

Gradient Boosting Approach for Sentiment Analysis for Job Recommendation and Candidate Profiling

Swapnil Singh

*Computer Engineering Department
Mukesh Patel School of Technology
Management and Engineering, NMIMS
University
Mumbai, India
swapnilsingh@outlook.in*

Deepa Krishnan

*Computer Engineering Department
Mukesh Patel School of Technology
Management and Engineering, NMIMS
University
Mumbai, India
deepa.krishnan@nmims.edu*

Pranit Sehgal

*Computer Engineering Department
Mukesh Patel School of Technology
Management and Engineering, NMIMS
University
Mumbai, India
pranit.sehgal88@nmims.edu.in*

Harshit Sharma

*Computer Engineering Department
Mukesh Patel School of Technology
Management and Engineering, NMIMS
University
Mumbai, India
harshit.sharma94@nmims.edu.in*

Tarun Surani

*Computer Engineering Department
Mukesh Patel School of Technology
Management and Engineering, NMIMS
University
Mumbai, India
tarun.surani55@nmims.edu.in*

Jayant Singh

*Computer Engineering Department
Mukesh Patel School of Technology
Management and Engineering, NMIMS
University
Mumbai, India
jayant.singh95@nmims.edu.in*

Abstract— Sentiment Analysis has increasingly been used nowadays in many applications to evaluate opinion of public about products, policies, movies, politics. It is also used by government and law enforcement to understand behavior of people. One of the potential applications of sentiment analysis is candidate profiling and job recommendation. In the proposed research work, we evaluated the performance of supervised machine learning algorithms on dataset generated by us from twitter and indeed. We illustrated the steps involved in preprocessing the dataset generated through web scraping and making it ready for feeding into supervised algorithms. From our experimental study it is observed that Gradient Boosting Classifier gave the highest classification accuracy of 78.08 percent and AUC score of 0.819 on the test dataset.

Keywords—Sentiment Analysis, Job Reviews, Gradient Boost Classifier, Machine Learning

I. INTRODUCTION

The first thing you do before going to a café, restaurant, hotel, resort, company, or even before buying a product, is to see its reviews. There is an overwhelming amount of data available on almost any item available to the public. Sentiment analysis is the pivotal tool, which when is applied to the mammoth data, gives us a holistic review of the body without having to go through the individual reviews manually. This is essential since the online reputation of a product is defined as the collective opinion of the online community on that product.

With the growing popularity of google reviews, Glassdoor, ZipRecruiter, LinkedIn Job Search, and Monster, a job applicant can surely get an empirical review of a company before applying for it. What proves to be a tedious task is to go through every review on all these platforms and formulate an overall sentiment of the public opinion regarding the company. That's where sentiment analysis comes into the picture and we draw a conclusion of the company by analyzing the tone of the reviews (positive or negative) and making this available to the job applicant.

Sentiment Analysis has become an essential part of the analytics that organizations must perform to understand how they are positioned in the market, owing to the rapid migration of customer interactions to digital formats such as forums,

social media posts, comments, reviews, and surveys. Like social media, e-commerce websites, and blogs become more popular among internet users daily, a massive amount of data is generated in the form of reviews, comments, opinions, and emotions from various resources such as websites, text surveys, and events, among others. The data obtained is extremely beneficial to both users and content owners. Users can decide by reading other people's reviews or opinions. Content owners can decide whether their product needs to be improved or if they should make another decision based on customer feedback. The process of analyzing this massive amount of data is both complex and time-consuming.

One would argue that a simple text classification using Natural language processing could work too, but the Lexical-based classification is ineffective. The reason for this is that some evaluations contain no subjective language but communicate a highly negative viewpoint, while others contain negative language but offer positive feedback. For example, "The company could be amazing" or "The client is trying to deliver his best performance". These kinds of reviews mislead the lexical-based technique and reduce the precision of the text classification algorithm. This is where sentiment analysis is anticipated to tackle the issue using novel classification techniques, such as polarity classification (positive or negative) or multiclass ones. We'll talk about different modern sentiment analysis techniques and provide a holistic review of them. We'll also implement and test our approach with others towards the end.

