

Classify Online Customer Reviews from Different Restaurants Based on Sentiment Analysis

Brian Zhou^{1*}, Qingyang Yang², Lingye Kong³

Corresponding author. Email: baz896@usask.ca

¹Arts and Science University of Saskatchewan, Saskatchewan, S7N5B2, Canada, baz896@usask.ca

²Pierce College, Tacoma, 98498, US, leoyyy06@outlook.com

³Pierce College, Tacoma, 98498, US, linyek623@gmail.com

Abstract. It is nearly difficult to manually analyze the data in sentiment analysis since the volume is growing quickly every day. People frequently express their opinions through social media platforms like Tripadvisor and restaurant review sites. Several machine-learning approaches might be used to automate the process of text analysis. Enormous amounts of textual data are produced, yet they are worthless without processing. This paper leverages sentiment analysis to classify customer reviews into three categories: positive, neutral and negative. We used many different classifiers, including Logistic Regression SVM, Random Forest, and Gaussian Naive Bayes, to sort the feedback into appropriate categories. The Random Forest and SVM turned out to have the highest and similar accuracy for the prediction of text data. Logistic regression has a very close but relatively lower accuracy in comparison, and Gaussian Naive Bayes yielded the lowest accuracy. In the end, we use a majority voting ensemble classifier which combines three classifier predictions with the most accuracy. We assessed their performance using a cross-validation technique, comparing their accuracy, precision, recall, and F1-score on the test set. The Majority voting turned out not to have a significant increase in accuracy.

Keywords: Sentiment analysis, Social media, Restaurant reviews, Data Processing, Different classifiers.

1 Introduction

The internet allows users from all across the world to express their opinions on social media by collecting and analyzing online opinions on social media about a product or service. It is possible to predict the overall public opinion regarding said product or service. This information can form a dataset that can help be leveraged as a restaurant recommendation engine through machine learning techniques. Additionally, NLP (natural language processing) techniques may be used. Review data is processed and analyzed using natural language processing techniques, including topic modeling, text mining, and sentiment analysis. Through the use of the NLP technique, one can comprehend the customers' opinion of the restaurant and their emotional inclinations and concerns. Virmani and Pillai used NLP to extract information from social networks [1]. During their research, they go over various text mining systems to analyze social network information. A method for improving the precision of sentiment analysis is shown by Tetsuya and Jeonghee [2]. Using a structural parser and sentiment lexicon, their system could

perform sentiment analysis on web pages and news stories with a high degree of accuracy (about 75% to 95%).

Many other research papers also explored sentiment extraction through social media for financial frameworks. For example, Houlihan and Creamer used social media sentiment to predict stock price direction moves [3]. By harnessing the power of NLP and machine learning, researchers have been able to extract valuable insights from social media and leverage them to provide some personalized recommendations to customers (few citations). Sehgal and Song used sentiment measures extracted from financial blogs to determine correlations between sentiment and stock prices and used sentiment to forecast asset prices [4]. They concluded that sentiment can be predicted with a high degree of model performance. Additionally, they were able to predict asset prices with up to 81% accuracy.

An example that can be cited is the Support Vector Machine (SVM). The technique of SVM is commonly employed to tackle problems related to classification and regression. Obiedat et al. proposed the replacement of the widest margin classifier, which is only capable of handling linearly separable data, with the support vector classifier [5]. This upgrade was introduced to enhance its capabilities. Another example is using Naive Bayes, like Rachmawan et al. did in their article [6].

1.1 Sentiment Analysis in Politics

Sentiment analysis is used not only in product reviews but also in news stories, political discussions, and the stock market. For instance, during political discussions, we may learn what individuals thought about particular political parties or candidates for office. Political positions may also be used to forecast election outcomes. Because so many individuals openly share and debate their thoughts on social media and microblogging sites, they are considered excellent sources of knowledge [7].

