

Classification of Online Toxic Comments Using Machine Learning

¹Rayudu Hinduja, ²P.Rama Krishna

^{1,2}Dept. of CSE, KIET, Kakinada, AP, India

Abstract

Conversational toxicity is a problem that might drive people to cease truly expressing themselves and seeking out other people's opinions out of fear of being attacked or harassed. The purpose of this research is to employ natural language processing (NLP) techniques to detect toxicity in writing, which might be used to alert people before transmitting potentially toxic informational messages. Natural language processing (NLP) is a part of machine learning that enables computers to comprehend natural language. Understanding, analysing, manipulating, and maybe producing human data language with the help of a machine are all possibilities. Natural Language Processing (NLP) is a type of artificial intelligence that allows machines to understand and interpret human language instead of simply reading it. Machines can understand written or spoken text and execute tasks such as speech recognition, sentiment analysis, text classification, and automatic text summarization using natural language processing (NLP).

I. Introduction

In the early days of the Internet, people communicated only by e-mail that encountered spam. At the time, classifying emails as good or bad, whether or not spam was a challenge. Internet communications and data flows have changed significantly over time, especially since the development of social networks. With the advent of social media, classifying content as "good" or "bad" has become more important than ever to prevent social harm and prevent people from engaging in antisocial behavior.

A. Purpose

Every day, vast amounts of data are released by social media networks. This massive volume of data is having a big impact on the quality of human existence. However, because there is so much toxicity on the Internet, it can be harmful. Because toxic comments limit people's ability to express themselves and have different opinions, as a result of the negative, there are no positive conversations on social networks. As a result, detecting and restricting antisocial behaviour in online discussion forums is a pressing requirement. These toxic comments might be offensive, menacing, or disgusting. Our goal is to identify these harmful remarks.

B. Project Overview

To classify noxious comments, this study will use different methods of machine learning. To handle the problem of text categorization, Some of the methods we use are logistic entertainment, random forest, SVM classifier, multi-navigation database, and XGBoost classifier. As a result, we maintain a data set using six machine learning algorithms, evaluate their accuracy, and compare them.. We choose the model with the highest accuracy and forecast the toxicity based on unseen data based on its accuracy level.

II. Literature Survey

A. Existing Systems

1. Convolutional Neural Networks (CNN)

CNN structures for categorization are Neural Network structures that can adapt to many stages. Each of these levels has a number of layers. The Embedding Layer is a CNN component that is used to solve problems with text categorization. The translation of the embedded text into CNN's description format is a function of the input layer. This approach converts each word in a text document into a standard compressed vector.

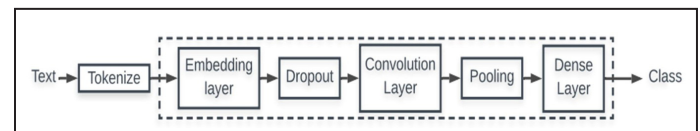


Fig: A Convolutional Neural Network

2. LSTM (Long Short Term Memory)

Long-term memory is abbreviated as LSTM. According to memory, LSDM is a kind of continuous neural network that works better than conventional continuous neural networks. LSTM is certainly better when it comes to memorizing specific patterns. Like other NNs, LSTM can have multiple hidden layers, storing the relevant data in each cell as it passes through each level, and rejecting inappropriate data. LSTM has a memory function that allows you to store data sequences. It also has another function: it works to delete unwanted data, and since we all know that text data has a lot of unnecessary data, LSTM can delete it to save time and money on calculations. As a result, LSTM's ability to memorize data sequences while deleting additional data makes it an effective text classification tool.

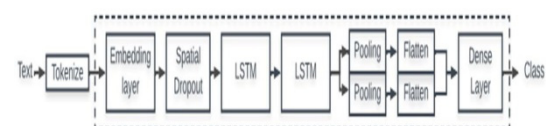


Fig:A Long Short-Term Memory Network.

