

**VISVESVARAYA TECHNOLOGICAL UNIVERSITY, BELGAUM**  
“Jnana Sangama”, Belgaum-590018



**Minor Project (IS6C06)**  
— Schmaltz Surveyor —

*Submitted in partial fulfillment of the requirement for the completion of 6th  
semester of*

**BACHELOR OF ENGINEERING  
IN  
INFORMATION SCIENCE & ENGINEERING**

by

**Nithyashree Arunachalam                  4NI19IS058**

**Pradyoth P                  4NI19IS062**

**Shashank BU                  4NI19IS086**

**Tejasvini SJ                  4NI19IS106**

Under the Guidance of

**Mrs. Nandini BM**

**Assistant Professor**

**Department of ISE, NIE Mysuru**



**The National Institute of Engineering  
Mysuru-570008**

**Department of Information Science and Engineering  
NIE, Mysuru - 570008**

**2021-2022**

**THE NATIONAL INSTITUTE OF ENGINEERING  
MYSURU - 570008**

**Department of Information Science and Engineering**



**CERTIFICATE**

Certifies that the project work titled “Schmaltz Surveyor” is a bonafide work carried out by

**Nithyashree Arunachalam (4NI19IS058)**

**Pradyoth P (4NI19IS062)**

**Shashank BU (4NI19IS086)**

**Tejasvini SJ (4NI19IS106)**

in partial fulfillment for the requirements of the sixth semester BE in Information Science & Engineering prescribed by The National Institute of Engineering, Autonomous Institution under Visvesvaraya Technological University, Belagavi. It is certified that all corrections/suggestions indicated for Internal Assessment have been incorporated. The Project report has been approved as it satisfies the academic requirements concerning the project work prescribed for the sixth semester in Minor Project (IS6C06).

Signature of Guide

Signature of HOD

Signature of Principal

(Mrs. Nandini BM)

(Dr. S Kuzhalvaimozhi)

(Dr. Rohini Nagapadma)

## **ABSTRACT OF THE PROJECT**

Social media has gained immense popularity and has become a major global platform to stay connected as well as express opinions. A huge amount of content is created on various topics and comments are posted on these platforms daily. The feedback received on a certain piece of content can be either negative or positive. Other than this, receiving negative feedback, on various occasions might affect the mental health of the content creators and, in some cases, might also lead to cyberbullying. Social media is a necessity in today's time to stay connected, informed, and relevant and therefore such issues must be tackled.

Sentimental Analysis (also known as Opinion Mining in this case) reads people's sentiments or emotions towards particular things or topics. Sentiment analysis is a machine learning tool that will help analyze and categorize the texts (written content and comments) as positive or negative. Sentiments cannot be directly predicted by analyzing the text and therefore the texts are vectorized or converted into numeric values. The text is converted using vectorization methods in NLP (Natural language Processing) such as Bag of Words or TF-IDF(Term Frequency–Inverse Document Frequency). This data is then put through machine learning algorithms such as Support Vector Machine, Logistic Regression, Random Forest and K Nearest Neighbours. Based on the result after using these methods, the most accurate machine learning model is selected. Python would be the programming language across the whole project.

The number of positive or negative texts and comments can be used to analyze people's reception or opinion on what particular piece of content is getting. Further, it can also be used by social media platforms to detect negative comments and delete them or classify content based on the number of positive or negative comments.

**Keywords:** Social Media, Sentiments, Analysis, Twitter.

## **ACKNOWLEDGEMENT**

The success and the outcome of this project required a lot of guidance and assistance from many people and we are incredibly fortunate to have got this all along with the completion of project work.

We express our profound thanks to **Dr. Rohini Nagapadma**, Principal, NIE, Mysuru for his much-needed moral support and encouragement.

We are grateful to **Dr. S Kuzhalvaimozhi**, Professor & HOD, Information Science and Engineering, NIE for her support and encouragement in facilitating the progress of this work.

We sincerely extend our thanks to our Project Coordinator and Project Guide **Mrs. Nandini BM**, Assistant Professor in the Department of ISE, for her guidance, technical expertise, encouragement, and timely help in making this project a reality.

We would like to extend our sincere regards to all the teaching and non-teaching staff of the Department of ISE, NIE Mysuru for their assistance and timely support.

We would also like to give credit to the authors of the various resources which were made available on the Internet for our reference.

**Nithyashree Arunachalam (4NI19IS058)**

**Pradyoth P (4NI19IS062)**

**Shashank BU (4NI19IS086)**

**Tejasvini SJ (4NI19IS106)**

## **TABLE OF CONTENTS**

<b>Sl. No.</b>	<b>Chapter</b>	<b>Page No.</b>
1.	<b>Introduction</b> 1.1 Objective 1.2 Purpose 1.3 Existing System 1.4 Proposed System	5
2.	<b>Literature Survey</b>	8
3.	<b>System Requirements</b>	9
4.	<b>System Design</b>	10
5.	<b>System Implementation</b> 5.1 Exploration of Dataset 5.2 Data Preprocessing 5.3 Applying Classifiers 5.4 Evaluation of the Models used	12
6.	<b>Testing</b>	29
7.	<b>Conclusion</b>	30
8.	<b>Future Enhancements</b>	31
9.	<b>References</b>	32

---

## **Chapter 1:**

### **Introduction**

Social Media has become a necessity to remain relevant. It gives everyone a platform to express their opinions and views on a larger scale. Since opinions are always going to be conflicting, much content on the internet can be segregated into positive and negative. Thousands and millions of data are produced regularly by these social media sites which can help us to look at the overall effect i.e we can calculate the amount of negative and positive content present, be it the opinions of a person towards a particular topic or towards a video.

