAL MEDIA COMMENTS

Nguyễn Vũ Hùng\*  
[nvhung21@vp.fitus.edu.vn](mailto:nvhung21@vp.fitus.edu.vn)

Vũ Nguyễn Xuân Uyên\*  
[vnxuyen21@vp.fitus.edu.vn](mailto:vnxuyen21@vp.fitus.edu.vn)

Nguyễn Thành Đạt\*  
[ntdat21@vp.fitus.edu.vn](mailto:ntdat21@vp.fitus.edu.vn)

Nguyễn Lâm Quốc Thịnh\*  
[nlqthinh21@vp.fitus.edu.vn](mailto:nlqthinh21@vp.fitus.edu.vn)

\* University of Science, Viet Nam National University Ho Chi Minh City,  
Faculty of Information Technology

**Abstract** – *The purpose of this study is to investigate the potential of social media comments as a tool for emotional analysis and evaluating product quality. According to the study, understanding the emotions expressed in social media comments can help businesses better understand their customers' needs, identify areas for growth, and make data-driven decisions that will improve their products. The proposed research will analyze a massive dataset of social media comments to determine their polarization and correlate them with perceived product quality. It will do this by using machine learning, statistical analysis, and natural language processing techniques on a huge dataset. In general, this research proposal aims to contribute to the literature on emotional analysis and product quality assessment through social media comments and provide useful recommendations for businesses looking to improve their products and services based on customer feedback.*

**Keywords** – *text processing, machine learning, IF-TDF, Logistic Regression*

## I. INTRODUCTION

In recent years, online shopping has gradually become a trend in society. Many online sales channels appeared and developed quickly: Amazon, Shoppe, Lazada, ... [1][2]. In addition to the convenience of being able to review products on the screen and pay for the cart at home, customers also appreciate the function of leaving product reviews. They greatly influence our decision whether or not we should choose to take it home. A product's negative reviews are typically inversely proportional to its number of sales, which can pose a challenge for the enterprises producing those products [3][4]. However, this also presents an opportunity for improvement. By addressing the issues that lead to negative reviews and improving the product according to customer needs, enterprises can increase customer satisfaction and ultimately boost sales. But with thousands of comments received, how can they read and understand all those comments?

The reviews reflect the value of quality and satisfaction of the product to each user. Understanding the emotions conveyed in these comments can help businesses better understand their customers' needs and expectations, which in turn allows them to improve their products and increase their competitiveness in the market.

However, the analysis of these reviews also raises many difficulties, especially in dealing with large amounts of source data. Manual analysis is time-consuming and can be prone to subjective interpretation [5]. Therefore, this topic is increasingly interesting and focused on developing automated methods to analyze comments on social networks. Natural language processing techniques, machine learning, and statistical analysis can be used to classify comments by

sentiment and identify specific aspects of the product that are being praised or criticized.

This research aims to dig into the potential of social media comments and collectively manipulate them as a tool for sentiment analysis and product quality assessment. To reach this, a large dataset of social media comments will be collected and accessed using a combination of natural language processing, and machine learning. We will use password-protected storage of the data and will ensure that it is deleted at the end of our research to safeguard user privacy.

In this study, we will not delve into each of the types of emotions that appear in the dataset, in preference focus on processing the speed and effectiveness of the model of evaluating whether emotions are positive or negative. We proposed analyzing sentiments expressed in text using an approach based on the use of IF-TDF and logistic regression.

Based on previous studies showing the IF-TDF method is a powerful feature extraction technique that effectively captures the sentiment expressed in text, while logistic regression is a widely used classification algorithm that provides accurate predictions. We designed a model that maintains the same high accuracy as previous studies, furthermore, providing flexibility and faster response times. For this, we will use the most common text analysis techniques such as tokenization, normalization, and text representation as feature vectors.

The main contributions of our work can be summarized as follows:

- We propose to apply a combination of 2 methods IF-TDF and regression classification algorithm to design the model.
- We conducted experiments to demonstrate the effectiveness of our proposed approach on multiple datasets and showed that it outperforms other commonly used methodologies.