II. LITERATURE SURVEY

Sentiment Analysis is a sub-domain of NLP (Natural Language Processing), which uses texts to mine sentiments and opinions [1]. Though the concept of Sentiment Analysis looks like a simple process, it involves many sub-tasks like subjectivity detection and sarcasm detection, which are very difficult as even humans are not able to cope well with them. [2].

These thoughts and sentiments are highly relevant to our daily lives; therefore, it is necessary to evaluate this user-generated data to automatically monitor public opinion and aid in decision-making [3]. Recent years the number of articles focusing on sentiment analysis and opinion mining has

increased dramatically. As a result, sentiment analysis has become the fastest-growing and most active academic field [4]. The expanding use of the Internet has made the web the most significant and universal information resource. In forums, blogs, wikis, social networks, and other digital resources, millions of individuals express their thoughts and feelings [3]. Hemmatian and Sohrabi [5] surveyed classification strategies for opinion mining and sentiment analysis. In this work, the authors analysed and classified numerous ways for aspect extraction and opinion mining in order to acquire a better understanding of their advantages and disadvantages. Chaturvedi [1] The authors conducted a literature evaluation of handcrafted and automatic models for subjectivity detection. The procedure of generic sentiment analysis is demonstrated in Figure 1. After data is collected and extracted from multiple sites, it is converted to text and then processed using NLP algorithms. Text preprocessing, feature extraction, and feature selection make up the processing step. Several approaches to sentiment analysis (such as machine learning) can be used to perform the classification stage, as described in the next section. Thus, the output can be provided in a variety of forms. There are numerous programmes available to perform sentiment analysis and manage all of these phases.



Fig. 1. Generic Process of Sentiment Analysis

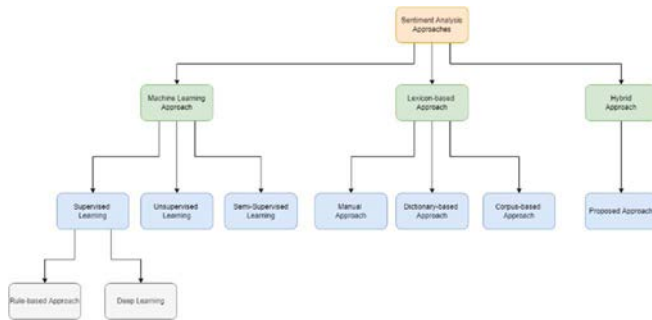


Fig. 2. Sentiment Analysis Approaches

Recently, numerous ways based on conventional or novel word embedding techniques have been developed, and distributions can be bolstered with readily available resources such as sentiment lexicons [6], [7].

Sentiment analysis is a dynamic and active research topic that has several applications. The goal is to enhance the performance of sentiment analysis and to offer solutions to the challenges surrounding this topic. In addition, the application of sentiment analysis to new fields is a major incentive and makes this study even more important. However, selecting the appropriate approach for sentiment analysis is essential. This section is designed to provide an overview of the most prevalent methodologies for sentiment analysis. Existing approaches to sentiment analysis can be classified based on a range of parameters; the majority of research categorises these methods into three categories. Methods based on Machine Learning, Lexicon-Based methods, and Hybrid approaches are examples of such techniques.

A. Machine Learning Approaches

On the basis of train and test datasets, Utilizing machine learning algorithms to determine sentiment polarity (e.g.,

negative, positive, and neutral). Machine learning-based algorithms can extract domain-specific patterns from text, leading to enhanced categorization results. However, these approaches frequently require huge training datasets to perform well. However, a trained classifier on a particular dataset performs less effectively than for another domain [8]. Support Vector Machine is a non-probabilistic classifier capable of separating data linearly or nonlinearly and of handling discrete and continuous variables. In a range of situations, it conducts categorization more precisely than the vast majority of other algorithms due to its solid theoretical foundation. Combining SVM with other algorithms also yielded positive results, as demonstrated by Al Amrani [9]. They suggested a technique that combines SVM and Random Forest are two machine learning algorithms. Their research reveals that hybrid strategies outperform individual algorithms. Recently, the application of deep learning (DL) to sentiment analysis has grown in popularity. DL is an emerging subject of machine learning that offers methods for supervised feature representation learning. Deep learning refers to neural networks with multiple layers of perceptrons that mimic the structure of the human brain. Many neural network models, such as Recurrent Neural Network, are included in DL [10].

