# Using Machine Learning Approach for Sentiment Analysis of Pashto Text

**Javed Iqbal[1*], Jamal Abdul Nasir[1], Muhammad Zubair Asghar[1]**
javedmarwat7@gmail.com, jamalnasir@gu.edu.pk, mzubair@gu.edu.pk
[1] Institute of Computing and Information Technology, Gomal University, Dera Ismail Khan, Pakistan 29050
*Corresponding Author: javedmarwat7@gmail.com

*Abstract: The internet occupies the glob and the people using different sources for sharing their views, opinions and sentiments about the company products, politics, and situation happened in daily life time to time. The sentiment analysis of rich resources languages like English, French and Chinese has gained highly position while the resource poor languages like Pashto was totally ignored in the area of sentiment analysis by researchers. The Pashto dataset has collected from various Pashto online news channels and blogs of social media about the different categories and these dataset are manually annotated by the Pashto language experts. The dataset are categorized into positive and negative sentiments. After processing the Pashto text and then apply the different machine learning techniques with different feature extraction, achieved the satisfactory results. Initially we use four machine learning algorithms Support vector Machine (SVM), Naïve Bayesian (NB), KNN and Decision Tree (DT) for sentiment analysis of Pashto text with their feature extractions. The Naïve Bayesian technique with TF-IDF shows the best performance with accuracy 82%, Precision 80 %, Recall 80% and F1-Score 80% among the rest of machine learning classifiers. The performance of machine learning classifiers can be enhanced by increasing the dataset and other classifiers of machine learning also be considered.*
*Keywords: Pashto sentiment analysis, Support Vector Machine, Naïve Bayesian, Decision Tree, K-Nearest Neighbor (KNN), feature extraction.*

## 1. Introduction

The people communicate their ideas about the object and events either positive or negative sentiments and massive data are produces by the social medial users on daily basis [1]. The opinion of the public about the Products, events, politics, sports, objects in the shape of blogs and reviews either may be positive or negative which help in decision making process [2]. Sentiment analysis classified the public opinions and judgments in the way of positive and negative sentiments [3]. Online news and social media is one of the rich source of sentiments, where public share their idea through comments, blogs and reviews [4].

Pashto is an Indo-European language and about 80-million people around the world speak and understand the Pashto [5]. Pashto is an official language of Afghanistan and Khyber Pakhtunkhwa Pakistani state but inspite of all such facts Pashto has not gain the attentions of researchers [6]. Actually, sentiment analysis is the procedural studies of human behavior, attitude, and ideas and can be categorized into three extents such as sentence level, aspect level and document level.

The documents will be consider as positive sentiment where the number of positive sentences are greater than negative sentences and the document termed as negative when the number of negative sentences are greater than positive sentences in document. In sentence level, the sentence consider as positive if the number of positive words are more in the sentence and the whole sentence cover the positive sentiment and the sentence termed as negative if the words in sentence are more negative and provide the negative meaning. The aspect level sentiment analysis classification are accomplished by examining the features of text already recognize their

126

nature either positive or negative sentiments [8].The researchers performed their research on rich resources languages like English, Arabic, Chinese and Urdu in the area of sentiment analysis but no attention given to the poor resources languages like Pashto [9].

Processing of Pashto language is dense due to some facts because Pashto is customized style of Arabic language; Pashto has more than 20-dialects, unbound word sequences and complicated morphological shape.

Pashto Sentence level sentiment analysis is concentrated in this research. Various step involved in the Pashto sentiment analysis procedure. First of all, the Pashto dataset is collected from various online Pashto news channels and online blogs. Then these Pashto sentences are classified into positive and negative sentiments by using various classifiers.

Machine learning classifier tools are used for the classification of Pashto sentences by using the Python's environment. For classification of Pashto sentences four popular supervised machine learning classifiers are used namely, Support Vector Machine (SVM), Naïve Bays (NB), Decision Tree (DT) and K-Nearest Neighbor (KNN). Several experiments are performed to improve the accuracy, Precision, Recall and F1-score that are discussed in section-3 and results are discussed in section-4.

The rest of paper arranged as: Section-2 composed of related work in sentiment Analysis. The material and methodology are discussed in section-3. Experiments and results obtained in section-4 and at end conclusion and future work are discussed in section-5.

## 2. Related Work

Several researches are performed on sentence level sentiment classification and outstanding survey presented by the researches in detail with supervised machine learning techniques in rich resources languages not in Pashto language.