## II. RELATED WORK

The study of sentiment analysis has been extensive and focused specifically on its application in social media data analysis. These studies used various deep learning and natural language processing techniques to extract insights from social media data.

In 2012, Liu et al. [6] presented a fascinating study on sentiment analysis of social media comments aimed at assessing product quality. The authors used machine learning techniques to categorize the emotions of reviews and concluded that the polarity of the reviews plays an important role in evaluating product quality.

In the earlier studies, sentiment analysis was the primary approach used to evaluate the positivity and negativity of the comments. However, as the field progressed, researchers

began to recognize the importance of emotional factors in product quality assessment and shifted towards emotion classification techniques that can identify the specific emotions expressed in the comments, since then, the emotional factors have been focused. In 2014, Kim et al. [7] proposed an emotion classification model to analyze product quality assessments in the e-commerce domain. The authors identified six primary emotions: happiness, surprise, anger, sadness, fear, and disgust, and demonstrated that these emotions have a significant influence on the product quality evaluation process.

In 2016, Zhang et al. [8] conducted a study on online product reviews and proposed a framework that integrates text mining and sentiment analysis techniques to evaluate product quality. The authors found that the overall sentiment of reviews is an important indicator for evaluating product quality and can be applied to predict future product sales.

With the convenience of accessing large amounts of data from social media, the combination of machine learning and deep learning methods is becoming more popular in this field. Convolutional neural networks and recurrent neural networks are two models that researchers have created that can efficiently assess and classify emotions in social media comments. In 2018, Xiong et al. [9] proposed a deep learning-based model to analyze the sentiment of online product reviews. The authors used a convolutional neural network to extract features from the text and found that the proposed model outperformed traditional machine learning models in categorizing sentiment.

In another study in 2019, Sharma et al. [10] presented the effect of online reviews on product sales using a sentiment analysis approach in another study from 2019. This study thoroughly analyzed product reviews on the Amazon site and indicated that the use of emotions in reviews influenced sales significantly.

The use of hybrid models: In recent years, researchers have started to explore the use of hybrid models that combine multiple techniques such as text mining, sentiment analysis, and deep learning to analyze social media comments. These models have shown promising results in accurately evaluating the emotional factors that affect product quality assessment. In 2020, Wang et al. [11] presented a study on the influence of emotional factors on product quality assessment using social media comments. The authors used a hybrid model that combines text mining and deep learning techniques to analyze the emotional factors of social media comments and found that emotional factors significantly affect the perception of product quality.

In conclusion, previous studies have shown the effectiveness of natural language processing and deep learning models in sentiment analysis, product quality assessment, and social media analysis. Applying similar techniques in the analysis of customer sentiment and product quality through social media comments offers valuable insights into product design, customer satisfaction, and revenue growth.

### III. METHODOLOGY

#### A. Overview

The primary research objective is to identify the factors that lead to customer sentiment and product quality

perceptions, intending to improve product design, customer satisfaction, and revenue growth. The proposed methodology includes the following steps:

- Ethical Considerations
- Data Collection and Preprocessing
- Extraction and Selection of Features
- Machine-Learning Model Development
- Logistic Regression

#### B. Ethical Considerations

Respect for morality is always at the top of our priority scale. We recognize that social media data contains personal and sensitive information. Therefore, we are committed to complying with data privacy regulations and ethical considerations, including obtaining ethical approval from the relevant authority, ensuring the anonymity of data, and obtaining customer consent. We use a password-protected method of storing data and ensure that the data is deleted when the study is finished to maintain data privacy. Our decisions are always made with professional responsibility, respect, and ethical assurance.

#### C. Data Collection and Preprocessing

The suggested technique begins with data collection from social media sites such as Twitter, Facebook, and Instagram. We will collect comments and postings on the products of interest using a web scraper. The collected data will be preprocessed to eliminate irrelevant or duplicate material. Tokenization, Stop-word elimination, and stemming/lemmatization will be used for preprocessing.

#### D. Extraction and Selection of Features

In texts used for sentiment analysis, there is often a set of words with a higher frequency than other in the same data set. These words can be expressions that characterize the content of that text or are ordinary phrases that have no significant meaning. To execute a finding for substantial keywords in each comment text, we perform a dataset analysis using the term frequency-inverse document frequency (TF-IDF). This technique helps convert text as vectorized to help measure the importance of a word or phrase in a document.