1.2 Making Investment Decisions using Social Media & Sentiment Analysis

Sentiment analysis entails examining certain attitudes or sentiments that investors have toward the markets. The best moment to use sentiment analysis as a contrarian signal is when price fluctuations have turned excessively bullish or bearish. Because they base their purchases and sales on the market's overall direction or the actions of a specific firm, retail traders frequently follow trends reactively. According to Kulka (2022), engaging in sentiment-based transactions against high emotional biases could potentially yield profitable outcomes [8]. For several of the chosen assets, we discovered that the trend in SM sentiment followed a pattern that was comparable to the stock market performance. A user might use the sentiment to guide investing decisions based on these findings [9].

1.3 Reputation and Sentiment Analysis

The collective opinion of the internet community on a specific thing is referred to as the item's online reputation. The reputation of a product, which is formed from the opinions of others, is significant information when customers must decide or choose a product. The results of sentiment analysis may be used to gauge customer opinion. Thus, providing guidance and recommendations in the selection of items in accordance with the wisdom of the community is the first use of sentiment analysis. A single overall rating may be misleading. Sentiment analysis

can aggregate reviewer comments and estimate ratings for particular product attributes [10]. According to Manning et al., we have suggested that we outline the architecture and application of the Stanford CoreNLP toolbox, an extendable pipeline that offers fundamental natural language processing [11]. The open-source NLP community for research as well as users in business and government utilize this toolbox frequently. We contend that this results from a primary interface, a solid and high-quality analytical component, a simple, accessible design, and the absence of a significant quantity of accompanying baggage.

1.4 Sentiment analysis for measuring customer satisfaction

According to Al-Otaibi et al., we can determine that social media analytics is the practice of collecting information from social media sites and using that information to conduct analysis and make business choices [12]. To enhance marketing and customer care efforts, social media analytics are most frequently used to harvest customer opinions. In order for businesses to make wise business decisions, it is helpful to understand what consumers think and what they have to say about the value of the goods and services. In this post, we exploited Twitter data to uncover popular sentiment that was concealed in the data. The unigram is employed as a feature extraction technique, and the SVM algorithm is applied to categorize tweet sentiment as positive or negative. The algorithm obtained excellent accuracy of over 87% in the trials utilizing a huge collection of training datasets. The use of the unigram as a technique for feature extraction. The algorithm obtained excellent accuracy of over 87% in the trials utilizing a massive collection of training datasets.

1.5 Challenges of NLP

It is challenging to create a program that can comprehend natural language. Large and including an unlimited number of sentences, most natural languages are. Furthermore, natural language has a lot of ambiguity. In various contexts, identical phrases often possess divergent implications, while numerous words, such as can, bear, fly, and orange, exhibit multiple meanings. Consequently, developing software that can comprehend natural language poses significant challenges. We may make decisions about how words are used to create deeper meanings by looking at a language's syntax. Chowdhary uses the example "the seller bought the merchant a dog" to emphasize the need to be explicit about what is being sold to whom [13].

1.6 Ethics of NLP

We do need to consider the ethics of NLP. Bias is among the most important ethical issues surrounding NLP. Large datasets are used to train NLP models, and the caliber of the data they use determines the caliber of their output. The NLP model may pick up and propagate any bias present in the training data, which might produce biased or unjust results. For instance, even inadvertently, a hiring process based on NLP can bias against individuals based on their ethnicity or gender [14]. Though the problem of prejudices in NLP might be deceptively subtle, going unnoticed while yet being deadly, it is a problem that can be resolved. Self-empowering and autonomous technology must be handled carefully and subject to ongoing ethical scrutiny. The utility of technology and its propensity for sneaky invasions of people's privacy, which is still one of their fundamental rights, must be constantly discussed. A level of compromise is frequently used to resolve the ethical seesaw between 100% privacy and 100% security [15].

