# Sentiment Analysis Using Machine Learning Approach

Andreea-Maria Copaceanu
*The Bucharest University of Economic Studies, Romania*
*andreea.copaceanu@gmail.com*

## Abstract

*Customers feedback is a valuable asset for businesses, that can be used in order to improve their performance. One of the fastest spreading areas today in computer science - Sentiment Analysis, helps to extract precious information from textual data, in order to identify the feeling of a statement. This research aims to build a classifier to predict customers' satisfaction, based on Amazon reviews dataset, for different brands of mobile phones. The paper proposes a comparison between four text classification algorithms - Naïve Bayes, Support Vector Machine, Decision Tree and Random Forest, using different feature extraction techniques, such as Bag of words and TF-IDF. In addition, the models are evaluated using accuracy, precision, recall and F-score metrics. Our experiments revealed that Support Vector Machine achieves the best results and is very suitable for classification of the sentiment on product reviews.*

## 1. Introduction

Online reviews and recommendations have a big impact on customers purchasing decisions, especially now, when people tend to express their opinions and feelings more than ever, on virtual communities and social networks.

According to a survey report, 93% of consumers refer to online reviews before taking their purchase decisions (Kaemingk, 2021). This can be explained by the fact that in general, our decisions are influenced by other opinions, when dealing with something new (Alharbi, 2021).

Purchase is always an interaction between two entities, customers and business owners. Customers can use reviews to make better decisions about what products to buy, while businesses, on the other hand, benefit from reviews in terms of gaining useful information about customers satisfaction on their products. This information can be then used to evaluate their marketing strategies, to improve their products and to enhance their performance (Al-Sheikh, 2018).

Classifying large amounts of unstructured data from Internet is a challenging task. Hence, the sentiment analysis, along with Natural Language Processing techniques, flourished in recent years, to provide a framework for analysis of textual data obtained from reviews. These techniques predict the polarity of the opinions (positive, negative, or neutral), assisting customers to have a conclusion about a product. On the other hand, companies can understand in this way the level of satisfaction of their customers (Alharbi, 2021).

Sentiment analysis is a natural language processing problem, which implies the detection and retrieval of knowledge from textual data. The sentiment analysis follows a sequence of steps such as the reviews collection, the lowercase conversion, punctuation and additional spaces removal, stop words removal, tokenization, lemmatization, feature extraction and finally classification (Dadhich, 2021).

Amazon is one among the most important e-commerce retailers, used every day for online shopping. The Amazon ranking system ranges from 1 to 5, where "1" is the worst rating and "2" is the highest rating (Roshan, 2020).

In order to assess the overall semantics of consumer feedback, this paper explores the sentiment classification into positive or negative feelings, for online reviews, using specific methods applied in this domain.

This research aims to build a classifier to predict consumers satisfaction, whose performance will be evaluated, based on the dataset of the mobile phone reviews. This has the potential to help companies to improve their products and on the other hand, to help potential customers to make better purchasing decisions.

The paper is structured as follows. This paper begins with the introduction. Section 2 discusses the related work in the previous literature. Section 3 explains both research methodology and implementation respectively. Section 4 reports the experimental results in terms of performance metrics for various classifiers. Lastly, section 5 concludes the findings of the paper and exposes the future scope.

## 2. Literature review

Various studies focused on the problem of identifying customers opinions on different products using Amazon reviews.

In the following, these papers are reviewed in terms of pre-processing techniques, feature extraction methods, proposed methodologies, and evaluation metrics.

(Guia, 2019) applied supervised machine learning algorithms such as Naïve Bayes, Support Vector Machine, Decision Trees and Random Forest, to predict the reviews sentiment, based on Amazon Reviews: Unlocked Mobile Phones dataset. In the classification process, the authors used only Rating and Review attributes and removed the instances with neutral reviews. After applying preprocessing steps such as converting the dataset into lowercase, tokenization, removal of punctuation and stop words, the data was splitted into 80% for training and 20% for test. The results for the application of classifiers, show that the Support Vector Machine classifier is the most accurate, with the highest values for all metrics, followed by Random Forest classifier. The authors presented also a statistical study in terms of the impact of brand and price in the reviews polarity. They concluded that ZTE has the most positive reviews rate with 82,9 % positive reviews and in terms of price, more positive reviews were obtained for higher prices, which can be explained by the quality of the phones.