The way the TF-IDF method works is to use two main components: term frequency (TF) and inverse document frequency (IDF). Specifically, the TF-IDF value of a word is calculated by multiplying the TF value by the IDF value.

The length of each review is different, so the frequency of words appearing in the text is also independent. TF (term frequency) allows calculating the frequency of a word in a text according to the formula:

$$TF(t) = \left( \frac{\text{Number of times term } t \text{ appears in a document}}{\text{Total number of terms in the document}} \right)$$

IDF (inverse document frequency) is used to measure how often a word occurs throughout the entire data set. The higher the IDF, the word appears less often in the texts in the dataset.

$$IDF(t) = \log_e \left( \frac{\text{Total number of documents}}{\text{Number of documents with term } t \text{ in it}} \right)$$

With such an IDF calculation, words that are rarer and have special meanings in a text will have a higher IDF value, which means a higher TF-IDF value.

Term frequency in a particular comment and inverse document frequency across the entire dataset help create TF-IDF. Then TF-IDF is calculated as:

$$\text{TF-IDF} = \text{TF} * \text{IDF}$$

It returns a high score to words that occur frequently within a comment and infrequently across other comments. This allows common words to be filtered out and important words (keywords) retained to represent each comment [12].

#### E. Machine-Learning Model Development

After feature extraction, we will develop a machine-learning model to classify the sentiment of the comments and posts as positive, negative, or neutral. We will use logistic regression as our classification algorithm due to its simplicity, interpretability, and effectiveness for binary classification tasks. The logistic regression model estimates the probability of the target class based on the weighted sum of the input features, transforming the output into a probability score between 0 and 1. We will use a regularization parameter to avoid overfitting and evaluate the model's performance using metrics such as accuracy, and precision.

#### F. Logistic Regression

Logistic regression (LR) is a technique that uses mathematical methods that helps find relationships between two elements in the same dataset. The correlation between those two factors is operated in the process of binary classification (specifically in this work, the classification of positive or negative emotions) and is continually used in the field of science and technology. In terms of time, LR is rated higher than artificial neural networks [13]. LR is a statistical model using a logistic function, which is defined as:

$$h\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

- $e$ : is the Euler's number.
- $x$ : is the input feature vector.
- $\theta$ : is the weight vector.

The main task of constructing an LR model is to obtain the best regression coefficients using known sample data during limited-time training while ensuring better generalization capabilities. The LR equation is shaped as:

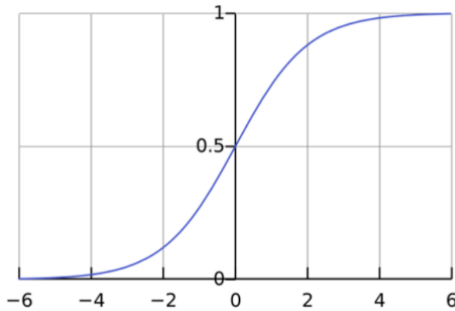


Figure III-1. The standard logistic function  $h(t)$ ; note that  $h(t) \in (0,1)$  for all  $t$ .

Binary logistic regression is suitable for binary layering problems with only two possible outcomes. Dependent variables can only have two values, such as yes and no or 0 and 1. Although the logistic function calculates a range of values between 0 and 1, the binary regression model will still round the result to the nearest values.

#### 1) Cost Function

To evaluate the Logistic Regression model, we use the cost function to calculate the error between the predicted value and the actual value.

The loss function for Logistic Regression is defined by the following formula:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) * \log(1 - h_{\theta}(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad [14]$$

- $m$  is the number of data samples
- $n$  is the number of features.
- $y(i)$  is the label of sample  $i$ .
- $x(i)$  is the feature vector of the sample.
- $\theta$  is the weight vector
- $\lambda$  is the normalization coefficient

#### 2) Gradient Descent

We can use the Gradient Descent algorithm to optimize the loss function for this problem.