Nevertheless, DL is complex and computationally expensive. a number of papers have examined in detail deep learning algorithms for sentiment analysis [11]. This model processes a stream of inputs using memory cells. RNNs are commonly employed in NLP applications such as sentiment analysis due to their capacity to capture and remember information about extended sequences [12]. The rule-based method searches a text for opinion words and classifies them based on the number of positive and negative words. It takes numerous classification concepts into account, such as dictionary polarity, negation words, booster words, idioms, emoticons, mixed attitudes, etc.

B. Lexicon-Based Approaches

One of the two basic methodologies used for sentiment analysis is the lexicon-based (also known as knowledge-based) method. It requires a lexical resource known as an opinion lexicon (a collection of predefined words) that distributes scores to words based on their semantic orientation as negative or positive terms. Based on the semantic orientation of a review's words and phrases, the lexicon-based technique computes the review's sentiment polarity. The "semantic orientation" quantifies subjectivity and opinion inside a text. The final orientation of a document is determined by computing the orientation values of its constituent words. A document is tokenized into single words or microphrases, and emotion values are then assigned to each element from the lexicon. Using formulas or algorithms (such as sum and average) it is possible to assess the overall sentiment of a document. It is feasible to circumvent the issue of various meanings for many words by designing a domain-specific sentiment lexicon or use a lexicon adaptation technique.

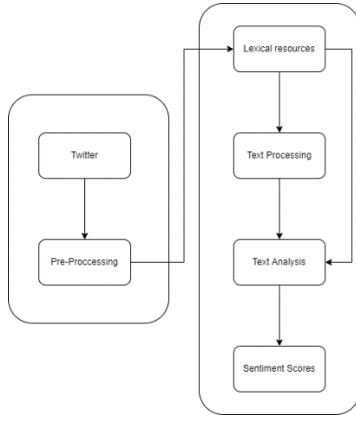


Fig. 3. Lexicon Based Sentiment Analysis

Sanagar and Gupta [13] offered an approach for genre-level sentiment lexicon adaption. This innovative methodology utilises unlabeled data to learn the source and target domain sentiment lexicons, in contrast to current adaption methods that rely on labelled data. As in the work of Sanagar and Gupta, transfer learning algorithms can be used to learn new domain-specific lexicons [14]. The authors proposed a method for unsupervised sentiment lexicon learning that can be applied to new genre-specific domains. After learning polarity seed words from several source domain corpora, genre-level knowledge is transferred to the target domains. In addition, the performance of the lexicon-based method decreases in comparison to the machine learning method when a big training dataset is provided.

C. Hybrid Approaches

The hybrid method combines lexical and machine learning methods. It combines the speed of lexical analysis with the adaptability of machine learning techniques to handle ambiguity and include the context of sentiment words [15]. The major objective of the hybrid strategy is to inherit the high accuracy of machine learning and the consistency of the lexicon-based strategy. The hybrid strategy integrates techniques from the two preceding approaches to circumvent their limits and reap their benefits. The lexical approach scores are used as input features for the sentiment classifier to achieve this objective. Consequently, lexical sentiment plays a vital role in the hybrid technique, which is known to yield improved results. Few models use a mixed approach to sentiment analysis. The bulk of them labelled word polarity for use in sentiment analysis classifiers using lexicon-based approaches.

III. PROPOSED METHODOLOGY

A. CV Parser

CV Parser is the backbone of our application and this starts with the CV extraction module which is split into two sections; CV Training and CV Parsing. The final goal of this module is to extract relevant text from any CV uploaded by the job applicant and sort the different data (of the CV) into appropriate labels. We have generated a corpus of 2000 CVs, manually sorted the resume's various entities, and given them labels. Details like the name of the applicant, his experience in years, his skill set, and his soft skills are all labelled and stored in a simple dictionary. These sorted CVs form the training set and have used supervised machine learning to build the model and make it capable of extracting data from any CV and placing it under the right tag or label. Next comes

CV parsing, here we extract all text information from any CV that the job applicant uploads and sort the applicant's details into different columns. This is made possible by using the Fitz module in pynumpy.