The research conducted by (Aljuhani, 2019), studied the performance of different machine learning algorithms, such as Logistic regression, Naïve Bayes, Stochastic gradient descent and convolutional neural network respectively, using different features extraction techniques such as BOW and TF-IDF, each of them with three variations depending on the number of grams used. The authors divided the data into 70% for training, 15% for testing and 15% for development. The results revealed that convolutional neural network provided the best results.

In their papers, both (Aljuhani, 2019) and (Bansal, 2018), used unbalanced and balanced datasets. While (Aljuhani, 2019) categorized both balanced and unbalanced data into, five and four stars as positive rating, one and two starts as negative rating, and three stars as neutral rating, (Bansal, 2018) categorized balanced and unbalanced data separately. (Bansal, 2018) used balanced data, meaning that the number of negative reviews (1 and 2 stars) is equal to the number of positive reviews (4 and 5 stars) and removed neutral reviews. For unbalanced data, they categorized (1 and 2 stars) as negative reviews and (3, 4 and 5 stars) as positive reviews. (Bansal, 2018) applied deep learning methods such as, CBOW and skip-gram, with different machine learning algorithms: SVM, Naïve Bayes, Logistic Regression and Random Forest. The experimental results showed that Random Forest using CBOW achieved the best accuracy.

(Shaheen, 2019) performed a sentiment classification on mobile phone reviews dataset, using seven different classifiers and based on their results, the Random Forest classifier outperformed all other classifiers, with an accuracy of 85% for the given dataset, followed by LSTM and CNN. The authors also exposed the distribution of reviews with respect to their ratings, showing that most reviewers have rated 4 stars and 3 stars. Also, the study concluded that there is a direct correlation between rating and price.

(Ravi, 2019) implemented four algorithms, Naïve Bayes, Support Vector Machine, Random Forest and K-Nearest Neighbor, using different sizes of training and test data. They concluded that Random Forest classifier produced the best accuracy metrics.

(Qaiser, 2021) focused on the comparison of machine learning methods applied in Sentiment Analysis such as Naïve Bayes, Decision Tree, Support Vector Machine and the modern method, Deep Learning. The ML techniques were applied to a dataset of 4289 rows, about technological impact on employment, remaining with 1047 rows, after completing the preprocessing steps (Qaiser, 2021). They concluded that the deep learning method performed the best, with an accuracy of 96,41%, followed by Naïve Bayes and Support Vector Machine with 87,18% and 82,05% respectively.

(Tan, 2018) conducted a study on a dataset of 34660 instances, from customer reviews of Amazon products. The dataset was divided into a training set of 60%, a validation set of 20% and a test set of 20%. They implemented machine learning algorithms such as Naïve Bayes, Support Vector Machine with Linear Kernel, Support Vector Machine with RBF Kernel, KNN-4, 5, 6 and deep neural networks, such as Recurrent Neural Network. They concluded that the Long Short-Term Memory generates the most accurate predictions.

## 3. Research methodology

This section presents the methodology and techniques used for the classification of mobile phone reviews. Figure no. 1 illustrates the phases of this research, starting with the dataset of Amazon reviews, until each review is classified into positive or negative.

*Figure no. 1. Phases of research*
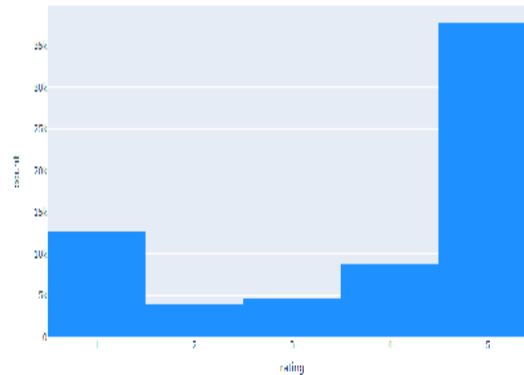


*Source:* Author prelucration

## A. Dataset

Our dataset (Amazon Cell Phones Reviews | Kaggle, 2019), consists in 67986 instances, fetched between 24th November 2003 to 25th December 2019. The data was retrieved from Amazon.com and focuses on reviews for both unlocked and locked carriers, related to ten brands: Apple, ASUS, Google, HUAWEI, Motorola, Nokia, OnePlus, Samsung, Sony and Xiaomi.