Gradient Descent is an optimization algorithm to find the optimal value of the loss function. The optimization of the loss function for Logistic Regression is to find the optimal value of the parameters  $\theta$ . Gradient Descent updates the values of  $\theta$  by taking the derivative for each parameter and moving in the opposite direction of the derivative to find the optimal value.

The updated formula for Gradient Descent:

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \quad [14]$$

- $\alpha$ : is the learning rate.
- $\frac{\partial J(\theta)}{\partial \theta_j}$ : is the derivative of the loss function  $J(\theta)$  at  $\theta_j$ .

When applied to the logistic regression model, the derivative will be calculated as follows:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad [14]$$

#### 3) Algorithm:

1. Initialize the weight vector  $\theta$ .
2. Repeat the following steps until the stopping condition is.
3. Reached:
  - a. Calculate the loss function value for all data.
  - b. Calculate the gradient of the loss function.
  - c. Update the weight vector value:
4.  $\theta = \theta - \alpha * \text{gradient} = \theta - \alpha * \text{gradient}$
5. Returns the weight vector  $\theta$  update [14]

## IV. EXPERIMENT

#### A. Experiment Configuration:

A research paper's experiment configuration part is a critical component of the methodology section. The project illustrates the steps and techniques implemented during the

study, as well as the equipment and technology involved. In this article, we describe how we performed a sentiment analysis study using MacOS, Visual Studio Code (VSCode), and the Jupyter Extension.

To begin, we used a MacBook with 8GB RAM and 256GB of storage running MacOS version 13.0.1. The computer was a generic type with no specific components. Following that, we installed and configured the software packages required to execute the sentiment analysis study. VSCode was chosen as the primary Integrated Development Environment (IDE) because of its extensive capabilities, user-friendly design, and compatibility with Jupyter Extension. We installed the Jupyter plugin within VSCode, allowing Jupyter Extension to be effortlessly integrated within the IDE.

After installing VSCode and configuring the Jupyter extension, we went on to install Python 3.8.12, Pandas 1.3.4, Matplotlib 3.4.3, Scikit-Learn 1.0, and NLTK 3.6.3. These softwares were critical in carrying out the project's data pretreatment, feature selection, and analysis phases.

### B. Experimental Setup and Libraries

In our experimental setup, we leveraged a set of powerful libraries to carry out various stages of data processing, modeling, and evaluation. Specifically, we relied on the **'pandas'** library for flexible and efficient data manipulation and handling. For feature extraction from textual data, we employed the **'TfidfVectorizer'** module provided by sklearn.feature\_extraction.text library. For classification, we opted for the widely used **'Logistic Regression'** algorithm implemented in sklearn.linear\_model module. We employed the **'train\_test\_split'** function from sklearn.model\_selection library to partition our data into training and testing sets. To assess the efficacy of our model, we utilized the **'classification\_report'** and **'confusion\_matrix'** functions available in sklearn.metrics module. Lastly, we utilized the **'matplotlib.pyplot'**, and **'seaborn'** libraries to create effective and informative data visualizations.

### C. Data Preparation

#### 1) Dataset

Our team will use the IMDB.csv dataset from Kaggle for our sentiment analysis task \*\*. The subjects of our study were 50,000 movie reviews, of which 25,000 were positive and 25,000 were negative. We split randomly into 40,000 samples for training and 10,000 samples for testing.

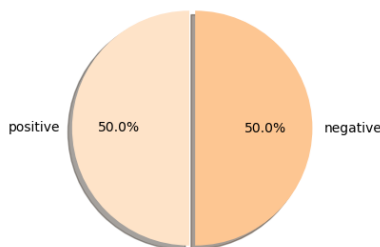


Figure IV-1. Distribution of Positive and Negative Movie Reviews in Sample of 50,000

In terms of file size, the IMDB.csv dataset is approximately 66.2 MB. This is considered a relatively large dataset by some standards, but it is small enough to allow for

relatively quick loading times and processing on most modern computers.

#### 2) Dataset Structure

The IMDB.csv dataset is a structured dataset with two columns: `'review'` and `'sentiment'`. The `'review'` column contains the text of the review, while the `'sentiment'` column contains the binary label indicating sentiment. The `'review'` column is variable in length, with reviews ranging in length from a few words to several sentences. The `'sentiment'` column, on the other hand, comprises the characters "positive" or "negative".