B. Data set

We collected employee reviews from two sources, www.indeed.com and www.seek.com.au. BeautifulSoup, a Python package, was utilised for extracting data from HTML and XML files. We created scripts for Review scrapping from the mentioned websites. In total we scraped approximately 100000 reviews from both the mentioned sources where reviews ranged from various companies in different domains and various job roles, where the data scraped included Company Rating, Review and Pros and Cons. Figure 5 we can see the various attributes that were scraped for the data set In Figure 6 the frequency distribution of Employee Rating is plotted. Figure 7 shows the word cloud of the 100 most common words used for reviews of a given company.

C. Sentiment Analysis on Company Reviews

Sentiment analysis is our next module and we are using that to create a holistic review of different companies (around 20) and give them to the job applicant. This module is divided into 3 sections; Web scraping, pre-processing, and Natural Language Analysis.- We begin by scraping reviews from 2 websites, Indeed and seek.com.

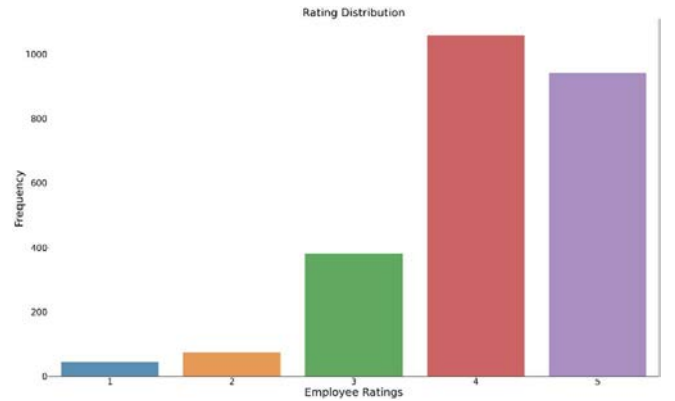
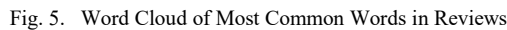


Fig. 4. Frequency Distribution of Employee Rating

Every information in a single review is extracted and stored in a dictionary format. Our training set consisted of 1 Lakh twenty thousand reviews that were considered and analyzed. Next comes preprocessing, now that the data is collected, it is still raw data and requires cleaning. We proceed by starting with the deletion of unnecessary columns, followed by removing contractions in the sentences (neutral words like I've, I'll). After this, we use the python module for language filtering and use only English language reviews. We then perform Tokenization on the remaining text. Tokenization is the act of dividing raw text into smaller chunks by deleting superfluous words while preserving the original sentence's meaning. Next up, we remove punctuation and capital words along with stop words. Lastly, we perform lemmatization of the data that is left. Lemmatization is the process of reducing any given word to its lemma. What we are left with are words that give us an idea of what exactly the review tells us. The data which is left is stored In .pkl format which is used further for analysis Moving on to the last sector of sentiment analysis we have Natural Language Analysis. We use the text blob library in python which is primarily used for

Top 100 Most Common Words



Next, we did Twitter Sentiment Analysis using Twitter API and tweepy python library. First, we used our Twitter API details like API key and Access Token to authenticate with tweepy. Then we made an authentication call using tweepy to retrieve all the latest tweets from the specific user's Twitter account. Next, we clean the data by removing retweets, mentions, hashtags, etc. Then using the textblob library we gave polarity and subjectivity to each tweet that we got from the tweepy call earlier. Then based on the polarity we scored the tweets and categorized them into three sentiments – positive, neutral, and negative. Then, we computed the positivity percentage, drew a graph to provide a comparison of all the tweets, and created a word cloud to display the user's most often used terms.