B. Proposed System

We will use five machine learning algorithms to classify the data the online harmful comments: logistic regression, random forest, SVM classifier, Naive Bayes, and XGBoost classifier. We'll use logistic regression to determine whether or not a comment is harmful because it either belongs to the poisonous group or does not. SVM classifiers can be used to separate data values, and can also be used in XGBoost and Random Forest systems. Because we can classify ideas into a wide range of destructive and non-toxic, we will use the concept of results in both directions. Because our data are independent of each other and have nothing in common with two different concepts, we use the Nave Bayes classifier to classify them. Because the data is marked, we can immediately use a controlled machine learning algorithm.

1. Logistic Regression

One of the few approaches to the classification of utility data is logistic regression. Suppose you have a medical history of a patient with a tumor. It is necessary to determine if the tumor is malignant (dangerous). The concept also suggests whether it is toxic or not. 1 means positive class and 0 means negative class. Positive(1) mean that it is toxic and negative(0) means that it is non-toxic. This is calculated by the sigmoid function. Suppose we have only one piece of information. It has different features. Suppose we have the right weight matrix. Depending on the data point, we now need to provide a class tag (classify it as 1 or 0). The weight vector is applied, and the input vector is multiplied by it, which gives a scalable output. To get a value between 0 and 1, this output is placed into a Sigmoid function. We'll name this likelihood of the projected class 1 for now. If the probability is greater than 0.5, the predicted class is 1. If the probability is less than 0.5, the predicted class is 0.

$$[x_1 \quad x_2 \quad x_3 \quad \dots \quad x_n]$$

$$sigmoid(x) = \frac{1}{1 + e^{-x}}$$

Fig : Sigmoid Function

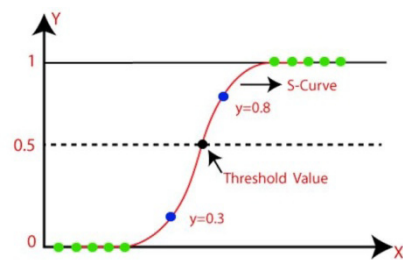


Fig : Logistic Regression curve

2. Random Forest

Random Forest is a controlled machine learning approach to classification and regression problems. Uses an explicit majority for classification and a moderate delay to create a decision tree of different data. One of the main features of the random forest algorithm is the ability to manage a set of data with classified and continuous variables, such as regression and classification. As for the rating, it performs better than its competitors.

STEPS

- Step 1: Random Forest takes random entries from the Q-record database..
- Step 2: A separate end tree is created for each model.
- Step 3: Each end tree will have an output.
- Step 4: The final decision on classification and regression is based on an appropriate or intermediate majority.

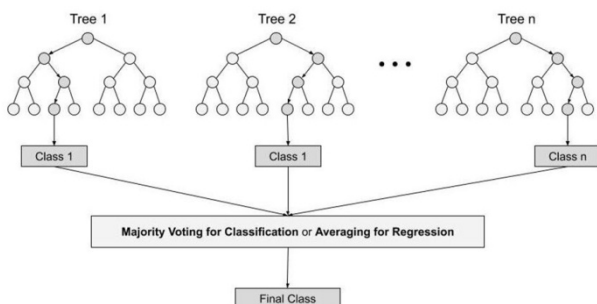


Fig : Random Forest

3. SVM Classifier

The reference vector machine (SVM) is a machine learning technique with classification and regression monitoring. SVM defines a cloud page that classifies the boundaries between two data sets. SVM is often used to classify data, although it can also be used for regression. This is a fast and reliable method that works well with low data. Another advantage of SVM is that it can explore various input functions without increasing system problems, using different types of core functions.