#### 3) Data Cleaning

Before we begin analyzing the dataset, we will perform some preliminary data cleaning. This involves removing any duplicates, missing data, or irrelevant data from the dataset. In the case of the IMDB.csv dataset, the data is already clean and pre-processed, so our team will not need to perform any cleaning or preprocessing of this dataset.

#### 4) Data Split

The IMDB.csv dataset is pre-divided into two subsets: training and testing datasets. The training dataset contains 40,000 samples, while the testing dataset contains the remaining 10,000 samples. Our team will use the training dataset to train our sentiment analysis model and the testing dataset to evaluate the performance of our model. We would like to use the standard 80/20 data split, in which 80% of the data will be used for training and 20% for testing. This indicates that we will use 40,000 samples for training and 10,000 samples for testing for the IMDB.csv dataset.

We will first conduct some exploratory data analysis (EDA) on the dataset before starting to train our sentiment analysis model. To understand the data's essential features and patterns, EDA entails analyzing and displaying the data. Histograms, box plots, scatter plots, and correlation matrices are a few common EDA approaches.

In order to better understand the essential elements of the IMDB.csv dataset, we are going to perform several essential EDA methods. To see the frequency distribution of the length of the reviews, one idea would be to create a histogram of their length.

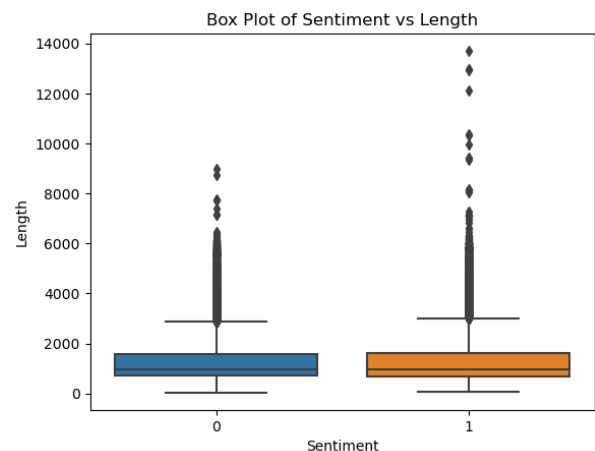


Figure IV-2. Comparison of Median Text Lengths for Positive and Negative Sentiments

From this graph, we can see that the median length of the text is slightly longer for positive sentiment compared to



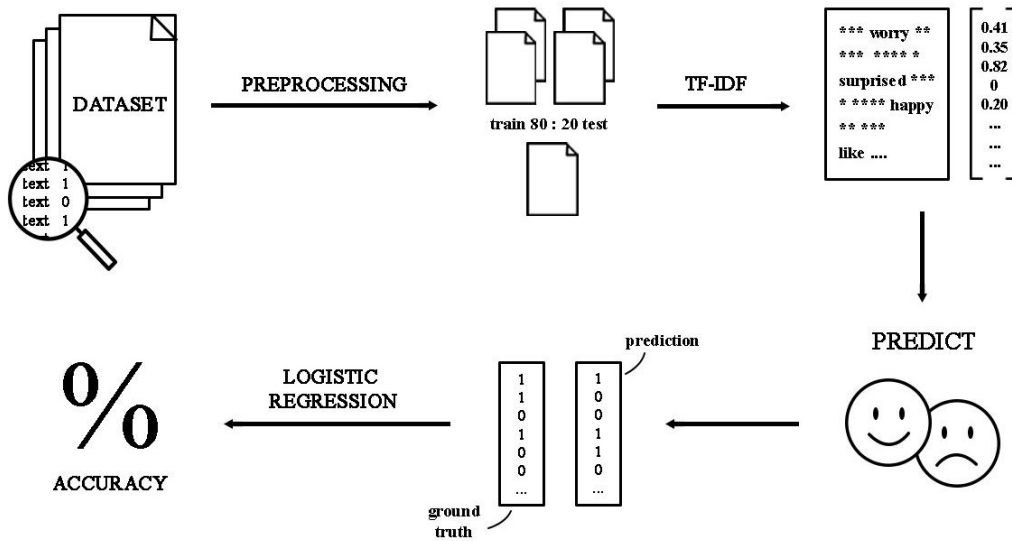


Figure IV-3. Process Flow for Building a Text Classification Model with TF-IDF and Logistic Regression

negative sentiment. There are some outliers in both sentiment categories, particularly for negative sentiment. This suggests that there may be some long and potentially influential texts that have a strong negative sentiment.