Supervised machine learning algorithms on the huge dataset which extracted from the twitter and classify such data as positive or negative [12]. They proposed the SVM, NB and KNN techniques for the classification of twitter dataset. The proposed classifier techniques show that Naïve Bayesian method show better result than SVM and KNN. They proposed that the classification can be improved by increasing the dataset for sentiments. Kundi et al [13], developed the lexicon-based framework which classify the tweets into positive, negative or neutral categories and also predict the score of slangs used in the tweets. They proposed framework for sentiment classification consist of major six modules i.e. tweets capturing module, pre-processing module, lexicon module, subjective text identification, Additional knowledge module. The accuracy achieved from the proposed system about 92% on binary classifications and 87% on multi-class classification. To extend the framework for testing of the datasets and the system will more improve for precision in negative cases and recall for neutral cases. Arabic reviews are collected from YouTube for sentiment analysis [14]. They applied various machine learning algorithms like NB, SVM, Decision Tree and function on the Arabic corpus. The result achieved from the machine learning algorithm shows that the performance of NB is better than other machine learning algorithms and the accuracy of NB about 91.5 %. Compared the various machine learning techniques used by the Arabic language for sentiment analysis on twitter data and they proposed three machine learning technique Decision Tree (DT), Naïve Bayesian (NB) and Support Vector Machine (SVM) for sentiment analysis for Arabic twitter data [15]. The obtained result 78% of F1-Score show that Decision Tree performance better than NB & SVM. In future, to extend the Arabic lexicon, sort out more features, using the advance machine learning technique to minimize the time and memory

space and combining the machine learning techniques for better performance. Urdu readers classified documents and categorized the Urdu applications like email, spam detection, web classification, ontology mapping etc for better understanding the Urdu documents easily. They proposed the Naïve Bayes classification technique for the classification of Urdu documents and satisfactory and an efficient result obtained from the proposed system [16]. Sentiment analysis on sentence level of Urdu language and the data is collected from different online blogs such sentences are manually annotated and finally classified the sentences as positive, negative or neutral [17] and proposed the SVM, Decision Tree and KNN for the sentiment classification of Urdu language. The result achieved from the proposed system is satisfactory and performance of KNN is better one. Improve the sentiment classifier more for better accuracy and test the result on other statistical methods for best performance. Sentiment analysis Sihdhi language structure and analyzed the Sindhi corpus dataset through sentiment analysis tools. The supervised machine learning techniques SVM and KNN are suggested for the analysis of Sindhi sentiment analysis corpus dataset. The satisfactory results are obtained from the supervised machine learning techniques on the Sindhi corpus dataset. To improve the results in future, increase the sentiment analysis of Sindhi language for organization [18]. Pashto is an Arabic script-based language having complex Morphological structure, Kamal et al. [19], proposed the various supervised machine learning classifiers including (SVM), Naïve Bayesian (NB), Logic Based Machine Translation (MLT) and j48 on lexical features with a corpus of 600 sentences acquired from social media sites like Facebook, Twitter and news website like VOA and BBC. The results obtained from the different machine classifiers applied on the Pashto sentiment show that SVM has obtained the highest accuracy rate. However, limited Pashto corpus and unigram based approach was adopted but by improving the efficiency introducing the n-gram approach, enhance the Pashto corpus and availability of Pashto corpus to everyone in future.

## 3. Methodology

The phenomenon adopted in this research composed of several steps and each step is further divided into sub-phases.



**Fig 3.1 System Architecture for Pashto Sentiment Analysis**

### 3.1. Data Collection

The first and compulsory step for sentiment analysis of text is the collection of dataset. The languages having rich resources like English, French, Arabic, Chinese etc. datasets are publically available but the poor resources languages like Pashto, generally no dataset available for sentiment analysis on web. For this propose, the required Pashto dataset are extracted from various sources like online Pashto news channels, Pashto websites, Pashto media resources etc. These Pashto blogs belongs to various categories like politics, sports, education, entertainments, religion, current affairs etc.

128

### 3.2. Annotation of Dataset

The Pashto blogs are manually extracted from various online news channels. These sentences are annotated by Pashto literature experts with positive and negative sentiments. In this research the sentences are categorized into categories positive and negative. The annotators classified the sentences according to the following annotation rules.