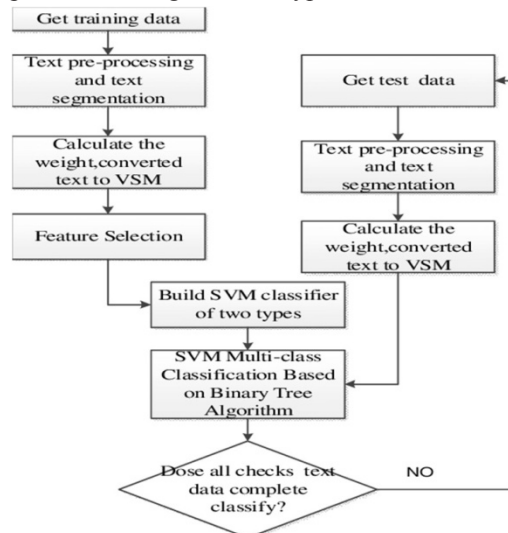


Fig : Flow chart of SVM

4. Naive Bayes

The Bayes theorem is a mathematical procedure for estimating conditional probabilities. A probability, as you may know, is the possibility of an event occurring. We call an event's probability if it has a chance of occurring.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Probability of B occurring given evidence A has already occurred Probability of A occurring
Probability of A occurring given evidence B has already occurred Probability of B occurring

Fig : Formula for Naïve Bayes

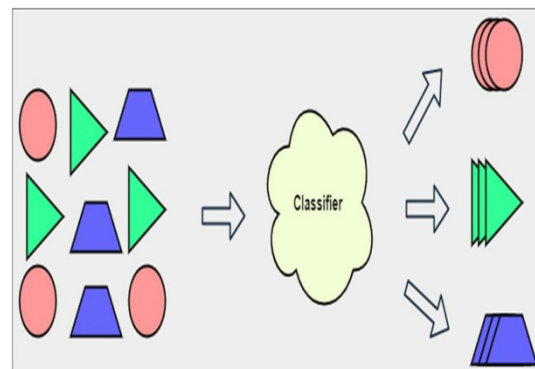


Fig : Naïve Bayes Classification

The Bayes theorem is the foundation of the Naive Bayes algorithm. It's primarily utilized for categorization tasks. In classification, we teach the model what each class belongs to using a labeled dataset, and then the model learns and classifies or labels a new dataset that has never been seen before. All of the features or

variables in the Naive Bayes model are treated as independent of one another.

5. XGBoost Classifier

Gradual gain, such as Random Forest (another end-of-tree algorithm), is a controlled machine learning approach for topics such as classification (men, women) and regression (expected value). The most common names to implement this method are Gradient Boosting Machines (GBM) and XGBoost. Gradient Boosting is a training group similar to Random Forest. This means that it links the final model to a set of individual models. These individual models have a poor prognosis and are very consistent, but combining several of these weak models into one group will give the best results. Random forest-related end trees are the most popular weak sample used in gradient reinforcement machines.