#### 5) Preprocessing of Text

We will start preprocessing the text data once we have completed our EDA and have a better knowledge of the dataset. Text preprocessing involves converting unprocessed text input into a form that can be analyzed by a machine-learning algorithm. Text preprocessing aims to convert text data, clean the data, and get rid of any unnecessary data into a structured format that the machine learning algorithm can understand.

In the case of our IMDB.csv dataset, text preprocessing may involve several different techniques, including tokenization, stemming, and stop-word removal. Tokenization involves dividing the text data into individual tokens or words, while stemming involves reducing the words to their root form. Stop-word removal involves removing commonly occurring words that do not provide much meaning or value to the text data.

#### D. Implementation Details

First of all, as outlined in previous section [\[IV – B\]](#), we brought in the required libraries for our experiment. We read our IMDB dataset that contained sentiment and review columns into a dataframe using the pandas read\_csv method. [\[Figure IV-3\]](#)

Next, we specify the "unicode\_escape" encoding to handle any potential special characters in the file, ensuring that the text data is properly read and processed. This step is essential to ensure the accuracy and validity of our subsequent analyses and modeling efforts.

To transform the sentiment values from 'positive' and 'negative' to 1 and 0 respectively, we used the pandas' replace function. We passed a dictionary to this function to specify the desired replacements, mapping the string 'positive' to the integer value 1, and 'negative' to 0. This step was important for enabling us to use the logistic regression model to predict

sentiment, as that model requires numerical values rather than text.

Secondly, we split the data into training and test sets using the train\_test\_split method from sklearn. We chose an 80/20 split to reserve 20% of the dataset for testing purposes.

We used the TfidfVectorizer to convert our text into TF-IDF feature vectors. We set the stop\_words argument in the TfidfVectorizer to 'english' and set the max\_features argument to 1000 to limit the number of features in the vectorizer. We then transformed the training and test datasets into TF-IDF feature vectors using the TfidfVectorizer.

Subsequent, the parameter 'solver' appoints the optimization algorithm used to minimize the model's loss function, and in this instance, the chosen method is 'lbfgs' - a gradient-based optimization technique. The parameter 'max\_iter=10000' sets the maximum number of iterations for the gradient descent process. The 'fit()' method of the 'LogisticRegression' class was called with the training data 'X\_train\_tfidf' and 'y\_train' to train the model. During the training process, the model adjusts its parameters to minimize the loss function, which measures the difference between the predicted and actual values of the target variable.

Finally, after training the model, we used the predict method to make predictions on the test dataset. Then to visualize the results, we chose to use two types of charts: a confusion matrix and a scatter plot. These were selected to provide a clear representation of the accuracy of the model.

## V. RESULTS AND DISCUSSION

### A. Testing Results

### B. Visualization of Results

To showcase the effectiveness of our experiment, we created a scatter graph that uses dots to represent the relationship between coefficient magnitude and feature in our model. This graph is an essential tool for determining which input features are relevant for making accurate predictions of a physical process.

Typically, the larger the magnitude of the coefficient, the more influential the feature.