The dataset contains the following attributes:
1. "asin": ID of the product
2. "name": name of the reviewer
3. "rating": rating of the product
4. "date": date of the review
5. "title": title of the review
6. "review": text of the review
7. "helpfulVotes": rating of the review's helpfulness

In our analysis, we will focus only on the rating and review features, as these are the most useful and relevant for model building. In order to have an overview of the reviews dataset, the ratings distribution is shown in Figure no. 2. The classes are imbalanced, as classes 2, 3 and 4 have very small amount of reviews, compared to class 5.
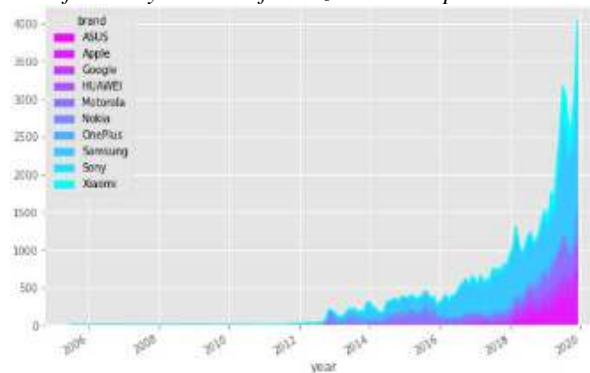
*Figure no. 2. Rating distribution for Amazon reviews dataset*



*Source:* Author computation

In terms of popularity, Xiaomi and Samsung are the most rated brands, over the years. (Fig. 5)

*Figure no. 3. Distribution of monthly number of Amazon reviews per brand*



*Source:* Author computation

Before applying the preprocessing tasks, we are going to label the dataset, as follows:
1. Rating with a value less than 3, is labeled as "Negative".
2. Rating with a value greater than 3, is labeled as "Positive".

In a 5-star rating scale, 3-star ratings are considered as neutral reviews, which means that the reviews are neither positive nor negative. So, we remove the 3-star rating reviews from our dataset.

**B. Data Preprocessing**

The performance of a classifier can be highly increased by preprocessing the data. Considering this, the preprocessing phase, applied to our dataset, included the following steps:
1. Convert the uppercase letters into lowercase.
2. Remove all the URLs starting with HTTP.
3. Remove all the special characters.
4. Remove all single characters.
5. Remove single characters from the start.
6. Substitute multiple spaces with single space.
7. Remove all the punctuation.
8. Remove the stop words, such as "the", "a", "in", using Stopwords Corpus for English words, from NLTK library.
9. Tokenization, the process of splitting the original text in the form of sentences into words.

10. Lemmatization, the process of transforming the word into its significant base structure, using Wordnet Lemmatizer from NLTK library.

After cleaning the text data, the dataset was splitted into 80% for training set, used to learn the models and 20% for testing set, used to calculate the model's performance.

## C. Word clouds of reviews for Mobile brands

Word cloud is a widely data visualization technique used for representing text data, in which the size of words indicates their frequency or importance. Both types of reviews contain some common words like "buy", "battery" or "one". Figure no. 4 shows that the most frequent words encountered in positive reviews are: "great", "good", "love", "use", "life", etc. On the other hand, the most frequent negative reviews words are "return", "screen", "charge", "back", as seen in Figure no. 5.