The impact of social media is very intense. Severely negative comments have the potential to hurt and ruin the mental health of a person. As much as it is essential to express our individual opinions to keep a conversation going, there is no need to express and give out negative comments to random content creators and other people who are doing nothing but entertaining billions of people around the globe anonymously.

Machine learning is an emerging scientific field in data science dealing with the ways in which machines learn from experience. Machine learning (ML) is a type of artificial intelligence that allows software applications to become more accurate at predicting outcomes without being explicitly programmed to do so. Machine learning algorithms use historical data as input to predict new output values.

Sentimental analysis, also known as opinion mining, is a method of analyzing and classifying data on the basis of the emotion they express with the help of machine learning models and algorithms present.

The main aim of this project is to use machine learning and be able to predict a positive and negative comment or post in the future and create a software of sorts that can be employed which blocks a user from typing negative comments. An already classified dataset is taken which is used to train the model and then the model will be further will able to predict the sentiment.

Schmaltz means excessive sentimentality. Therefore, this project is called ‘Schmaltz Surveyor’ which on its own suggests that the sentiments are going to be surveyed here.

## **1.1 Our Objective**

The main objective of this project is to be able to predict comments as positive or negative. This is done using an ML model which has been trained on a dataset that already has comments labeled based on positive or negative. This model will be applied at the backend while the frontend will have a functionality wherein the user will not be able to post a negative tweet/content. The user's text will go through the model at the backend that will predict its sentiment and based on it the user will be able to post it.

In the first half of the project, a set of 4 classifiers namely:

1. Logistic Regression
2. Support Vector Machine
3. K Nearest Neighbor
4. Random Forest

which are predominantly used to categorize will be used to train the dataset. The models then predict a set of tweets which have already been labeled and then compare its prediction against the actual results to calculate the accuracy and create a confusion matrix. The model with the highest accuracy score is taken as the main model to predict the sentiment.

## **1.2 Purpose**

The purpose of this project is to analyze the amount of positive and negative influence present on the social media platform(Twitter in this case) and further be able to predict the sentiment of tweets/posts made by the users in future. This sentiment analysis will help us in developing an easy to use software which will prevent the user from putting out negative tweets out there. By crunching large volumes of data, machine learning technology can help to reduce the amount of negativity present on such social media platforms and make it a healthy and more positive environment for the users.

### **1.3 Existing System**

Currently the only functionality available on platforms such as twitter is that of reporting offensive tweets/comments or completely blocking the users creating chaos. After the tweet/comment is reported, it is in most cases manually reviewed and then the decision taken if it should be removed or not.

### **1.4 Proposed System**

The proposed system is to develop an easy to use software in which when a user is about to post something , that piece of text goes through the trained model in the backend and its sentiment is predicted. If the sentiment is negative the user is barred from posting that piece of content.

The main page will contain a dialogue box in which the user will type his/her text. On clicking on the post button, the text will go through the chosen model and its sentiment is predicted. If the sentiment is positive the tweet is posted whereas if it is a negative sentiment the statement does not get posted and the user gets a warning.

---

## **Chapter 2:**

### **Literature Survey**

1. Xing Fang, and Justin Zhan[1], worked on Sentiment Analysis on Amazon Online Products Review by collecting data from amazon.com. They figured out the issues of categorical sentiment polarity using Machine Learning Algorithms. They used many libraries that hold Naive Bayesian, SVM, and Random Forest.
  2. Faizan[2] built a model for the analysis of feeling using the KNN algorithm with unigram, bigram, and n-gram features. They then performed training and testing of their model on the US Airlines data set, for which they attained an accuracy of 65.33%.
  3. Chirag Kariya and Priti Khodke's[3] paper explains various steps involved in the analysis of Twitter sentiments along with the various tools that are used to perform Twitter sentiment analysis. Amongst the various algorithms available, the KNN algorithm is used to increase the efficiency of sentiment analysis whereas Naive Bayes for simple and efficient sentiment analysis by classifying the tweets as either positive, negative, or zero.
  4. Akshay Amolik et al.[4] proposed sentiment analysis, and they accurately classified tweets by using Feature Vector and classifiers like Naïve-Bayesian and SVM. Exception of lower recall and accuracy, Naïve Bayesian has better precision as a comparison to SVM. SVM gives better results when it comes to accuracy. With the increase of training data, the accuracy of classification will also increase.
-

## **Chapter 3:**

### **System Requirements**

#### **Hardware Requirements:**

- Processor- 5th gen Intel Core i3 equivalent or greater
- RAM- 8GB or above
- Keyboard and a Mouse
- Monitor

#### **Software Requirements:**

- Windows 11, macOS or GNU/Linux.
  - Visual Studio Code
  - HTML, CSS and Javascript
  - Python 3.10
  - Flask
  - Chrome
  - Jupyter Notebook
  - Scikit-learn
  - Natural Language Toolkit (NLTK)
  - Modules used were Pickle, numPy, pandas, re, Tweepy, matplotlib, WordCloud.
-