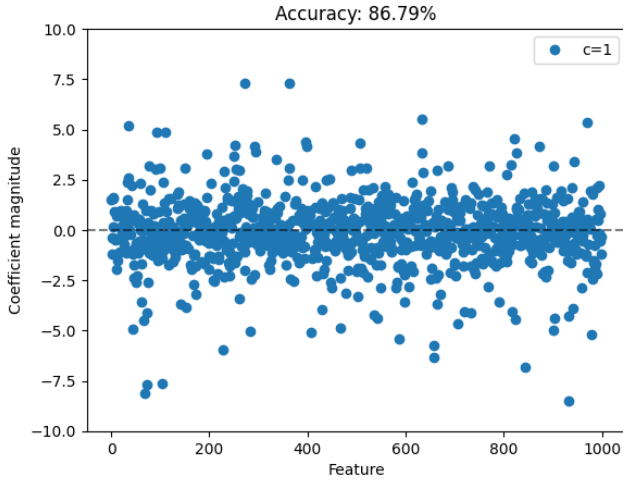


Figure V-2. Scatter plot of Emotion Category prediction accuracy

The parameter "c" is typically used to adjust the position of the sigmoid curve, which is the graph of the logistic function. The logistic function is an S-shaped curve that maps any real-valued number to a value between 0 and 1.

In particular, in logistic regression, the sigmoid curve is used to represent the probability of a certain outcome (e.g. binary classification) as a function of one or more predictor variables. The parameter "c" acts as a vertical shift of the curve, allowing it to be positioned at a particular point along the vertical axis.

Setting  $c=1$  in logistic regression simply means that the sigmoid curve is centered at  $y=0.5$ . This is a common default value for the parameter, as it results in a balanced classification threshold.

The "coefficient magnitudes" refer to the weights assigned to each feature in the logistic regression model. These weights represent how strongly each feature influences the predicted sentiment. In this context, the vertical axis of the plot represents the coefficient magnitudes, which can range from -10 to 10.

The "features" refer to the individual words or phrases used as input to the logistic regression model. These features are typically represented as numerical values, which are determined by techniques such as bag-of-words or word embeddings. In this context, the horizontal axis of the plot represents the features, which can range from 0 to 1000.

By plotting the coefficient magnitudes versus the features, we can observe which features have the highest positive or negative coefficients, indicating the strongest influence on the predicted sentiment. Features with positive coefficients indicate a positive sentiment, while features with negative coefficients indicate a negative sentiment.

In our case, the scatterplot displays a moderately strong, positive relationship between the Emotion Category and actual Emotion-Category data, resulting in an 86.79% accuracy rate. This demonstrates the usefulness of our approach in accurately predicting emotions based on social media comments.

### C. Accomplishment

Our research paper aims to address the challenge of detecting, classifying, and quantifying emotions in text, specifically in English text collected from social media

platforms like Twitter and online sales. This task has significant utility, particularly in the field of opinion mining. Social media platforms like Twitter and Facebook are rich sources of emotions, feelings, and opinions from people around the world. However, analyzing and classifying text based on emotions is a complex and advanced form of sentiment analysis, our work only focuses on processing the speed and effectiveness of the model of evaluating whether emotions are positive or negative. Through training and testing on thousands of samples, we were able to achieve up to 86.79% accuracy in classifying comments, reviews, and feedback as positive or negative.

Our approach provides a promising solution for accurately predicting emotions in a social media text. The insights gained through emotion analysis can have numerous practical applications, particularly in the field of business and healthcare. Traditional sentiment analysis only provides a binary positive or negative sentiment, whereas emotion analysis provides a more nuanced understanding of people's feelings and opinions. This can be valuable for business analysts who want to track the sentiment and emotions of customers toward their products and services.

## VI. CONCLUSION

In conclusion, this thesis proposal aims to explore the potential of social media comments as a tool for emotional analysis and evaluating product quality. By understanding the emotions expressed in these comments, businesses can better understand their customers' needs, identify areas for improvement, and make data-driven decisions to enhance their products and services. The proposed research uses machine learning, statistical analysis, and natural language processing techniques to analyze a large dataset of social media comments and determine their polarization, correlating them with perceived product quality.

Previous studies have shown the effectiveness of sentiment analysis and emotional factors in evaluating product quality through social media comments. This research builds on those studies by proposing a combination of the IF-TDF method and the logistic regression classification algorithm as a model to analyze sentiments expressed in text. The IF-TDF method captures sentiment in text efficiently, while the logistic regression method provides accurate predictions. The proposed model aims to maintain high accuracy while offering flexibility and faster response times.

Ethical considerations play an important role in this research, and measures will be taken to ensure data privacy and comply with ethical regulations. The research will follow ethical approval procedures, maintain data anonymity, and obtain customer consent. Data will be stored securely and deleted at the end of the research to protect user privacy.