2 Data Collection

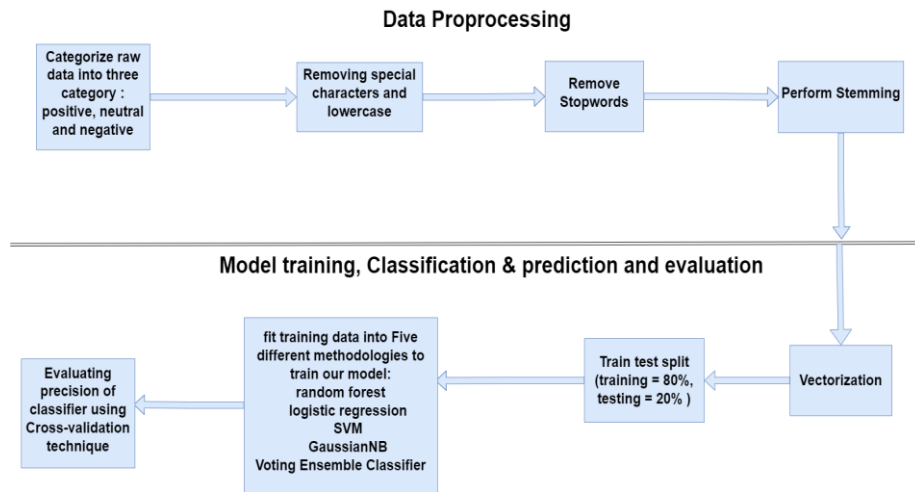
The dataset was retrieved from Kaggle, which contains restaurant reviews across 20 US cities in Washington, Texas, California, New York and New Jersey. Kaggle serves as a platform that facilitates data modeling and analysis competitions. Data from various sources, including businesses, researchers, and statisticians, can be shared on the platform, fostering a competitive environment where data mining and statistics professionals strive to develop the most optimal models. All reviews in the Kaggle dataset are scraped from the TripAdvisor website using BeautifulSoup, a popular Python API used to scrape data off web pages. BeautifulSoup is a Python library that facilitates the extraction of specific content from webpages while also enabling the removal of HTML markup and subsequent preservation of the extracted information. The web scraping tool is designed to make it easier to collect and arrange data from online-based documents. TripAdvisor is an internet-based platform that serves as a comprehensive resource for travel-related information and reservations. It offers a wide range of user-generated content, including reviews, photographs, and forums, which provide valuable insights and guidance pertaining to numerous hotels and resorts worldwide. After repeatedly cleaning and deleting repetitive data using Python in Spyder software, we were shocked to find 2636 restaurants in five important US states with various names in this dataset, which also has 3060 reviews in Table 1.

Table 1. Statistics about the data, fields, field description, unique value.

Field	Description	Unique
Name	Name of the restaurant	2636
Location	Physical Location	2653
Type	Type of cuisine	78
Reviews	Star rating between 0-5.0	3060
Comments	Comments from customers	2292

3 Methodology

The objective is to develop models that can effectively categorize reviews into three distinct classes: positive, neutral, and negative. We will perform the following steps in Figure 1.



Note: This illustrates all of the steps taken for the methodology.

Fig. 1. Flow diagram.

Step 1- Data Preprocessing. First of all, we categorize comments into three categories (sentiment labels): positive, neutral and negative, according to the score of review, which is contained in the raw data. Reviews with scores higher than four are considered as positive comments, those with scores equal to 4 are considered as neutral comments and reviews with scores lower than four are considered as negative comments.

Table 2. Review.

review	score of Reviews	sentiment label
Both times we were there very late, after 11 PM. At that time, many diners (forget restaurants!) you get warmed-over food and lousy service. Not so here - food was uniformly very good and the service quite good. There weren't many people but it...More	4 out of 5	neutral
On Tuesdays Moore's has 2 for 1 burgers. The burgers are good and large. The service however was not so good. My wife ordered bacon with her burger but did not get it. When we mentioned this to the waitress so we wouldn't be charged...More	3.5 out of 5	negative
Shlomo is the nicest Man!!! He is the owner and he always personally greets me and knows my order, Chicken kebab with spicy rise no salad. I just tried the hummus this past	5 out of 5	positive

week, it is much better than the store bought stuff. The portions...More		
--	--	--

Then, we have deleted any information that included missing fields, special characters, or lowercase letters, and we concentrated on the following fields: comments, which are customer feedback. We then remove stop words, such as "a," "an," "and," "but," "the," "that," "of," and "from" from the review. By removing stop words widely presented in multiple classes, the classification processes are made more accessible. Then we apply a stemming algorithm to remove the words' prefixes and suffixes, leaving only their stem. Stemming aids in combining all of a word's derivatives that do not vary conceptually into a single idea. For instance, terms like "eating" and "eaten" are all regarded as "eat" in documents including such words [16]. Performing stemming prevents us from doing classification of unnecessary words since they represent the same meaning as their root words. This reduces dimensions and increases performance significantly.