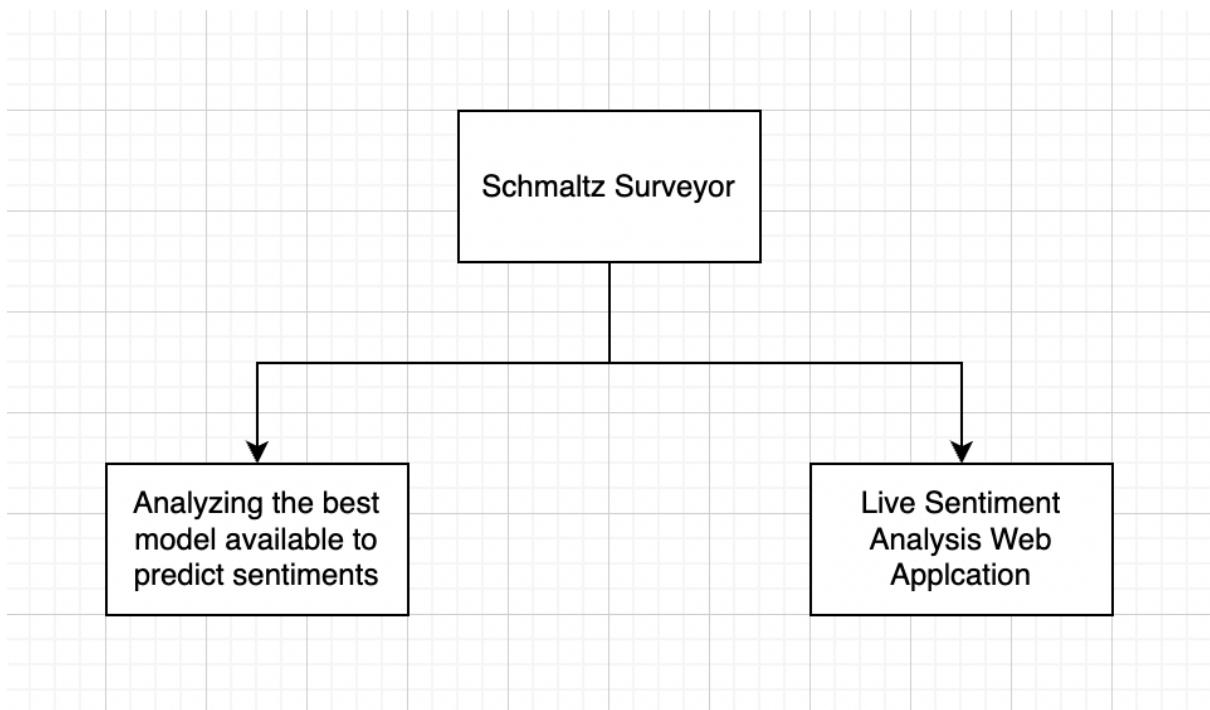
## Chapter 4:

### System Design

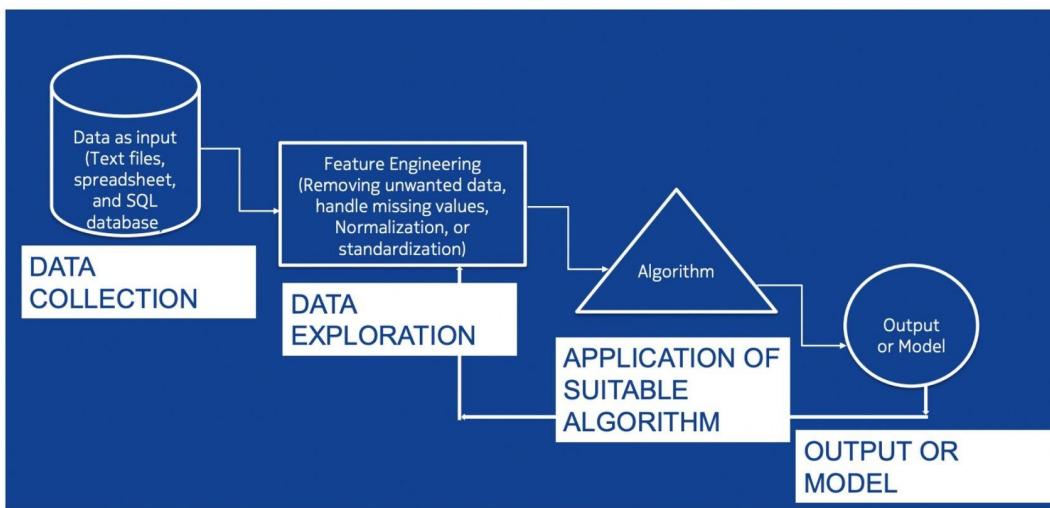
The project has been mainly divided into two parts.

The first part mainly deals with analyzing the efficiency of the methods through which Sentiment Analysis has been performed like we've seen in the Literature Survey and combining the results to find the best model for our use case.

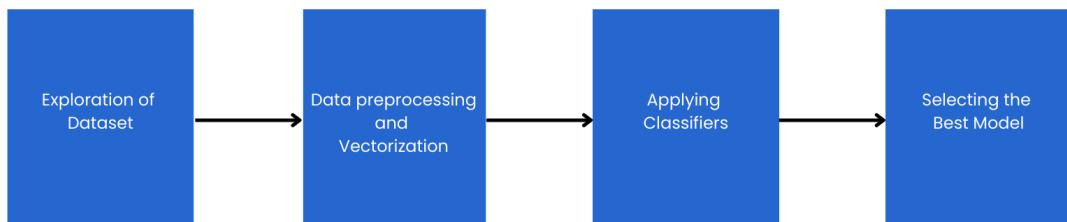
The second part mainly deals with utilizing the information we got in the first part, deploying it at the backend, and making the model we chose to predict sentiments as and when we ask for from the User Interface provided through the Frontend, basically a Live Sentiment Analysis Web Application. The application provides an output that shows the sentiment user is trying to display through that particular message.