- The sentence is marked as positive if the whole sentence presents the positive sentiment and when the sentence composed of positive and neutral words then marked as positive sentence [20].

- The sentence is marked as negative if the whole sentence expressed the negative sentiment and when the sentence composed more words of negative than positive and neutral then it may be marked as negative [21].

- When the sentence expressed the dispute concept then sentence marked as negative [22].

### 3.3. Data Preprocessing

Preprocessing procedure is the basic and necessary phase to accomplish the task of NLP in which cleaning the Pashto text extracted from various media resources. By removing the irrelevant elements in Pashto text such as exceptionable tags (%, $, ^, ? <,>, @,# etc), numeric digits, URL, email, punctuation etc to enhancing the correctness of models.

The preprocessing technique of Pashto text is very important to perform the NLP assignment. By improving the accuracy of models, the un-necessary items in Pashto text like undesired symbols ($,?, @,# etc), numbers & digits, URLs, Tokenization, email addresses, punctuation marks etc. must be removed from the text and text must be transformed in proper format before processing. However, some other preprocessing phase should be performed to enhance the accuracy and capability of the models for Pashto text.

#### i. Noise Reduction

By eliminating the irrelevant signs and labels like %, $, #, @, <, > etc., from the Pashto dataset to construct the quality and effective level text before further processing.

#### ii. Removing the Stop Words

Eliminating those words which are mostly used in the sentences and have no participation in sentiment analysis procedure to built accurate dataset for classification. The mostly stop words used in Pashto text are ("وي", "به", "هم", "په", "کښي", "ده", "د"), so eliminate these words before processing.

#### iii. Tokenization

Split the long and more complicated sentences into tiny sentences and further these sentences splitting into words called token. The bellow table shows the conversion of Pashto sentence into related token.

<div dir="rtl">

افواج پاکستان په نره د ترهه ګرو مقابله کړي ده۔

</div>

This Pashto sentence will be tokenized as:

**Table 3.1: Tokenization of Pashto Sentence**

| Sentence in Pashto Text | Tokenized Text | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | t1 | t2 | t3 | t4 | t5 | t6 | t7 | t8 | t9 |
| افواج پاکستان په نره د ترهه ګرو مقابله کړي ده۔ | افواج | پاکستان | په | نره | د | ترهه ګرو | مقابله | کړي | ده |

#### iv. Stemming

Removing the prefix and postfix words and convert the word into root form. No straightforward rule available for Pashto stemming technique, so Pashto stemming procedure is very complex due to its morphological shape and manually performed.

Table-3 shows some stemming words in Pashto text.

### Table: 3.2: Stemming of Pashto Words

| Pashto Word | Stem |
|---|---|
| ببخونده | خوند |
| ببد ماغه | د ماغ |
| خبرونه | خبر |

### v. Lemmatization

Transformation of words into its actual root form rather breaks the word. No specific procedure available for Pashto text. So lemmatization process performed manually upto some limit. Table-4 presents the Pashto lemmatization.

### Tabel-3.3: Lemmatization of Pashto words

| Pashto Word | Lemmatization |
|---|---|
| توده | تود |
| زرګي | زره |

## 4. Classification Models

Various machine learning models are proposed to obtain the effectual results from the Pashto dataset for the sentiment analysis. In this research, four leading models of machine learning Support Vector Machine (SVM), Naïve Bayesian (NB), K-Nearest Neighbor (KNN) and Decision Tree (DT) are proposed for sentiment analysis of Pashto text. After implementation of these models, we proposed the best model with feature selection on the basis of performance on Pashto dataset.

### 4.1. Splitting Module

By creating the benchmark dataset of Pashto language, manually annotated by the Pashto literature human experts labeled as P & N. Experts of Pashto language label each sentence as positive (P) or negative (N) sentiment categories. The dataset contain about 5500 Pashto news reviews, 2225 sentences of dataset are positive while 2225 composed of negative sentences, as show in table 4.1.

### Table 4.1: Pashto Statistical Dataset

| Attributes | Size |
|---|---|
| Positive News | 2225 |
| Negatives | 2225 |
| Total Pashto News | 5500 |
| Tokens | 7584352 |
| Avg: Tokens | 97.7 |

### 4.2. Train Set

By trained the proposed model, train set is used, which balanced the parameter on machine learning models. 80% dataset is used for the proposed models to evaluate the results.