Step 2- Model training, Classification & prediction and evaluation, vectorization.

In machine learning, the process of vectorization is an essential part of the feature extraction phase. The objective of converting textual information into numerical vectors is to extract discernible features from the written content, which can be utilized for training the model [17]. Most Deep Learning Architectures and Machine Learning algorithms are unable to interpret raw strings or plain text. To perform various tasks such as classification, regression, and grouping, numerical data is required as input. Moreover, extracting pertinent information from the extensive volume of textual data available is imperative to develop useful applications [18].

In the fields of natural language processing (NLP) and information retrieval, TF-IDF (Term Frequency-Inverse Document Frequency) is a popular method for text vectorization. Documents of text are converted into numerical vectors for use in subsequent machine-learning processes.

TF-IDF is a metric that quantifies the importance of keywords in identifying and classifying texts; alternatively, it offers the keywords themselves. Take the case of a blogger who has hundreds of writers contributing to his site and who has recently recruited an intern whose primary responsibility will be to create fresh blog entries every day. Interns are notorious for not correctly tagging blog articles, which leads to a lot of unorganized content. The TF-IDF method, which can detect tags mechanically for bloggers, works particularly well under such conditions. Without having to worry about tags, bloggers and interns will save a tonne of time [19].

3.1 Model Training, Classification and Prediction

After conducting a thorough examination of the data, we are now able to proceed with the training of our algorithm, which aims to classify sentiment into three separate groupings: positive, negative, and neutral. Then we split the 80% data into train and 20% data for test. To make predictions on the test data, we will also fit four distinct classifiers into the training dataset, which are GaussianNB, SVM, Logistic Regression, as well as Random Forest. Lastly, we will use a cross-validation technique to improve our evaluation result.

3.2 GaussianNB

The Naive Bayes method assumes that the predictors are chosen randomly and have no systematic effect on the target class. The Naive Bayes model often produces a decent outcome [20], despite its unrealistic assumption that all predictors are independent. The Gaussian Naive Bayes classifier is used when the predictor's values are continuous and are presumed to follow a Gaussian distribution.

3.3 SVM

The straightforward Support Vector Machine (SVM) will be used, since it is often used for classification and linear regression jobs. There are many different uses for Support Vector Machines (SVMs), including web page analysis, intrusion detection systems, face recognition, email categorization, gene classification, and handwriting recognition. In this inquiry, as well as many others, machine learning has made use of Support Vector Machines (SVMs). Linear and non-linear data can be effectively analyzed using regression and classification techniques. One of the reasons for utilizing Support Vector Machines (SVMs) is their capability to uncover complex relationships within data without requiring extensive manual adjustments. It's an excellent choice when working with smaller data sets comprising tens of thousands to hundreds of thousands of features. According to McGregor, algorithms that process intricate and minuscule data yield more accurate results than alternative algorithms [21].

3.4 Logistic Regression

Sperandei posits that logistic regression is employed to derive Odds ratios when multiple explanatory variables are present [22]. The process is similar to multiple linear regression, except that the answer variable is a binomial. The result is how each variable will affect the chance that a significant event will be seen. One notable advantage lies in potentially avoiding confounding effects by examining the interrelationships among all variables.

3.5 Random Forest

The machine learning technique known as Random Forest is very powerful. The foundation of this approach is the Ensemble Learning method, specifically bagging. The Random Forest algorithm possesses several advantages. The bagging technique forms the fundamental basis for the Random Forest algorithm, which utilizes ensemble learning. A large number of trees are generated from the subset of data, and their combined output is then analyzed. By employing this approach, the problem of overfitting in decision trees is mitigated, resulting in a reduction in variance and an improvement in accuracy [23]. Furthermore, this classifier can be utilized to address both classification and regression problems.