## System Design of any Machine Learning Application:



The above figure shows the typical design of any Machine Learning Application in a generalized manner. Based on the same, we've created our Project Flow as shown.



1. Exploration of Dataset
2. Data preprocessing and Vectorization
3. Applying Classifiers
4. Selecting the best Model.

# **Chapter 5:**

## **System Implementation**

The data set was initially taken from Kaggle and it contains three columns:

1. ID: Id number of the tweet
  2. Tweet
  3. Label: Positive sentiment is represented by 1 while negative sentiment is represented by 0

Let us focus on the System Implementation of the first part as said in the previous chapter, System Design.

## Exploration of Dataset

This step involved finding a dataset suitable for all our needs and requirements while keeping a large amount of data, so we could train the model better.

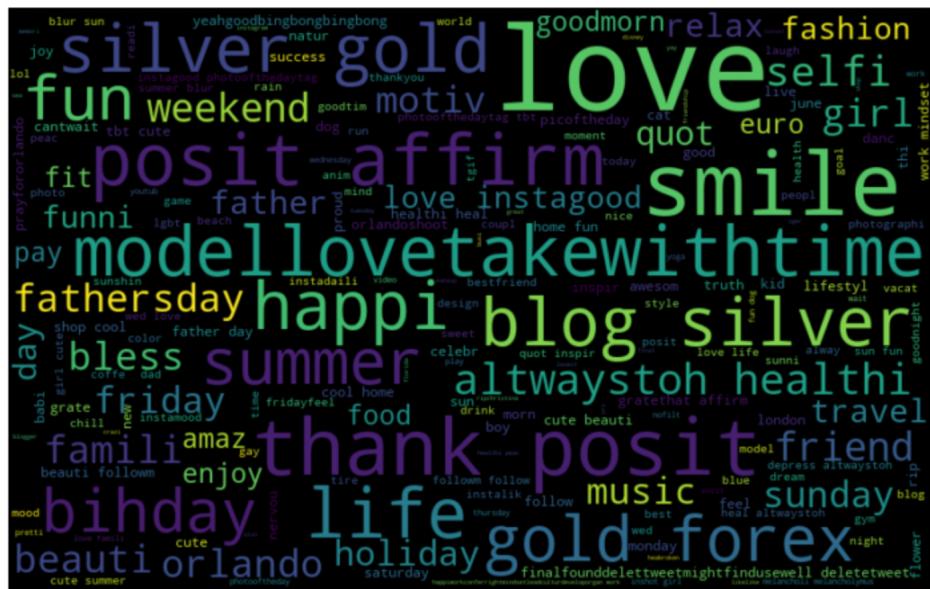
The Dataset referred to for training the models is a labeled dataset of 31,962 tweets. The dataset is provided in the form of a CSV file with each line storing a tweet id, its label, and the tweet.

Exploratory Data Analysis was also performed, where we got the following outputs.

## 1. Common Frequent Word Count Visualization



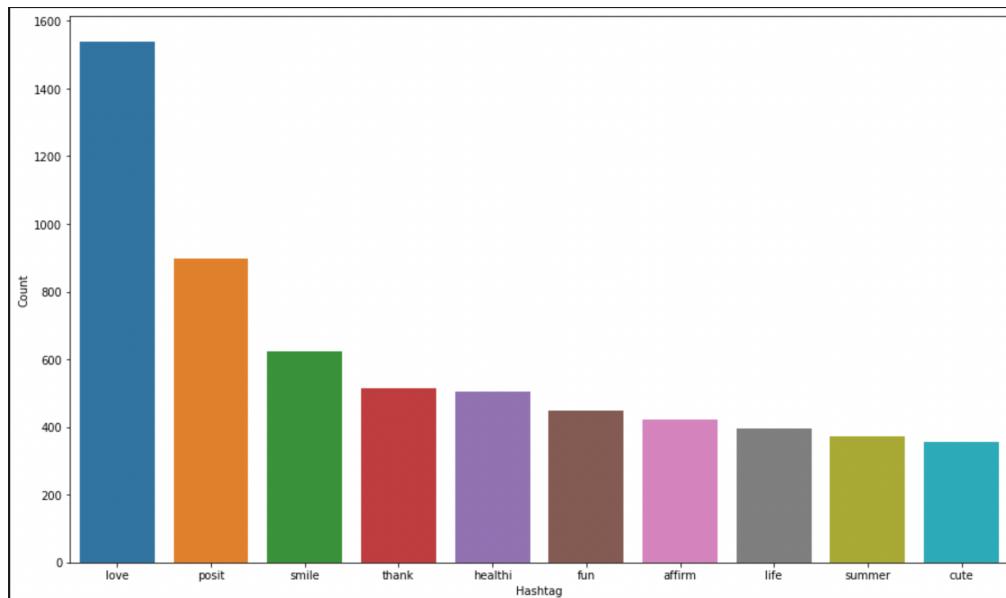
## 2. Frequent Word Visualization for positive Labeled tweets



### 3. Frequent Word Visualization for negative Labeled tweets



Now, we could also get the information about the top 10 used hashtags as shown in the graph below.