*Figure no. 4. Word Cloud vizualization for Amazon positive reviews*



*Source:* Author computation

*Figure no. 5. Word Cloud vizualization for Amazon negative reviews*



*Source:* Author computation

## D. Feature extraction

When dealing with text features, the original text needs to be converted into a document-term, since the machine learning algorithms do not support text features. Thus, after the preprocessing stage, data will be vectorized, using the following methods: Bag of Words (BoW) and TD-IDF (Term Frequency-Inverse Document Frequency). The result from each method will be a matrix, that represents the text as vectors, which can be fed to the machine learning algorithms to build classification models.

## E. Topic modelling

Latent Dirichlet Allocation (LDA) is an example of a model which is used to classify text in a document, to a topic. It builds a topic per document and shows the most relevant words per each topic. Our LDA model was created using the Gensim library. For visualize topics along with the most relevant words, pyLDAvis library was used. In Figure no. 6, there are shown top 30 most relevant words for our topics.

*Figure no. 6. Topic modelling visualization using LDA*



*Source:* Author computation

## F. Classification Models

In language processing, most part of classifications are performed using supervised machine learning, and this will be the subject of this paper. Text classification can be done using different algorithms. In this context, the algorithms implemented for classification are named classifiers. This section describes four of the most used supervised classifiers in text classification.

- **Naïve Bayes** technique will select the best class for a document, based on the probability that the terms in the document belong to that class. From a mathematical point of view, the probability of classifying the document in a class c is (Martin, 2017):

$$C_{NB} = \text{argmax } P(c) \prod_{1 \le k \le n_d} P(t_k|c) \tag{1}$$

Where,
- $P(c)$ is the probability that a document belongs to class c (based on training data), also called previous class probability.
- $P(t_k | c)$ is the probability that a term t from position k in document d, will be found in documents of class c.
- $n_d$ is the number of terms in document d.

- **Support Vector Machine** tries to find the optimal hyperplane which could separate the data into two classes in the case of binary classification. From a geometric point of view, given two types of points in a space, it tries to minimize the distance from one of the points to the other. This minimization problem is equivalent to the following problem (Fan, 2018):

$$\min \frac{1}{2}\|w\|^2 \tag{2}$$

Where w is the direction of a vector x.

- **Decision Tree** is a hierarchical model of supervised learning, in which local regions are represented as a series of recursive separations by decision nodes. This classifier has a tree-like structure, in which each internal node tests an attribute, each branch represents the test result, and each terminal node indicates the class label. The data equation is as follows (Ronaghan, 2018):

$$\sum_{i=1}^{c} -f_i \log(f_i) \tag{3}$$

Where $f_i$ is the frequency of label i at a node and C is the number of unique labels.

- **Random Forest** is an integration of several decision trees. Random Forest builds several decision trees in the training phase, the final prediction being a prediction based on the result of the predictions of all decision trees. Each time a division into a new decision tree is considered, a random selection of m predictors is made, as potential candidates, from the total number of trees.

The equation of the Random Forest classifier is as follows (Guillot, 2017):

$$f^B(x) = \frac{1}{B} \sum_{b=1}^{B} T_b(x) \tag{4}$$

## G. Evaluation Metrics

Evaluation metrics play an important role to measure the classification performance. In order to evaluate the results of the four algorithms four of the most popular measures are used: Accuracy, Precision, Recall, and F1 score. These four metrics are explained in the following.

- **Accuracy** predicts how often the classifier makes the correct prediction. Accuracy is the ratio between the number of correct predictions and the total number of predictions. (Martin, 2017):
- **Precision** measures the exactness of a classifier; how many of the return documents are correct. A higher precision means less false positives, while a lower precision means more false positives. Precision is the ratio of numbers of instance correctly classified from total. (Martin, 2017)
- **Recall** calculates the sensitivity of a classifier; how many positive data it returns. Higher recall means less false negatives. Recall is the ratio of number of instances accurately classified to the total number of predicted instances (Martin, 2017).
- **F-score** is the weighted harmonic mean of precision and recall (Martin, 2017).

## 4. Findings

In our research, the reviews have been classified as positive and negative, based on the star rating.