3.6 Voting Ensemble Classifier

In addition to utilizing a single classifier, we also use a Voting Classifier that combines three classifiers that shows the highest precisions (SVM, Logistic Regression, and Random Forest) by employing 'hard' voting, which determines the final prediction depending on which classifiers received the majority of the votes. In an ensemble classifier, using an even number for majority voting might result in a possible tie if an equal number of classifiers vote for distinct classes. We can guarantee a majority choice by selecting an odd number of classifiers in the

ensemble, even if two classifiers vote for one class and the remaining classifier(s) vote for a different class. In this manner, the ensemble can provide a conclusive forecast, and the majority opinion can be confidently regarded as the outcome. The Voting Classifier accepts a list of tuples, where each tuple comprises the classifier's name and its actual instance. Each classifier's predictions are combined, and the majority class is chosen as the final prediction. The cross-validation scores and performance data for the Voting Classifier, which reflects the three classifiers' combined majority vote. Comparing this ensemble technique to individual classifiers may help to increase performance and robustness overall.

3.7 Cross-validation

We evaluate the classification accuracy with cross-validation techniques as we try different classifiers. Evaluations and Outcomes are presented in the result section. Cross-validation is a helpful tool. We use cross-validation techniques to the precision of the classifiers that we used. The utilization of data has enabled us to enhance its efficacy and gain substantial insights into the performance of our algorithm. When developing a machine-learning model using designated data, it is common practice to partition the data into training and validation/test sets. The training set is used to develop the model, whereas the test set is used to evaluate the model using data that it has never seen before. The conventional approach involves a clear-cut division of 80%–20%, sometimes with alternative ratios such as 70%–30% or 90%–10%. Multiple splits are performed during the process of cross-validation. K from 3 to 10 may be used [24].

4 Results

With model performance metrics between 94.6% and 94.8% and nearly equivalent precision, recall, and F1-scores, Random Forest, Logistic Regression, SVM, and the Majority Voting Ensemble Classifier perform similarly (see Table 2). They are effective at this categorization job. In comparison to the other classifiers, Gaussian Naive Bayes has a significantly lower mean accuracy (87.7%). While the performance of this model is satisfactory, it does not achieve the same level of accuracy and F1-score as the other three models. The classifiers performed well, with mean cross-validation accuracies ranging from 0.877 to 0.948. On the test set, every classifier had excellent accuracy, precision, recall, and F1 scores, ranging from 0.905 to 0.97. The Majority Voting Ensemble classifier increased accuracy to 0.97, matching the top-performing individual classifiers. Additionally, the ensemble method demonstrated consistent F1-scores, recall, and accuracy.

Overall, based on the following table 3, the majority voting ensemble classifier, SVM, Random Forest, Logistic Regression, and Logistic Regression all showed great consistency and performance, with mean cross-validation accuracies ranging from 0.946 to 0.9481. The majority voting ensemble approach improved accuracy on the test set without degrading precision, recall, or F1-score, which makes it effective for sentiment analysis of restaurant reviews. These results suggest that the accuracy and robustness of sentiment analysis models for restaurant reviews may be improved by combining a variety of classifiers using ensemble approaches. We just tested the first 1000 sample datasets. Prior to making a definitive decision, it is recommended to conduct additional experiments and perform further tests using a larger and more comprehensive dataset to evaluate the stability and generalizability of the model's performance.

When deciding, it's crucial to take into account additional aspects, including model complexity, interpretability, and processing resources.

Table 3. Result of classification and prediction.