### Table 4.2 Shows the Dataset of Trained Model

| Pashto Text | Sentiments |
|---|---|
| انشا الله د پاکستان مستقبل روښنانه دے | Positive |
| د پنجاب کښنی کارروائي دوران دوه تر هګر ګرفتار کرے | Negative |
| پاکستان یو پرامن ګاونډيتوب باندي يقين لري | Positive |

| | |
|---|---|
| د جهاز غورځيدو پيښنه کښې يوسل لس تنان مړه شويدي | Negative |

## 4.3 Test Set

The test set is used to provide the neutral assessment to fit the final model on the training data. Test set is used to corroborate the factual anticipating capability of the machine learning models. 20% Pashto dataset used for testing purpose in this research.

**Table 4.3 Represents The Sample Of Test Dataset.**

| Pashto Text | Sentiment |
|---|---|
| انشا الله د پاکستان مستقبل روښانه دے | Positive |
| د پنجاب کښې کارروائي دوران دوه تر هګر ګرفتار کړے | Negative |
| پاکستان يو پرامن ګاونډيتوب باندي يقين لري | Positive |
| د جهاز غورځيدو پيښنه کښې يوسل لس تنان مړه شويدي | Negative |

## 5. Performance Evaluation Parameters

The parameters used for the assessment prophecy are accuracy, precision, recall and F1-score. To obtain the best accuracy from the proposed Pashto dataset is our major goals along with solid other three parameter. The

Accuracy, precision, recall and F1- score is formulated as:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Where

TP= True Positive
TN=True Negative
FP= False Positive
FN=False Negative

$$Precsion = \frac{TP}{TP + FP}$$
$$Recall = \frac{TP}{TP + FN}$$
$$F1\ score = 2 * \frac{precision * recall}{Precision + Recall}$$

To compute the TP, TN, FP, FN for binary classification (positive & negative) can be determining as:

➢ True Positive (TP): For a Positive classification, TP is the number sentences that actually belong to Positive/Negative class and are also correctly predicted as Positive /Negative by classifier.

➢ True Negative (TN): TN is the number of sentences that do not belong to Positive/Negative category and are also not predicted by a classifier.

➢ False Positive (FP): FP represents the number of sentences whose actual labels do not belong to Positive category but are predicted as Positive by a classifier, and vice versa.

➢ False Negative (FN): FN represents the number of sentences whose actual labels belong to Positive category but are predicted as Negative by a classifier, and vice versa.

### 5.1. Performance of Machine Learning Models Using TF-IDF Embedding

Different results achieved through experiments by using the TF-IDF embedding for machine learning classifiers with accuracy, precision, recall and F1-score as shown in table 5.1.

**Table-5.1:  Results Using TF-IDF Embedding Technique for Sentiment Classification**

| S. No. | Machine Learning Classifiers | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| 1 | Decision Tree | 74% | 71% | 71% | 71% |
| 2 | Naïve Bayesian | 82% | 80% | 80% | 80% |
| 3 | KNN | 69% | 71% | 73% | 69% |
| 4 | Support vector Machine | 69% | 70% | 71% | 70% |

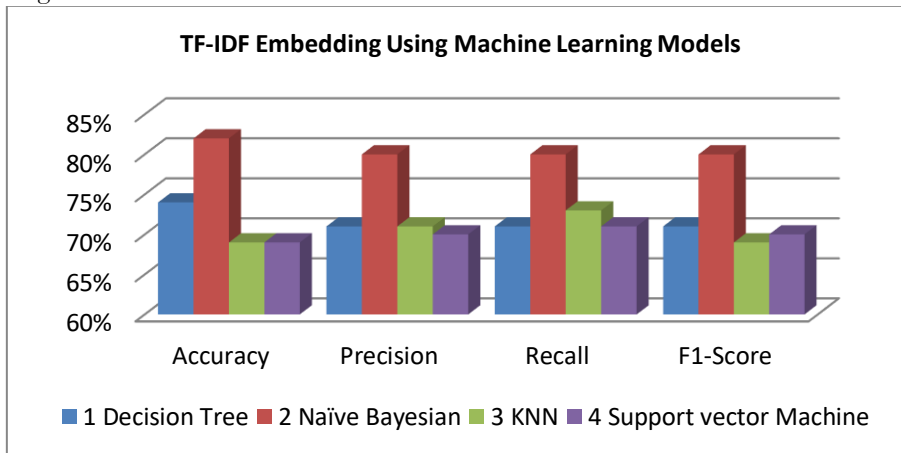The graphical representation of accuracy, precision, recall and F1-Score machine learning model by using TF-IDF feature extraction as:



**Fig: 5.1 Overall Results Comparison of Machine Learning Model using TF- IDF Embedding**

The Fig: 5.1 show the overall results comparison different machine learning classifiers by using TF-IDF embedding technique in which logistic regression have the highest accuracy rate 83%, Precision 87% and F1-Score while Recall rate 80% of Naïve Bayesian is highest.