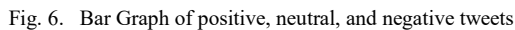


Fig. 9. Scraped Company Review Data set

For sentiment analysis, we have explored techniques like the Machine Learning approach, Lexicon Based approach, and Hybrid approach. Finally, we narrowed it down to the Machine learning approach. The experimentation was performed on i5 8500 (6 Logical Cores), 16 GB RAM 4GB, Video Memory Machine, working with 50000 reviews collected across the mentioned websites. Figure 10 illustrates the plots of e positive and negative word count. Figure 11 give the Confusion Matrix for compared models and it shows the percentage of missed classifications. and Figure 13 represents the ROC Curve for Gradient Boosting Classifier Regression. Table 1 compares various models such as Logistic Regression, Bayes Classifier, Gradient Boost Classifier and Support Vector Classifier used in our experimentation. It can be seen that Gradient Boosting Classifier gave the highest scores in both accuracy and AUC score. However, Logistic Regression and Support Vector Classifier also gave good performance scores even though lesser than GBC.



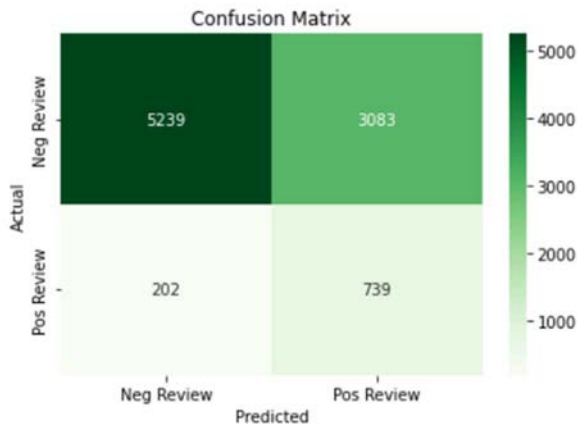


Fig. 11. Confusion Matrix: Gradient Booster Classifier

GBC is a good method to use when you have a mix of feature types, such as categorical, numerical, and so on. Furthermore, when compared to Neural Networks, it has a lower number of hyperparameters to tune. As a result, having the best setting model is more expedient. There is also the option of parallel training. With high performing machine, you can train multiple trees at the same time.

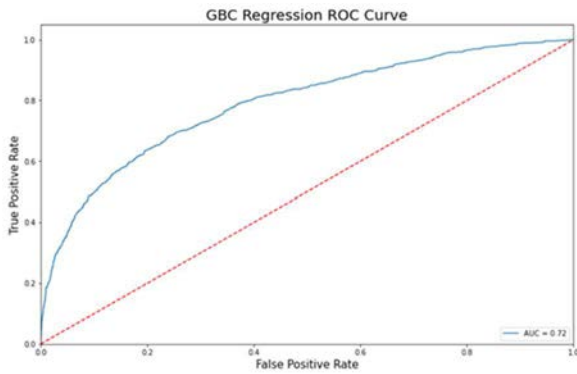


Fig. 12. ROC Curve – Gradient Booster Classifier Regression

TABLE I. COMPARISON OF VARIOUS MACHINE LEARNING ALGORITHMS

Model Name	Logistic Regression		Bayes Classifier		GBC		SVC	
Train Acc.	77.5		78.7		77.1		77.7	
Test Acc.	77.5		78.7		78.0		77.8	
Train AUC	0.79	$s = \pm 0.013$	0.57	$s = \pm 0.012$	0.79	$s = \pm 0.016$	0.79	$s = \pm 0.013$
Test AUC	0.79	$s = \pm 0.001$	0.57	$s = \pm 0.002$	0.81	$s = \pm 0.002$	0.79	$s = \pm 0.002$

V. CONCLUSION AND FUTURE SCOPE

The proposed work can detect the sentiments on the dataset generated by us with sufficiently high accuracy. Gradient Boosting Classifiers have outperformed other classifiers with a classification accuracy of 78.0 percent and an AUC score of 0.82. This model could be leveraged in job recommendation sites and by HR companies during recruitment drives. In future, a neural network model could be

explored for improving performance scores. Also, researchers can include sarcastic detection and multilingual support for consideration.