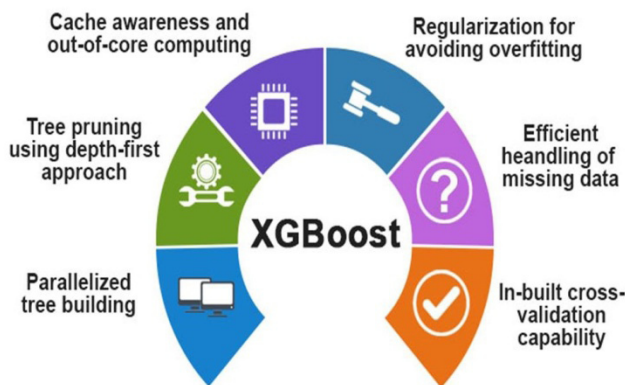


Figure :XG boost

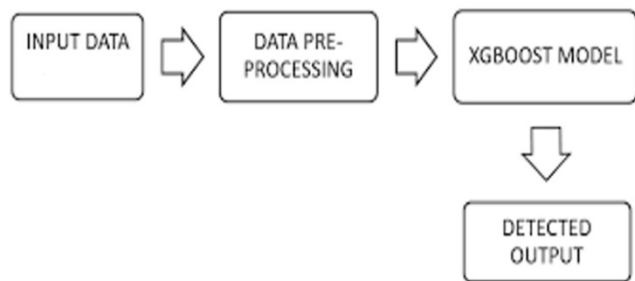


Fig :XGBoost process

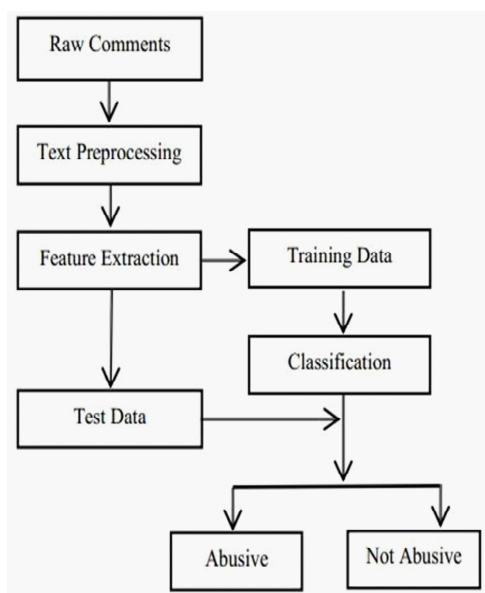


Fig : Flow chart for detecting toxic comments

III. Project Requirements

A. Python

Python is an advanced interactive and semantic programming language. It is written dynamically and has a ridiculously simple syntax. It is written in plain English and has no syntax like Java or C. Seriously, it follows a stroke. Using language keywords, such as when, how, in, etc., makes the code much clearer and simpler. It is a more semantic language because it is a structured language with all the complete structures necessary for a specialized synthetic language. Python is an interactive programming language. Python does not have a compiler, so the compiler executes the code directly at runtime. This is similar to how JS and PHP work.

1. Python is a universal programming language that can be used for a variety of applications. It was established in 1991 by Guido van Rossum as a computer language with a design philosophy that emphasises code readability and the use of large spaces.
2. Dynamic kind systems and memory management square measure 2 of Python’s options. Object-oriented and imperative programming, still as useful and procedural approaches, square measure all supported by the big and in-depthlibrary.
3. Python is a language for writing program descriptions. Use Python Prompt or Idle to build programs on Python Translator and test variables and environments. Python is an object-oriented programming language. Python supports object-oriented programming with full use of technology and its programming capabilities.

B. Numpy

NumPy is a Python module that uses arrays to perform a wide range of functions. To ensure faster execution, it can be implemented using the C and C++ libraries. It works faster with the NumPy module than similar performance without it. This ability to fill signs with pure math is still a correction in February. It will not be added automatically when Python is inserted. Must be installed at the same time as the NumPy pip installation command. The Python number is indicated by NumPy.

C. Pandas

Pandas are a free, open source architecture that facilitates and directly manages linked and encrypted data. It includes several data structures and methods for working with numerical data and time series data. This library is based on the NumPy Python library. Pandas are a fast-paced program, and its users benefit from its speed and performance.

D. Matplotlib

For 2D series diagrams, Matplotlib is a great Python imaging library. Matplotlib is a cross-platform NumPy-based data imaging software compatible with the entire SciPy stack. First launched in 2002 by John Hunter. One of the most important advantages of imaging is that it allows you to see large amounts of data in easy-to-understand images. Line, strip, skater, histogram and many other maps are available on Matplotlib.

E. Seaborn

Seaborn is a great Python imager for programming statistical graphics. It has attractive default styles and a color palette that make statistical maps very attractive. It uses matplotlib software and is closely related to the panda data structure. Seaborn’s mission

is to develop data visualization as an important part of data recognition and understanding. It provides a database-based API that allows you to switch between multiple visual representations of a variable for a better understanding of the data.

F. Scikit-Learn

Scikit-Learn is a free Python machine learning library. Supports both controlled and uncontrolled machine learning and includes a number of algorithms for classification, regression, grouping, and size reduction. This library was created using many libraries you know, such as NumPy and SciPy. It can also be used with other libraries, such as Panda and Seaborn.