There were several machine learning algorithms employed in this paper such as Multinomial Naïve Bayesian, Support Vector Machine, Random Forest and Decision Tree. Different feature selection techniques were applied on classifiers, such as TF-IDF and Bag of Words. For both BOW and TF-IDF we used three variations of grams, unigrams, bigrams and trigrams. Firstly, BOW was applied for each machine learning algorithm. Then, TF-IDF, was applied, with parameters min_df=5 and max_df=0.8, which means to ignore terms that appear in less than 5 documents, respectively to ignore terms that appear in more than 80% of documents. Lastly, we used the Scikit-Learn Pipeline method, which chains together TfidfTransformer and the CountVectorizer (Bengfort, 2018).

Tables no. 1, 2, 3, 4 and 5 show the results for all mentioned classifiers, using different features extraction techniques.

*Table no. 1 Results of Naïve Bayes for Amazon reviews dataset (train split)*

|                        | Accuracy | Recall | Precision | F-score |
|------------------------|----------|--------|-----------|---------|
| BOW + U                | 0.903    | 0.903  | 0.902     | 0.902   |
| BOW + B                | 0.896    | 0.896  | 0.897     | 0.891   |
| BOW + T                | 0.834    | 0.834  | 0.844     | 0.812   |
| TF-IDF + U             | 0.891    | 0.891  | 0.892     | 0.886   |
| TF-IDF + B             | 0.881    | 0.881  | 0.885     | 0.872   |
| TF-IDF + T             | 0.784    | 0.784  | 0.794     | 0.737   |
| Pipeline (BOW + TF-IDF)| 0.866    | 0.866  | 0.875     | 0.853   |

*Source:* Author's
computation

*Table no. 2 Results of Support Vector Machine with RBF kernel for Amazon reviews dataset (train split)*

|  | Accuracy | Recall | Precision | F-score |
|---|---|---|---|---|
| BOW + U | 0.857 | 0.857 | 0.850 | 0.829 |
| BOW + B | 0.842 | 0.842 | 0.854 | 0.822 |
| BOW + T | 0.747 | 0.747 | 0.799 | 0.656 |
| TF-IDF + U | 0.926 | 0.926 | 0.925 | 0.926 |
| TF-IDF + B | 0.887 | 0.887 | 0.888 | 0.881 |
| TF-IDF + T | 0.777 | 0.777 | 0.783 | 0.728 |
| Pipeline (BOW + TF-IDF) | 0.927 | 0.927 | 0.926 | 0.926 |

*Source:* Author computation

*Table no. 3 Results of Support Vector Machine with linear kernel for Amazon reviews dataset (train split)*

|  | Accuracy | Recall | Precision | F-score |
|---|---|---|---|---|
| BOW + U | 0.849 | 0.849 | 0.820 | 0.821 |
| BOW + B | 0.862 | 0.862 | 0.858 | 0.857 |
| BOW + T | 0.775 | 0.775 | 0.782 | 0.723 |
| TF-IDF + U | 0.917 | 0.917 | 0.916 | 0.916 |
| TF-IDF + B | 0.877 | 0.877 | 0.876 | 0.871 |
| TF-IDF + T | 0.780 | 0.780 | 0.783 | 0.734 |
| Pipeline (BOW + TF-IDF) | 0.927 | 0.927 | 0.926 | 0.926 |

*Source:* Author computation

*Table no. 4 Results of Decision Tree for Amazon reviews dataset (train split)*

|  | Accuracy | Recall | Precision | F-score |
|---|---|---|---|---|
| BOW + U | 0.853 | 0.853 | 0.853 | 0.853 |
| BOW + B | 0.826 | 0.826 | 0.823 | 0.824 |
| BOW + T | 0.792 | 0.792 | 0.786 | 0.763 |
| TF-IDF + U | 0.859 | 0.859 | 0.859 | 0.859 |
| TF-IDF + B | 0.832 | 0.832 | 0.828 | 0.829 |
| TF-IDF + T | 0.769 | 0.769 | 0.750 | 0.736 |
| Pipeline (BOW+TF-IDF) | 0.861 | 0.861 | 0.861 | 0.861 |