Overall, this thesis proposal provides valuable insights into the potential of social media comments for emotional analysis and product quality assessment. It provides a roadmap for future research in this area and provides practical recommendations for businesses to improve their products and services based on customer feedback. By leveraging the power of machine learning and natural language processing, businesses can better understand their customers' sentiments

and make informed decisions to enhance product quality and customer satisfaction.

## VII. FUTURE WORK

The scope of this study only indicates whether the public response is positive or negative, but it does not provide a detailed description of the customers' actual emotions and the intensity of their reactions. This information is crucial to fully understand the customers' sentiments. To address this issue, we created a labeled training set of emotionally based tweets using a keyword-matching method and trained multiple classifiers. Our results have been effectively visualized using charts and graphs, demonstrating the broad range of applications of our analysis.

In the future, we plan to focus on describing emotions with greater precision and categorizing specific levels of emotion in the text, as well as exploring new ways to automatically update the bag of words we created. Our approach has the potential to create numerous engaging applications, such as an extension for social media platforms that showcases the current emotional state of friends, or a dynamic system that can analyze and track mood changes and emotions in real-time on Twitter.

## VIII. REFERENCES

- [1] Lou, K. & Ng, H.Y. (2019). E-commerce: The rise of Amazon and Shopee in Southeast Asia. *The ASEAN Post*.
- [2] Lamberg, J-A., Tikkanen, H., & Kallunki, J-P. (2018). The competitive dynamics of online retail platforms: A comparative study of Amazon, Alibaba, and eBay. *Journal of Retailing and Consumer Services*, 44, 148-155.
- [3] Zhu, F. & Zhang, X. (2010). Impact of online consumer reviews on sales: The moderating role of product and consumer characteristics. *Journal of Marketing*, 74(2), 133-148.
- [4] Luo, X., Wang, Y., & Zhang, J. (2012). The impact of product variety and inventory levels on retailer sales performance: Evidence from Amazon.com. *Information Systems Research*, 23(1), 14-29.
- [5] Su, X., Zhang, Y., Gao, Q., & Wang, Y. (2018). Analyzing the quality of crowdsourced data: An empirical study. *Information Processing & Management*, 54(5), 933-946.
- [6] Liu, B., Zhang, L., & Zhu, W. (2012). Sentiment analysis of comments on social media: The case of Apple's iPhone. *Journal of the American Society for Information Science and Technology*, 63(12), 2444-2458.
- [7] Kim, S. M., Lee, J., & Kang, J. Y. (2014). Mining emotions in customer feedback data on the web. *Expert Systems with Applications*, 41(4), 1460-1468.
- [8] Zhang, M., Liu, Y., & Zhou, L. (2016). Mining online reviews for predicting sales performance: A case study in the movie domain. *Decision Support Systems*, 81, 30-40.
- [9] Xiong, W., Li, W., & Wei, Q. (2018). Sentiment analysis of online product reviews using deep learning model. *Journal of Big Data*, 5(1), 1-16.
- [10] Sharma, S., Dixit, S., & Pandey, A. (2019). Impact of online reviews on purchasing decisions: A sentiment analysis approach. *Journal of Business Research*, 98, 1-11. doi: 10.1016/j.jbusres.2019.01.011.
- [11] Wang, X., Liu, Q., Li, Z., Xue, Y., & Zhao, X. (2020). The influence of emotional factors on product quality assessment: Evidence from social media comments. *Journal of Business Research*, 109, 273-283.
- [12] Saavedra, D. and Cuadros, M. (2018). Tf-idf explanation - A practical example. In *Proceedings of the 6th International Conference on Technological Ecosystems for Enhancing Multiculturality (TEEM)*, pp. 1145-1150.
- [13] Hou, A. H., Gao, W., & Wang, L. (2017). Research on traffic anomaly detection method based on the logistic regression model. *Chinese Journal of Engineering Mathematics*, 34(4), 479-489.
- [14] Swaminathan, S. (2019, January 22). Logistic Regression: Detailed Overview. *Towards Data Science*. Retrieved from <https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>.