**5.2. Performance of Machine Learning Models Using Count Vector Embedding**

Various types of results obtained through experiments by using the count vector embedding for machine learning classifiers with accuracy, precision, recall and F1-score as shown in table 5.2.

**Table 5.2 Results using Count Vector Embedding Technique**

| Machine Learning Classifiers | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Decision Tree | 78% | 76% | 76% | 76% |
| Naïve Bayes | 65% | 67% | 66% | 65% |
| KNN | 67% | 68% | 70% | 67% |
| SVM | 64% | 32% | 50% | 39% |

   The graphical representation of accuracy, precision, recall and F1-Score machine learning model by using count vector feature extraction as:
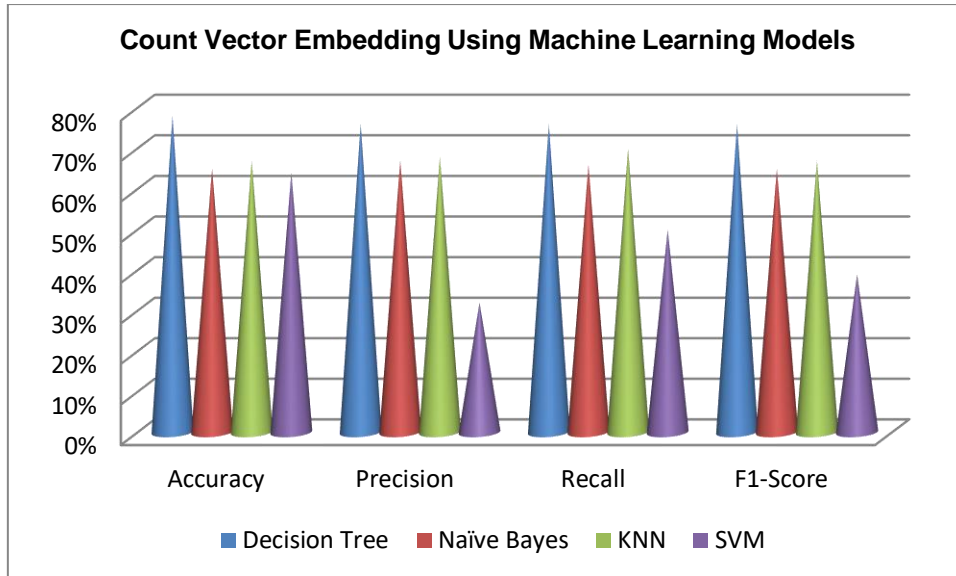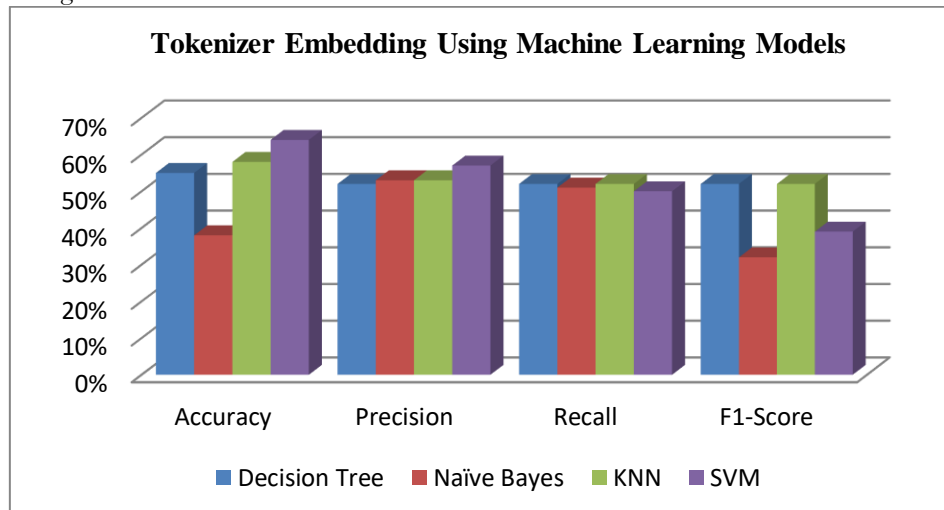
**Fig: 5.2 Overall Results Comparison of Machine Learning Model using Count Vector Embedding**

### 5.3. Performance of Machine Learning Models Using Tokenizer

Different results obtained through experiments by using the count vector embedding for machine learning classifiers with accuracy, precision, recall and F1-score as shown in table 5.3.