*Source:* Author computation

*Table no. 5 Results of Random Forest for Amazon reviews dataset (train split)*

|  | Accuracy | Recall | Precision | F-score |
|---|---|---|---|---|
| BOW + U | 0.909 | 0.909 | 0.908 | 0.906 |
| BOW + B | 0.866 | 0.866 | 0.863 | 0.860 |
| BOW + T | 0.788 | 0.788 | 0.802 | 0.743 |
| TF-IDF + U | 0.914 | 0.914 | 0.912 | 0.912 |
| TF-IDF + B | 0.865 | 0.865 | 0.862 | 0.859 |
| TF-IDF + T | 0.778 | 0.778 | 0.769 | 0.740 |
| Pipeline (BOW +TF-IDF) | 0.913 | 0.913 | 0.912 | 0.910 |

*Source:* Author computation

For Naïve Bayes, we can see in Table no. 1 that Bag of Words with unigrams achieved the highest accuracy with a value of 90,3%. As seen in Table no. 2 and Table no. 3, both types of kernel, RBF, respectively Linear kernel, were used to evaluate the Support Vector Machine. Both experimental evaluations demonstrate that the classifier obtains the best performance when using the pipeline approach, achieving 92,7% for accuracy and recall and 92,6% for precision and F-score. TF-IDF with unigrams achieved an accuracy of 92,6% for RBF kernel and 91,7% for linear kernel, which can be considered very good results.

Decision Tree achieved its highest accuracy of 86,1%, with the pipeline technique, followed by TF-IDF with unigrams, where an accuracy of 85,9% was obtained.

As shown in Table no. 5, the best performance for Random Forest was 91,4% for accuracy and recall and 91,2% for precision and F-score. So, the performance of the classifier is very high.

TF-IDF achieved the highest performance of 92,6% accuracy and recall, for unigrams. On the other hand, Bag-of-words obtained its highest accuracy of 90,9%, with unigrams. In contrast, overall lower results were obtained for bigrams or trigrams, compared to unigrams, when applied to both BOW and TF-IDF.

From all the experiments the best results of all four algorithms are obtained by Support Vector Machine and Random Forest.

## 5. Conclusions

Reviews are essential for both customers and companies. From consumers point of view, reviews help them to make better decisions when buying products. On the other hand, companies benefit from reviews, by knowing the level of consumers satisfaction about their products and acting accordingly. In this paper, we proposed a machine learning approach for text sentiment analysis.

We performed the sentiment analysis on mobile phone reviews dataset, using different types of machine learning algorithms, such as Naïve Bayes, Support Vector Machine, Decision Tree and Random Forest. We used different feature extraction approaches such as Bag of words and TF-IDF with unigrams, bigrams and trigrams and analyzed the classifiers results, based on four performance metrics: Accuracy, Precision, Recall and F1 score. We also proposed the Latent Dirichlet Allocation (LDA) model for topic extraction, which shows document topics along with the most relevant words for each topic. We described the basic theory behind the models, approaches used in our research and the performance metrics for the conducted experiments. We went through different research papers on sentiment analysis over text-based datasets.

Overall, we were able to achieve promising results for classifiers, based on the performance metrics obtained. We found that the pipeline approach, which combines TfidfTransformer and CountVectorizer, achieves the highest metrics results, for almost all the classifiers. We can also observe that, the unigrams applied to Bag of Words and TF-IDF, gives better results compared to bigrams or trigrams, for all classifiers. The highest accuracy is 92,7 %, obtained by Support Vector Machine classifier, for both types of kernel: inear and RBF. Even if Random Forest has its highest accuracy of 91,4%, it can be considered the most complete classifier, with high values for all the metrics. Our results show that Naïve Bayes is also a classifier to consider, being just slightly lower than the Random Forest classifier, with its highest accuracy of 90,3%.

As future work, we plan to continue to study other algorithms applied in Sentiment Analysis field and to evaluate them. Also, another future direction would be to collect more data, in order to test the performance of classifiers on a massive dataset and see if there are improvements in results. Another point to consider in the future, would be to adjust and explore more parameters for the classifiers, which could contribute to even better results. And not lastly, in the future we intend to explore more methods of linguistic analysis, such as semantic analysis.