IV. System Design

A. Architecture

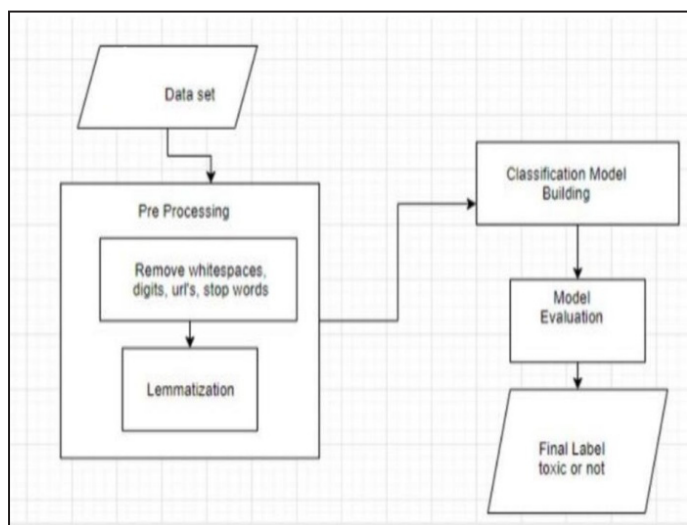


Fig : Architecture

We upload our dataset first, then perform preprocessing such as stemming and removing stop words, before applying our machine learning models and calculating accuracy for each model. Finally, the most accurate model is chosen. The best model is then used to estimate toxicity based on unknown data.

V. Outputs

Distribution of Iweets in the Dataset

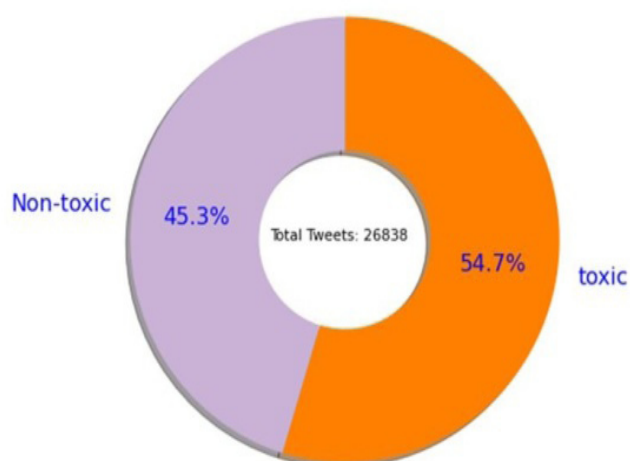


Fig : Distribution of Tweets in Dataset

	Algorithms	Accuracy
0	Naive-Bayes	0.791330
1	SVM	0.789219
2	Logistic Regression	0.800149
3	Random forest	0.583778
4	Xgboost	0.755186

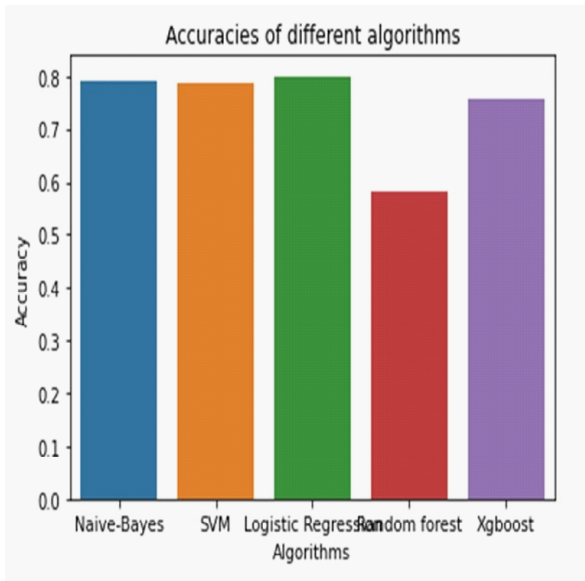


Fig : Accuracies of Algorithms

Now we predict on unseen data by taking our best algorithms Logistic Regression.

```

text='im mad at you bitch' # enter your text for testing
# get the prediction for the text
text=count_vect.transform([text])
pred=lr.predict(text_) # predicting
prob=np.amax(lr.predict_proba(text_))# getting probability
print(pred,prob)
  
```

['Offensive'] 0.9951340848415849

```

text='i am attracted to ur beauty' # enter your text for testing
# get the prediction for the text
text=count_vect.transform([text])
pred=lr.predict(text_) # predicting
prob=np.amax(lr.predict_proba(text_))# getting probability
print(pred,prob)
  