**Table 5.3 Results using Tokenizer Embedding Technique**

| Machine Learning Classifiers | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Decision Tree | 55 % | 52 % | 52 % | 52 % |
| Naïve Bayes | 38 % | 53 % | 51 % | 32 % |
| KNN | 58 % | 53 % | 52 % | 52 % |
| SVM | 64 % | 57 % | 50 % | 39 % |

The graphical representation of accuracy, precision, recall and F1-Score machine learning model by using Tokenizer feature extraction as:
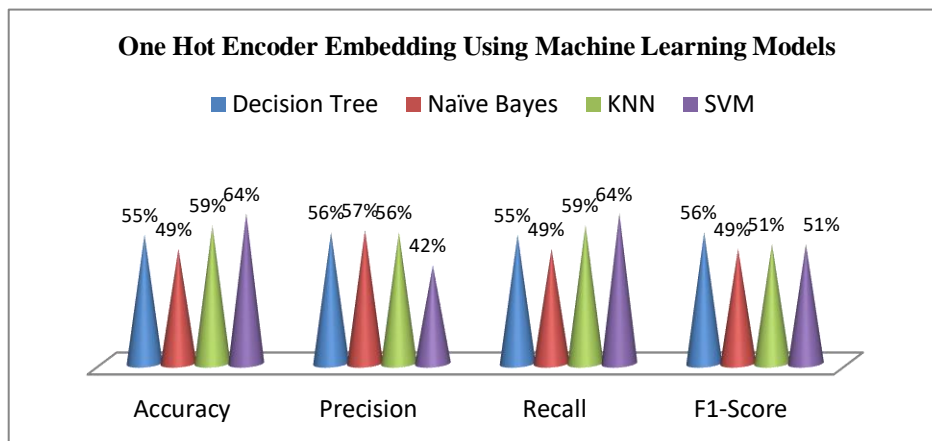


**Fig- 5.3 Tokenizer Using Machine Learning Models**
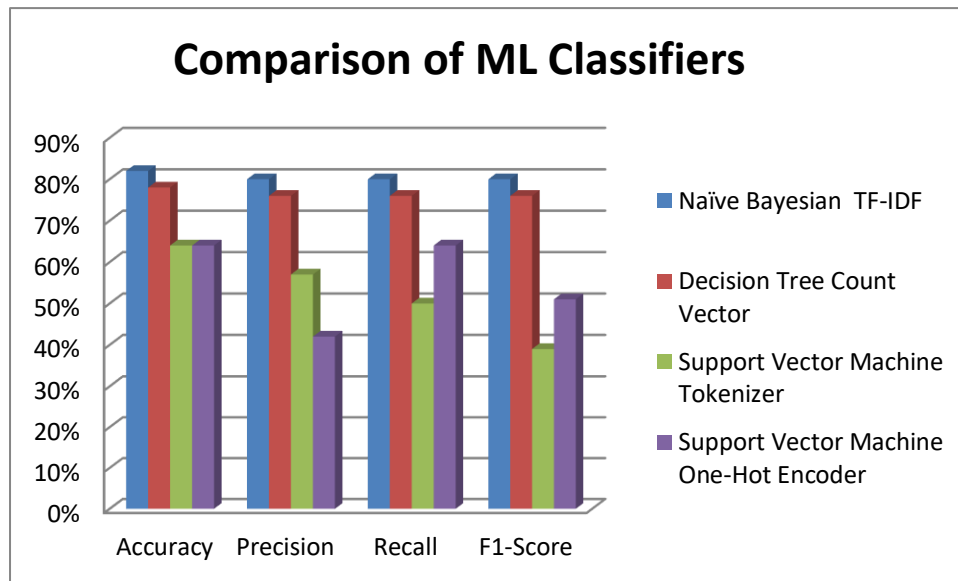
The Fig: 5.3 show the overall results comparison different machine learning classifiers by using Tokenizer embedding technique in which Gradient Booster has the highest accuracy rate 67%, Precision 67%, Recall 56% and F1-Score 52%.