REFERENCES

- [1] I. Chaturvedi, E. Cambria, R. E. Welsch, and F. Herrera, "Distinguishing between facts and opinions for sentiment analysis: Survey and challenges," *Information Fusion*, vol. 44, pp. 65–77, Nov. 2018, doi: 10.1016/J.INFFUS.2017.12.006.
- [2] A. Valdivia, M. V. Luzón, E. Cambria, and F. Herrera, "Consensus vote models for detecting and filtering neutrality in sentiment analysis," *Information Fusion*, vol. 44, pp. 126–135, Nov. 2018, doi: 10.1016/J.INFFUS.2018.03.007.
- [3] F. J. Ramírez-Tinoco, G. Alor-Hernández, J. L. Sánchez-Cervantes, B. A. Olivares-Zepahua, and L. Rodríguez-Mazahua, "A brief review on the use of sentiment analysis approaches in social networks," in *Advances in Intelligent Systems and Computing*, 2018, vol. 688, pp. 263–273. doi: 10.1007/978-3-319-69341-5_24.
- [4] M. v. Mäntylä, D. Graziotin, and M. Kuuttila, "The evolution of sentiment analysis—A review of research topics, venues, and top cited papers," *Comput Sci Rev*, vol. 27, pp. 16–32, Feb. 2018, doi: 10.1016/J.COSREV.2017.10.002.
- [5] F. Hemmatian and M. K. Sohrabi, "A survey on classification techniques for opinion mining and sentiment analysis," *Artif Intell Rev*, vol. 52, no. 3, pp. 1495–1545, Oct. 2019, doi: 10.1007/s10462-017-9599-6.
- [6] M. Kasri, M. Birjali, and A. Beni-Hssane, "Word2Sent: A new learning sentiment-embedding model with low dimension for sentence level sentiment classification," *Concurr Comput*, vol. 33, no. 9, May 2021, doi: 10.1002/cpe.6149.
- [7] A. Yadav and D. K. Vishwakarma, "Sentiment analysis using deep learning architectures: a review," *Artif Intell Rev*, vol. 53, no. 6, pp. 4335–4385, Aug. 2020, doi: 10.1007/s10462-019-09794-5.
- [8] G. Yoo and J. Nam, "A Hybrid Approach to Sentiment Analysis Enhanced by Sentiment Lexicons and Polarity Shifting Devices," May 2018. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-01795217>
- [9] Y. al Amrani, M. Lazaar, and K. E. el Kadiri, "Random Forest and Support Vector Machine based Hybrid Approach to Sentiment Analysis," *Procedia Comput Sci*, vol. 127, pp. 511–520, 2018, doi: 10.1016/j.procs.2018.01.150.
- [10] W. Li, F. Qi, M. Tang, and Z. Yu, "Bidirectional LSTM with self-attention mechanism and multi-channel features for sentiment classification," *Neurocomputing*, vol. 387, pp. 63–77, Apr. 2020, doi: 10.1016/j.neucom.2020.01.006.
- [11] S. Sohagiri, D. Wang, A. Pomeranets, and T. M. Khoshgoftaar, "Big Data: Deep Learning for financial sentiment analysis," *J Big Data*, vol. 5, no. 1, p. 3, Dec. 2018, doi: 10.1186/s40537-017-0111-6.
- [12] A. Aziz Sharfuddin, Md. Nafis Tihami, and Md. Saiful Islam, "A Deep Recurrent Neural Network

with BiLSTM model for Sentiment Classification,” in *2018 International Conference on Bangla Speech and Language Processing (ICBSLP)*, Sep. 2018, pp. 1–4. doi: 10.1109/ICBSLP.2018.8554396.

- [13] S. Sanagar and D. Gupta, “Automated genre-based multi-domain sentiment lexicon adaptation using unlabeled data,” *Journal of Intelligent & Fuzzy Systems*, vol. 38, no. 5, pp. 6223–6234, May 2020, doi: 10.3233/JIFS-179704.

- [14] S. Sanagar and D. Gupta, “Unsupervised Genre-Based Multidomain Sentiment Lexicon Learning

Using Corpus-Generated Polarity Seed Words,” *IEEE Access*, vol. 8, pp. 118050–118071, 2020, doi: 10.1109/ACCESS.2020.3005242.

- [15] I. Gupta and N. Joshi, “Enhanced Twitter Sentiment Analysis Using Hybrid Approach and by Accounting Local Contextual Semantic,” *Journal of Intelligent Systems*, vol. 29, no. 1, pp. 1611–1625, Sep. 2019, doi: 10.1515/jisys-2019-0106.