We've also split our Input dataset into Training Data and Test Data as we need the same for processing and applying classifiers later.

## Input Split

```
# feature extraction
from sklearn.feature_extraction.text import CountVectorizer
bow_vectorizer = CountVectorizer(max_df=0.90, min_df=2, max_features=1000, stop_words='english')
bow = bow_vectorizer.fit_transform(df['clean_tweet'])
```

Python

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(bow, df['label'], random_state=42, test_size=0.25)
```

Python

The test set already contains the actual values. Therefore after prediction, the predicted values will be measured against the actual values and the accuracy will be calculated.

---

## Data Preprocessing and Vectorization:

Step 1: Checking for missing values.

Step 2: Text Normalization.

Normalize the text data as texts from such online platforms usually contain inconsistent language and the use of special characters in place of letters. To tackle such inconsistencies in data, Regex.

Step 3: Lemmatization

Lemmatization is the process of grouping together the different forms of a word so they can be analyzed as a single item. For example, we do not want the Machine Learning algorithm to treat eating, eats, and eat as three separate words because they convey the same message. Lemmatization helps reduce the words to their root form.

Step 4: Removal of stop words

Stop words are a set of commonly used words in the language. Examples of stop words in English are “a”, “the”, “is”, “are” and etc. Stop words are commonly used in Text Mining and Natural Language Processing (NLP) to eliminate words that are so commonly used that they carry very little useful information.

	<b>id</b>	<b>label</b>	<b>tweet</b>	<b>clean_tweet</b>
0	1	0	@user when a father is dysfunctional and is s...	whenfatherdysfunctselfishdragkidintodysfunct#run
1	2	0	@user @user thanks for #lyft credit i can't us...	thank#lyftcreditcausheyofferwheelchairvan#dis...
2	3	0		bihday your majesty
3	4	0	#model i love u take with u all the time in ...	#modellovetakewithtime
4	5	0	factsguide: society now #motivation	factsguidsociety#motiv

[+ Code](#) [+ Markdown](#)

This is the data frame head, after Steps 1-4. This stands common for all classifiers and hence, we do it beforehand. As can be seen, the clean\_tweet column contains tweets without stop words, hashtags, and they have also been lemmatized.

## Step 5: Converting to Numerical Form

The dataset needs to be converted into the numeric form so that it can be put through classifiers.

TF IDF: TF-IDF means Term Frequency - Inverse Document Frequency. This is a statistic that is based on the frequency of a word but it also provides a numerical representation of how important a word is for statistical analysis.

Term frequency works by looking at the frequency of a particular term you are concerned with relative to the document. This is preferred over CountVectorizer method for vectorisation as TF-IDF has the ability to measure the weight or in layman terms the importance of the word in the document.

The document frequency of word i represents the number of documents in the corpus with word i in them. Let us represent document frequency for word i by  $df_i$ . With N as the number of documents in the corpus, the tf-idf weight for word i in document j is computed by the following formula:

$$w_{ij} = tf_{ij} \times \left(1 + \log \frac{1+N}{1+df_i}\right)$$

The following papers can be referred to see how TF-IDF is used:

1. Jalilifard, A., Caridá, V.F., Mansano, A.F., Cristo, R.S., da Fonseca, F.P.C. (2021). Semantic Sensitive TF-IDF to Determine Word Relevance in Documents. In: Thampi, S.M., Gelenbe, E., Atiquzzaman, M., Chaudhary, V., Li, KC. (eds) Advances in Computing and Network Communications. Lecture Notes in Electrical Engineering, vol 736. Springer, Singapore.  
[https://doi.org/10.1007/978-981-33-6987-0\\_27](https://doi.org/10.1007/978-981-33-6987-0_27)
2. Kim, SW., Gil, JM. Research paper classification systems based on TF-IDF and LDA schemes. *Hum. Cent. Comput. Inf. Sci.* 9, 30 (2019).  
<https://doi.org/10.1186/s13673-019-0192-7>

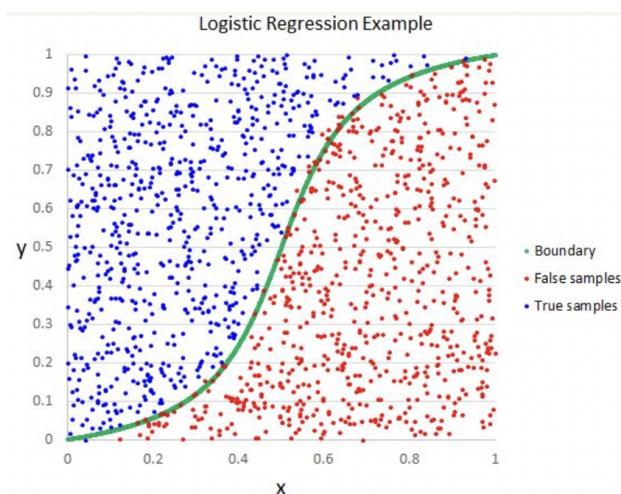
## Applying Classifiers:

We've used our dataset on four classifiers, namely

- **Logistic Regression**

Logistic regression is a statistical analysis method to predict a binary outcome, such as 0 or 1, based on prior observations of a data set.

A logistic regression model predicts a dependent data variable by analyzing the relationship between one or more existing independent variables.



The above example of Logistic Regression shows the Classification of False and True samples based on some parameter. The regression line divides them into two classes.

Now, we apply Logistic Regression to our dataset:

```
Logistic Regression

[25]
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size = 0.5, random_state=45)
from sklearn.feature_extraction.text import TfidfVectorizer

[26]
vectorizer = TfidfVectorizer()
X_train_vec = vectorizer.fit_transform(X_train)
X_test_vec = vectorizer.transform(X_test)
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

[27]
lr.fit(X_train_vec, y_train)
lr_score = lr.score(X_test_vec, y_test)
print("Results for Logistic Regression with tfidf")
print(lr_score)
y_pred_lr = lr.predict(X_test_vec)

... Results for Logistic Regression with tfidf
0.9449971841561855

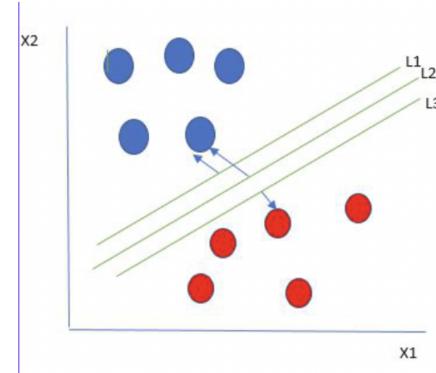
[28]
#Confusion matrix
from sklearn.metrics import confusion_matrix
cm_lr = confusion_matrix(y_test, y_pred_lr)
tn, fp, fn, tp = confusion_matrix(y_test, y_pred_lr).ravel()
print(tn, fp, fn, tp)
tpr_lr = round(tp/(tp + fn), 4)
tnr_lr = round(tn/(tn+fp), 4)
print(tpr_lr, tnr_lr)

... 267 856 23 14835
0.9985 0.2378
```

As can be seen, we have a result of accuracy percentage of ~94.5% with Logistic Regression.

- **Linear Support Vector Classifier or Support Vector Machine**

SVM performs classification by finding the hyper-plane that differentiates the classes we plotted in n-dimensional space.



It classifies positive and negative examples, here blue and red data points.

Largest margin is found in order to avoid overfitting. The optimal hyperplane is at the maximum distance from the positive and negative examples

Now, we apply SVM to our dataset.

## Using SVM

```
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer,CountVectorizer
from sklearn.metrics import confusion_matrix
]

x=df['clean_tweet']
y=df['label']
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x,y,test_size = 0.5, random_state=55)
vectorizer = TfidfVectorizer()
X_train_vec = vectorizer.fit_transform(X_train)
X_test_vec = vectorizer.transform(X_test)
]

from sklearn import svm
#params = {'kernel':('linear', 'rbf'), 'C':[1, 10, 100]}
svcl = svm.SVC(kernel = 'rbf')
#clf_sv = GridSearchCV(svcl, params)
svcl.fit(X_train_vec, y_train)
svcl_score = svcl.score(X_test_vec, y_test)
print("Results for Support Vector Machine with tfidf")
print(svcl_score)
]

Results for Support Vector Machine with tfidf
0.9537575871347225

y_pred_sv = svcl.predict(X_test_vec)
#Confusion matrix
from sklearn.metrics import confusion_matrix
cm_sv = confusion_matrix(y_test, y_pred_sv)
tn, fp, fn, tp = confusion_matrix(y_test, y_pred_sv).ravel()
print(tn, fp, fn, tp)
tpr_sv = round(tp/(tp + fn), 4)
tnr_sv = round(tn/(tn+fp), 4)
print(tpr_sv, tnr_sv)
]

413 721 18 14829
0.9988 0.3642
```

As can be seen, we get an accuracy of ~95.38% with SVM.

- **K Nearest Neighbors**

K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.

K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for Classification problems.

K-NN is a non-parametric algorithm, which means it does not make any assumptions on underlying data.

It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.



Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.

Now, we apply kNN to our dataset.

## k Nearest Neighbours

```
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer,CountVectorizer
from sklearn.metrics import confusion_matrix
X = df['clean_tweet']
y = df['label']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.5, random_state=65)
vectorizer = TfidfVectorizer()
X_train_vec = vectorizer.fit_transform(X_train)
X_test_vec = vectorizer.transform(X_test)

[29]
```

```
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train_vec, y_train)
knn_score = knn.score(X_test_vec, y_test)
print("Results for KNN Classifier with tfidf")
print(knn_score)
y_pred_knn = knn.predict(X_test_vec)