133

**Copyrights @ Roman Science Publications**          **Vol. 7 No. 1 June, 2022, Netherland**
**International Journal of Applied Engineering Research**

**5.4. Performance of Different Machine Learning Models Using One- Hot Encoder for sentiment classification of Pashto Text.**

Different results achieved through experiments by using the TF-IDF embedding for machine learning classifiers with accuracy, precision, recall and F1-score as shown in table 5.4.

**Table 5.4 Results using One-Hot Encoder Embedding Technique**

| Machine Learning Classifiers | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Decision Tree | 55% | 56% | 55% | 56% |
| Naïve Bayes | 49% | 57% | 49% | 49% |
| KNN | 59% | 56% | 59% | 51% |
| SVM | 64% | 42% | 64% | 51% |



**Fig 5.4: One Hot Encoder Using Machine Learning Models**

The Fig: 5.4 show the overall results comparison different machine learning classifiers by using One Hot Encoder embedding technique in which Gradient Booster has the highest accuracy rate 66%, Precision 63%, Recall 66% and F1-Score 60%.

**5.5. Recommendation of Machine learning techniques with best performance for sentiment classification of Pashto text**

By using the same data preprocessing for the proposed methodology of machine learning classifiers by using the different Embedding Models for sentiment classification of Pashto text the Naïve Bayesian classifier with TF-IDF is embedding, Decision Tree with Count Vector and Support Vector Machine show the best performance with Tokenizer and One-Hot Encoder embedding.

**Table: 5.5 : Comparison of Machine Learning Classifiers**

| Machine Learning Classifiers | Embedding Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| Naïve Bayesian | TF-IDF | 82% | 80% | 80% | 80% |
| **Decision Tree** | **Count Vector** | 78% | 76% | 76% | 76% |
| Support Vector Machine | Tokenizer | 64 % | 57 % | 50 % | 39 % |
| Support Vector Machine | One-Hot Encoder | 64% | 42% | 64% | 51% |

**Figure**



**.5.5.**

## Comparison of ML Classifiers

Comparison of Machine Learning Classifiers

### 5.5.1. Overall Performance and Efficiency of Machine Learning Technique for Sentiment Classification

The Naïve Bayesian classifier with TF-IDF embedding show the best performance & efficiency among the all Machine learning classifiers and their feature extractions

### 6. Conclusion and Future Work

Each of the four machine learning classifiers are executed on the Pashto dataset and all the divulge results are deeply investigated to enhance the results and identify the best machine learning classifier with features that obtain the better results than other regarding the accuracy, precision, recall, and F1-Score.

Limited research studies have been observed in the Pashto language in sentiment analysis zone. In this research, high classification has been obtained for Pashto sentiment analysis using various machine learning classifiers. Various result achieved from the experiments regarding the accuracy, precision, recall, and F1-Score.

Various experiments performed on the Pashto dataset through machine learning classifiers with their feature extractions, the Naïve Bayesian with TF-IDF embedding achieved the highest accuracy 82%, Precision 80%, Recall 80% and F1-Score 80%. The performance of decision tree with count vector embedding achieved the second highest accuracy 78%, Precision 76%, Recall 76% and F1-Score 76%. The SVM show the average performance and KNN show the poor performance among the above machine learning classifiers. This work provides the opportunity for researchers to investigate the poor resource languages. The limitation of this study is that only positive and negative classes are used and neutral class is not included. In future we will enhance the dataset and include the neutral class for sentiment analysis. However, we will implement rest of the machine learning classifiers on Pashto dataset to improve the results.