Classifier	Precision on Test Set	Recall on Test Set	F1-score on Test Set	Cross-validation scores	Mean accuracy with Cross-validation	Accuracy on Test Set
GaussianNB	0.938877	0.905	0.921627	[0.87 0.88 0.875 0.835 0.925]	0.877	0.905
SVM	0.9409	0.97	0.955228	[0.945 0.945 0.945 0.95 0.955]	0.9480000 000000001	0.97
LogisticRegression	0.9409	0.97	0.955228	[0.945 0.945 0.945 0.945 0.95]	0.946	0.97
RandomForest	0.9409	0.97	0.955228	[0.945 0.945 0.945 0.95 0.955]	0.9480000 000000001	0.97
Voting Ensemble Classifier	0.9409	0.97	0.955228	[0.945 0.945 0.945 0.95 0.955]	0.9480000 000000001	0.97

5 Conclusion

In conclusion, restaurant firms in the food sector may better understand consumer preferences and attitudes by using the informative and helpful approach of sentiment analysis of customer evaluations. Sentiment analysis applies machine learning and natural language processing approaches to identify positive, negative, or neutral comments.

To increase client happiness and the entire eating experience, the managers and owners of restaurants may use this study to pinpoint their businesses' major strengths and shortcomings. Positive feelings help a business capitalize on its assets and improve its image by offering insightful feedback on qualities that patrons value, such as exquisite cuisine, excellent service, and a welcoming atmosphere.

On the other hand, negative feelings highlighted problem areas, such as delayed service, subpar food, or sanitary problems. Restaurants may avoid bad reviews harming their reputation and patron loyalty by promptly resolving these concerns. Additionally, sentiment analysis may assist

businesses in monitoring alterations in client preferences over time as well as patterns in those preferences. In order to adapt services, menus, and marketing tactics to changing consumer demands, this information may be employed. It's critical to understand that sentiment analysis has its limits. It may be challenging to grasp context, use sarcasm, and communicate across cultural boundaries.

While Random Forest and Voting Ensemble Classifiers top the present evaluation, the best classifier ultimately depends on a number of other elements, including model interpretability, computing power, and the quantity of the training dataset. In addition, to confirm the classifiers' capacity to generalize, the performance should be further tested using an independent test dataset. Furthermore, to collect more complex sentiment data and boost model performance, one might experiment with feature engineering methods, sentiment lexicons, or embeddings. Moreover, to gain even greater accuracy and interpretability, it may be investigated to assemble numerous classifiers or to use sophisticated deep learning models like LSTM or Transformers.

References

- [1] Virmani, C., Pillai, A., & Juneja, D. (2017). Extracting information from social network using nlp. *International Journal of Computational Intelligence Research*, 13(4), 621-630.
- [2] Tetsuya Nasukawa and Jeonghee Yi. (2003). Sentiment analysis: capturing favorability using natural language processing. In *Proceedings of the 2nd international conference on Knowledge capture (K-CAP '03)*. Association for Computing Machinery, New York, NY, USA, 70–77. <https://doi.org/10.1145/945645.945658>
- [3] Houlihan, P., & Creamer, G. G. (2021). Leveraging social media to predict continuation and reversal in asset prices. *Computational Economics*, 57(2), 433-453.
- [4] V. Sehgal and C. Song, "SOPS: Stock Prediction Using Web Sentiment," *Seventh IEEE International Conference on Data Mining Workshops (ICDMW 2007)*, Omaha, NE, USA, 2007, pp. 21-26, doi: 10.1109/ICDMW.2007.100.
- [5] Obiedat, R., Qaddoura, R., Ala'M, A. Z., Al-Qaisi, L., Harfoushi, O., Alrefai, M. A., & Faris, H. (2022). Sentiment analysis of customers' reviews using a hybrid evolutionary svm-based approach in an imbalanced data distribution. *IEEE Access*, 10, 22260-22273.
- [6] Laksono, Rachmawan Adi, et al. (2019) "Sentiment analysis of restaurant customer reviews on tripadvisor using naïve bayes." *2019 12th international conference on information & communication technology and system (ICTS)*. IEEE.
- [7] Spoorthi, C., Kumar, P. R., & Adarsh, M. J. (2018). Sentiment analysis of customer feedback on restaurants. *Int. J. Eng. Res. Technol*, 6, 1-4.
- [8] Kulka, G. (2022). The Role of Sentiment Analysis in Investment Decision-Making. JK Investment Group. <https://www.jkinvestmentgroup.com/blog/the-role-of-sentiment-analysis-in-investment-decision-making-1>
- [9] Hasselgren B, Chrysoulas C, Pitropakis N, Buchanan WJ. Using Social Media & Sentiment Analysis to Make Investment Decisions. *Future Internet*. 2023; 15(1):5. <https://doi.org/10.3390/fi15010005>
- [10] Collomb, A., Costea, C., Joyeux, D., Hasan, O., & Brunie, L. (2014). A study and comparison of sentiment analysis methods for reputation evaluation. *Rapport de recherche RR-LIRIS-2014-002*.
- [11] Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., & McClosky, D. (2014). *The Stanford CoreNLP Natural Language Processing Toolkit* (pp. 55–60). Association for Computational Linguistics.
- [12] Al-Otaibi, S., Alnassar, A., Alshahrani, A., Al-Mubarak, A., Albugami, S., Almutiri, N., & Albugami, A. (2018). Customer Satisfaction Measurement using Sentiment Analysis. *International*