```

['Non-offensive'] 0.6150410062957921

```

text='ugly face idiot' # enter your text for testing
# get the prediction for the text
text=count_vect.transform([text])
pred=lr.predict(text_) # predicting
prob=np.amax(lr.predict_proba(text_))# getting probability
print(pred,prob)

```

```
['toxic'] 0.8856673402446275
```

```

text='you are pretty' # enter your text for testing
# get the prediction for the text
text=count_vect.transform([text])
pred=lr.predict(text_) # predicting
prob=np.amax(lr.predict_proba(text_))# getting probability
print(pred,prob)

```

```
['Non-toxic'] 0.6230064987506286
```

```

text='i will kill you' # enter your text for testing
# get the prediction for the text
text=count_vect.transform([text])
pred=lr.predict(text_) # predicting
prob=np.amax(lr.predict_proba(text_))# getting probability
print(pred,prob)

```

```
['toxic'] 0.556559847255282
```

Fig : Unseen data prediction

VI. Conclusion

We tested the accuracy of five machine learning methods: logistic regression, support vector machine, random forest, SVM classifier, and XGBoost classifier. After careful consideration, we can conclude that the logistic regression model has the best performance in terms of accuracy, as it is the most accurate model among other models. As a result, we choose our final model based on its precise nature. If there is a logistic regression model, we get the maximum. We will use the logistic regression model as our newest machine learning strategy, as it works best with our data.

VII. Future Scope

Other machine learning models can be used to calculate accuracy, loss, and loss recording to improve performance in future research. We consider RNN, BERT, multilayer perceptron and GRU, all of which are deep learning methods. As a result, we can explore a variety of alternative strategies to help improve the end result.

References

- [1] Rahul, H. Kajla, J. Hooda and G. Saini, "Classification of Online Toxic Comments Using Machine Learning Algorithms," 2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS), 2020, pp. 1119-1123, doi: 10.1109/ICICCS48265.2020.9120939 [2] M. Duggan, "Online harassment 2017," Pew Res., pp. 1–85, 2017, doi: 10.2419.4372.
- [3] M. A. Walker, P. Anand, J. E. F. Tree, R. Abbott, and J. King, "A corpus for research on deliberation and debate," Proc. 8th Int. Conf. Lang. Resour. Eval. Lr. 2012, pp. 812–817, 2012.
- [4] J. Cheng, C. Danescu-Niculescu-Mizil, and J. Leskovec, "Antisocial behavior in online discussion communities," Proc. 9th Int. Conf. Web Soc. Media, ICWSM 2015, pp. 61–70, 2015.
- [5] B. Mathew et al., "Thou shalt not hate: Countering online hate speech," Proc. 13th Int. Conf. Web Soc. Media, ICWSM

2019, no. August, pp. 369–380, 2019.

- [6] C. Nobata, J. Tetreault, A. Thomas, Y. Mehdad, and Y. Chang, "Abusive language detection in online user content," 25th Int. World Wide Web Conf. WWW 2016, pp. 145–153, 2016, doi: 10.1145/2872427.2883062.
- [7] E. K. Ikonomakis, S. Kotsiantis, and V. Tampakas, "Text Classification Using Machine Learning Techniques," no. August, 2005.
- [8] M. R. Murty, J. V. Murthy, and P. Reddy P.V.G.D, "Text Document Classification based on Least Square Support Vector Machines with Singular Value Decomposition," Int. J. Comput. Appl., vol. 27, no. 7, pp. 21–26, 2011, doi: 10.5120/3312-4540.
- [9] E. Wulczyn, N. Thain, and L. Dixon, "Ex machina: Personal attacks seen at scale," 26th Int. World Wide Web Conf. WWW 2017, pp. 1391–1399, 2017, doi: 10.1145/3038912.3052591.
- [10] H. Hosseini, S. Kannan, B. Zhang, and R. Poovendran, "Deceiving Google's Perspective API Built for Detecting Toxic Comments," 2017, [Online]. Available: <http://arxiv.org/abs/1702.08138>.