### 7. References:

1. Sohail, M., Imran, A., Rehman, H. U., & Salman, M. (2020, January). Anti-Social Behavior Detection in Urdu Language Posts of Social Media. In *2020 3rd International Conference on Computing, Mathematics and Engineering Technologies (iCoMET)* (pp. 1-7). IEEE.
2. Mukhtar, N., & Khan, M. A. (2018). Urdu sentiment analysis using supervised machine learning approach. *International Journal of Pattern Recognition and Artificial Intelligence*, *32*(02), 1851001.
3. Nemes, L., & Kiss, A. (2021). Social media sentiment analysis based on COVID-19. *Journal of Information and Telecommunication*, *5*(1), 1-15.
4. Ahmad, S., Asghar, M. Z., Alotaibi, F. M., & Awan, I. (2019). Detection and classification of social media-based extremist affiliations using sentiment analysis techniques. *Human-centric Computing and Information Sciences*, *9*(1), 24.
5. *Khan, S., Ali, H., Ullah, Z., Minallah, N., Maqsood, S., & Hafeez, A. (2019). KNN and ANN-based recognition of handwritten Pashto letters using zoning features. arXiv preprint arXiv:1904.03391.*
6. *Khan, S., Ali, H., Ullah, Z., Minallah, N., Maqsood, S., & Hafeez, A. (2019). KNN and ANN-based recognition of handwritten Pashto letters using zoning features. arXiv preprint arXiv:1904.03391.*
7. W. Medhat, A. Hassan and H. Korashy, Sentiment analysis algorithms and applications: A survey, Ain Shams Eng. J. 5 (2014) 1093–1113.
8. Tang, D. (2015, February). Sentiment-specific representation learning for document-level sentiment analysis. In *Proceedings of the eighth ACM international conference on web search and data mining* (pp. 447-452).
9. Kamal, U., Siddiqi, I., Afzal, H., & Rahman, A. U. (2016, November). Pashto Sentiment Analysis Using Lexical Features. In Proceedings of the Mediterranean Conference on Pattern Recognition and Artificial Intelligence (pp. 121-124). ACM.
10. Zada, B., & Ullah, R. (2020). Pashto isolated digits recognition using deep Convolutional neural network. *Heliyon*, *6*(2), e03372.
11. Basiri, M. E., & Kabiri, A. (2017, April). Sentence-level sentiment analysis in Persian. In *2017 3rd International Conference on Pattern Recognition and Image Analysis (IPRIA)* (pp. 84-89). IEEE.
12. Gautam, G., & Yadav, D. (2014, August). Sentiment analysis of twitter data using machine learning approaches and semantic analysis. In *2014 Seventh international conference on contemporary computing (IC3)* (pp. 437-442). IEEE.
13. Kundi, F. M., Khan, A., Ahmad, S., & Asghar, M. Z. (2014). Lexicon-based sentiment analysis in the social web. *Journal of Basic and Applied Scientific Research*, *4*(6), 238-48.
14. Elawady, M. (2015). Sparse coral classification using deep convolutional neural networks. *arXiv preprint arXiv:1511.09067.*
15. Altawaier, M., & Tiun, S. (2016). Comparison of machine learning approaches on arabic twitter sentiment analysis. *International Journal on Advanced Science, Engineering and Information Technology*, *6*(6), 1067-1073.
16. Asghar, M. Z., Sattar, A., Khan, A., Ali, A., Masud Kundi, F., & Ahmad, S. (2019). Creating sentiment lexicon for sentiment analysis in Urdu: The case of a resource-poor language. *Expert Systems*, *36*(3), e12397.
17. Mukhtar, N., Khan, M. A., & Chiragh, N. (2018). Lexicon-based approach outperforms Supervised Machine Learning approach for Urdu Sentiment Analysis in multiple domains. *Telematics and Informatics*, *35*(8), 2173-2183.
18. Ali, M., & Wagan, A. I. (2017). Sentiment summerization and analysis of Sindhi text. *Int. J. Adv. Comput. Sci. Appl*, *8*(10), 296-300.

19. Kamal, U., Siddiqi, I., Afzal, H., & Rahman, A. U. (2016, November). Pashto Sentiment Analysis Using Lexical Features. In *Proceedings of the Mediterranean Conference on Pattern Recognition and Artificial Intelligence* (pp. 121-124).

20. Maynard, D. G., & Bontcheva, K. (2016, May). Challenges of evaluating sentiment analysis tools on social media. In *Proceedings of the tenth international conference on language resources and evaluation (LREC 2016)* (pp. 1142-1148). LREC.

21. Ganapathibhotla, M., & Liu, B. (2008, August). Mining opinions in comparative sentences. In *Proceedings of the 22nd international conference on computational linguistics (Coling 2008)* (pp. 241-248).

22. Maynard, D. G., & Bontcheva, K. (2016, May). Challenges of evaluating sentiment analysis tools on social media. In *Proceedings of the tenth international conference on language resources and evaluation (LREC 2016)* (pp. 1142-1148). LREC.