## 6. References

- Alharbi, N. M., Alghamdi, N. S., Alkhammash, E. H. and Al Amri, J. F., 2021. Evaluation of Sentiment Analysis via Word Embedding and RNN Variants for Amazon Online Reviews. *Mathematical Problems in Engineering*, 2021, pp.1-10.
- Aljuhani, S. A. and Alghamdi, N. S., 2019. A comparison of sentiment analysis methods on Amazon reviews of Mobile Phones. *International Journal of Advanced Computer Science and Applications*, 10(6), pp. 608-617.
- Al-Sheikh, E. S. and Hasanat, M. H. A., 2018. Social media mining for assessing brand popularity. *International Journal of Data Warehousing and Mining*, 14(1), pp. 40-59.
- Bansal, B. and Srivastava, S., 2018. Sentiment classification of online consumer reviews using word vector representations. *Procedia computer science*, 132, 1147-1153.
- Bengfort, B., Bilbro, R. and Ojeda, T., 2018. *Applied text analysis with python: Enabling language-aware data products with machine learning*. Sebastopol: O'Reilly Media, Inc.
- Dadhich, A. and Thankachan, B., 2021. Sentiment Analysis of Amazon Product Reviews Using Hybrid Rule-based Approach. *International Journal Engineering and Manufacturing,* 2021, pp. 40-52.

- Fan, S., 2018. Understanding the mathematics behind Support Vector Machines, *Shuzhan Fan*, [online]. Available at: <https://shuzhanfan.github.io/2018/05/understanding-mathematics-behind-support-vector-machines> [Accessed 20 May 2021].
- Guia, M. S., Silva, R. and Bernardino, J., 2019. Comparison of Naïve Bayes, Support Vector Machine, Decision Trees and Random Forest on Sentiment Analysis. *Proceedings of the 11th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management (IC3K 2019)*, pp. 525-531.
- Guillot, D., 2017. Mathematical Techniques in Data Science Random forest. *University of Delaware* [online]. Available at: <http://www.math.udel.edu/~dguillot/teaching/MATH567-S2017/lectures-handout/MATH567-Week12-handout-Random_forest.pdf> [Accessed 15 May 2021].
- Kaemingk, D., 2021. Online reviews statistics to know in 2021, *Qualtrics XM // The Leading Experience Management Software,* [online]. Available at: <https://www.qualtrics.com/blog/online-review-stats/> [Accessed 20 May 2021].
- Martin, M., 2017. Predicting ratings of amazon reviews-techniques for imbalanced datasets. Dissertation (Master). HEC University of Liège.
- Qaiser, S., Yusoff, N., Ali, R., Remli, M. A. and Adli, H. K., 2021. A Comparison of Machine Learning Techniques for Sentiment Analysis. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, 12(3), pp. 1738-1744.
- Ravi, A. K., Khettry, A. R. and Sethumadhavachar, S. Y., 2019. Amazon Reviews as Corpus for Sentiment Analysis Using Machine Learning. *International Conference on Advances in Computing and Data Sciences*, 1045, pp. 403-411.
- Ronaghan, S., 2018. The Mathematics of Decision Trees, Random Forest and Feature Importance in Scikit-learn and Spark, *Towards Data Science,* [online]. Available at <https://towardsdatascience.com/the-mathematics-of-decision-trees-random-forest-and-feature-importance-in-scikit-learn-and-spark-f2861df67e3> [Accessed 12 May 2021].
- Roshan, R. P. J., 2020. Amazon Reviews Sentiment Analysis: A Reinforcement Learning Approach. Dissertation (Master). Griffith College Dublin.
- Shaheen, M. A., Awan, S. M., Hussain, N. and Gondal, Z. A., 2019. Sentiment Analysis on Mobile Phone Reviews Using Supervised Learning Techniques. *International Journal of Modern Education and Computer Science,* 2019, 7, pp. 32-43.
- Tan, W., Wang, X. and Xu, X., 2018. Sentiment analysis for Amazon reviews. *International Conference*, 2018, pp. 1-5.