- Journal of Advanced Computer Science and Applications, 9(2).
<https://doi.org/10.14569/ijacsa.2018.090216>
- [13] Chowdhary, K.R. (2020). Natural Language Processing. In: Fundamentals of Artificial Intelligence. Springer, New Delhi. https://doi.org/10.1007/978-81-322-3972-7_19
- [14] Bhattacharyya, S. (2023). The Ethical Considerations of Natural Language Processing (NLP) | Analytics Steps. www.analyticssteps.com. <https://analyticssteps.com/blogs/ethical-considerations-natural-language-processing-nlp>
- [15] Sunday, O. (2022). Ethical Dilemma in Artificial Intelligence: A Focus on Natural Language Processing bias. Information Matters.<https://informationmatters.org/2022/09/ethical-dilemma-in-artificial-intelligence-a-focus-on-natural-language-processing-bias/>
- [16] Asani, E., Vahdat-Nejad, H., & Sadri, J. (2021). Restaurant recommender system based on sentiment analysis. Machine Learning with Applications, 6, 100114.
- [17] Jha, A. (2021). Vectorization Techniques in NLP [Guide]. Neptune.ai. <https://neptune.ai/blog/vectorization-techniques-in-nlp-guide>
- [18] GOYAL, C. (2021). Text Vectorization and Word Embedding | Guide to Master NLP (Part 5). Analytics Vidhya. <https://www.analyticsvidhya.com/blog/2021/06/part-5-step-by-step-guide-to-master-nlp-text-vectorization-approaches/>
- [19] Qaiser, S., & Ali, R. (2018). Text Mining: Use of TF-IDF to Examine the Relevance of Words to Documents. International Journal of Computer Applications, 181(1), 25–29. <https://doi.org/10.5120/ijca2018917395>
- [20] Sharma, P. (2021). Implementation of Gaussian Naive Bayes in Python Sklearn. Analytics Vidhya. <https://www.analyticsvidhya.com/blog/2021/11/implementation-of-gaussian-naive-bayes-in-python-sklearn/#h-types-of-naive-bayes-classifiers>
- [21] McGregor, M. (2020). SVM Machine Learning Tutorial – What is the Support Vector Machine Algorithm, Explained with Code Examples. FreeCodeCamp.org. <https://www.freecodecamp.org/news/svm-machine-learning-tutorial-what-is-the-support-vector-machine-algorithm-explained-with-code-examples/>
- [22] Sperandei, S. (2014). Understanding logistic regression analysis. Biochemia medica, 24(1), 12-18.
- [23] Naresh Kumar. (2019). Advantages and Disadvantages of Random Forest Algorithm in Machine Learning. Blogspot.com. <http://theprofessionalspoint.blogspot.com/2019/02/advantages-and-disadvantages-of-random.html>
- [24] Shulga, D. (2018). 5 Reasons why you should use Cross-Validation in your Data Science Projects. Medium. <https://towardsdatascience.com/5-reasons-why-you-should-use-cross-validation-in-your-data-science-project-8163311a1e79>