[30]
...
... Results for KNN Classifier with tfidf
0.9369876728615231
```

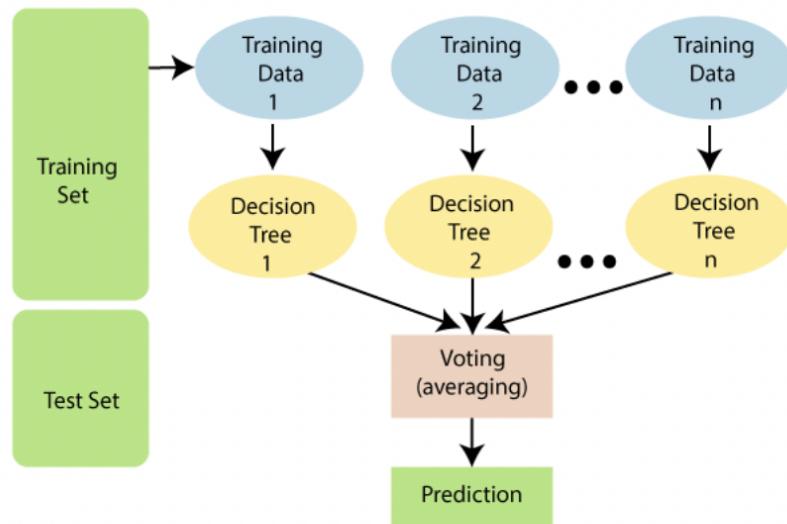
```
#Confusion matrix
cm_knn = confusion_matrix(y_test, y_pred_knn)
tn, fp, fn, tp = confusion_matrix(y_test, y_pred_knn).ravel()
print(tn, fp, fn, tp)
tpr_knn = round(tp/(tp + fn), 4)
tnr_knn = round(tn/(tn+fp), 4)
print(tpr_knn, tnr_knn)

[31]
...
... 152 1006 1 14822
0.9999 0.1313
```

As can be seen, kNN gives an accuracy of ~93.7%.

- **Random Forest**

Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset.



Instead of relying on one decision tree, the random forest takes the prediction from each tree, and based on the majority votes of predictions, it predicts the final output.

The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.

By the fact it prevents overfitting and relies on each tree, and majority of the trees votes, we can assume a high accuracy rate, however, we will check the same with our code on our dataset which can be seen in the figure in the next page.

Now, we apply Random Forest to our dataset.

### Random Forest Model

```
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn.metrics import confusion_matrix
X = df['clean_tweet']
y = df['label']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.5, random_state=65)
vectorizer = TfidfVectorizer()
X_train_vec = vectorizer.fit_transform(X_train)
X_test_vec = vectorizer.transform(X_test)

[32]

from sklearn.ensemble import RandomForestClassifier
test_classifier= RandomForestClassifier(n_estimators=200,random_state=0)
test_classifier.fit(X_train_vec,y_train)
predictions = test_classifier.predict(X_test_vec)

[33]

from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
print(confusion_matrix(y_test, predictions))
print(classification_report(y_test, predictions))
print("Results for Random Forest with tfidf")
rf_score = accuracy_score(y_test, predictions)
print(rf_score)
test_classifier

[35]
... [[ 512  646]
 [ 57 14766]]
      precision    recall   f1-score   support
          0       0.90      0.44      0.59      1158
          1       0.96      1.00      0.98     14823

      accuracy                           0.96      15981
     macro avg       0.93      0.72      0.78      15981
  weighted avg       0.95      0.96      0.95      15981

Results for Random Forest with tfidf
0.9560102621863463

RandomForestClassifier(n_estimators=200, random_state=0)

[36]

#Confusion matrix
text_classifier = confusion_matrix(y_test, predictions)
tn, fp, fn, tp = confusion_matrix(y_test, predictions).ravel()
print(tn, fp, fn, tp)
tpr_knn = round(tp/(tp+fn), 4)
tnr_knn = round(tn/(tn+fp), 4)
print(tpr_knn, tnr_knn)

[36]
... 512 646 57 14766
0.9962 0.4421
```

As can be seen, we have a result of accuracy percentage of ~95.6%

## Summary of the Outputs:

The Model used as Classifier	Accuracy Score (in %)
Support Vector Machines	95.38
Logistic Regression	94.50
k Nearest Neighbors	93.69
<b>Random Forest</b>	<b>95.60</b>

From the above data comparison, it is clear that Random Forest Model, is the way to go with the highest accuracy of approx. 95.60%

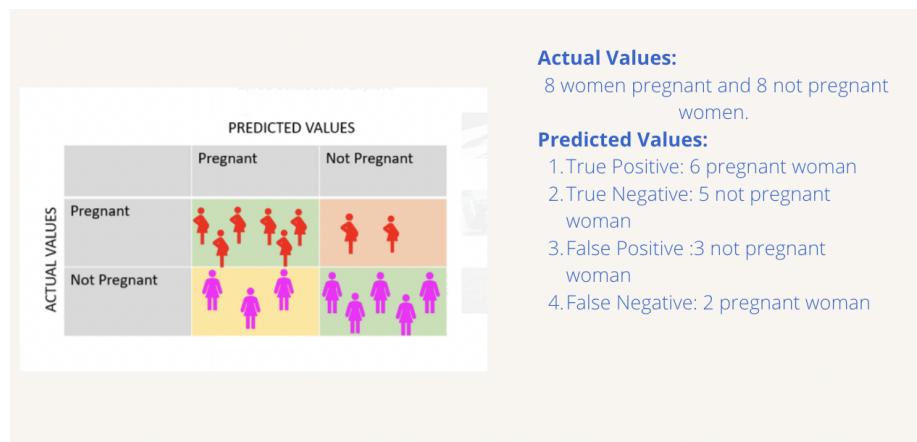
## Evaluating the Models Used:

Since this is classification, there is a lot of confusion, especially when there is a mismatch in prediction, giving rise to True and False predictions.

When there is sample data and we use modeling through which we can predict its class/labels. The “true” represents the records that the model was able to identify its class, whereas “false” represents the records that the model was not able to identify.

There are also “Positive” and “Negative” classes and depending on the objective of the study, we take one as a positive and the other as a negative.

To understand this better, let us look at an example, with the Actual Values and Predicted Values.



There is also the concept of Confusion Matrix, which helps us understand the problem better. A confusion matrix is a table that represents the summary of the prediction results on a classification problem.

		PREDICTED VALUES		ACTUAL VALUES	
		Positive	Negative	Positive	Negative
ACTUAL VALUES	Positive	True Positive (TP)	False Negative (FN)	Positive	True Positive (TP)
	Negative	False Positive (FP)	True Negative (TN)	Negative	False Negative (FN)

The position of the predicted values and actual values changes the position of False negative (FN) and False positive (FP) but True positive (TP) and True negative (TN) remains in the same place in the matrix placed diagonally to each other. But because of this, the situation becomes confusing.

True Positive(TP): Values that are actually positive and predicted positive.

False Positive(FP): Values that are actually negative but predicted to be positive.

False Negative(FN): Values that are actually positive but predicted to be negative.

True Negative (TN): Values that are actually negative and predicted to be negative.

Rate is a measuring factor in a confusion matrix.

True Positive Rate(TPR):  $\text{True Positive} / (\text{True Positive} + \text{False Negative})$

True Negative Rate(FPR):  $\text{True Negative} / (\text{True Negative} + \text{False Positive})$

Let us analyze the results for the same we got from our models.

The Model used as a Classifier	True Positive Rate	True Negative Rate
Logistic Regression	0.9985	0.2378
Support Vector Machines	0.9988	0.3642
k Nearest Neighbors	0.9999	0.1313
Random Forest	0.9962	0.4421

Finding a balanced output and predicting True Positives and Negatives is important, and hence, out of the models we've shown here, Random Forest is the one we choose as the best one.

The classification report for Random Forest is:

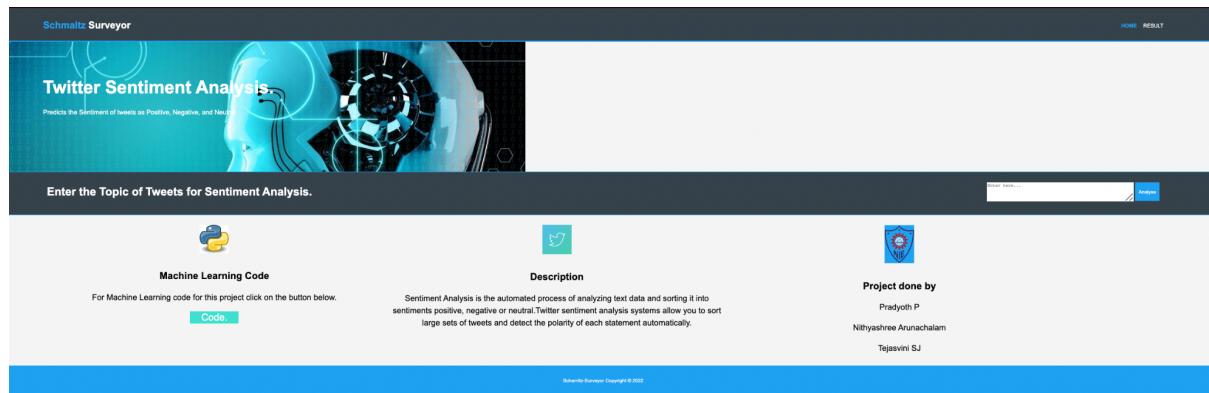
	precision	recall	f1-score	support
0	0.90	0.44	0.59	1158
1	0.96	1.00	0.98	14823
accuracy			0.96	15981
macro avg	0.93	0.72	0.78	15981
weighted avg	0.95	0.96	0.95	15981

**Results for Random Forest with tfidf**  
**0.9560102621863463**

The next step was to deploy a backend that does the same analysis and prediction LIVE when a user enters a keyword about tweets.

The dashboard can be seen below:

It was built using a Python framework, Flask, and modules such as Pickle, numPy, SkLearn, Tweepy, re, and pandas were used



The screenshot shows a web application interface for "Twitter Sentiment Analysis". At the top, there's a navigation bar with "HOME" and "RESULT" links. Below the header, a banner features a blue-toned image of a robotic arm and the text "Twitter Sentiment Analysis" along with a subtitle: "Predicts the Sentiment of tweets as Positive, Negative, and Neutral". A search bar with placeholder text "Enter the Topic of Tweets for Sentiment Analysis." is positioned above a main content area. The content area contains three sections: "Machine Learning Code" (with a Python icon and a "Code" button), "Description" (with a Twitter icon and a detailed text block about sentiment analysis), and "Project done by" (with icons for three team members: Pradyoth P, Nithyashree Arunachalam, and Tejasvini SJ). The bottom of the page has a footer with copyright information: "Schmalz Surveyor Copyright © 2020".

## **Chapter 6:**

### **Testing**

The main aspect of the testing involved for this project was to test it with different Vectorizers and Classifiers.

We had the option to choose two vectorizers, Count Vectorizer and TF-IDF Vectorizer, and based on the use case, both were tested and TF-IDF was chosen.

TF-IDF is better than Count Vectorizers because it not only focuses on the frequency of words present in the corpus but also provides the importance of the words. We can then remove the words that are less important for analysis, hence making the model building less complex by reducing the input dimensions.

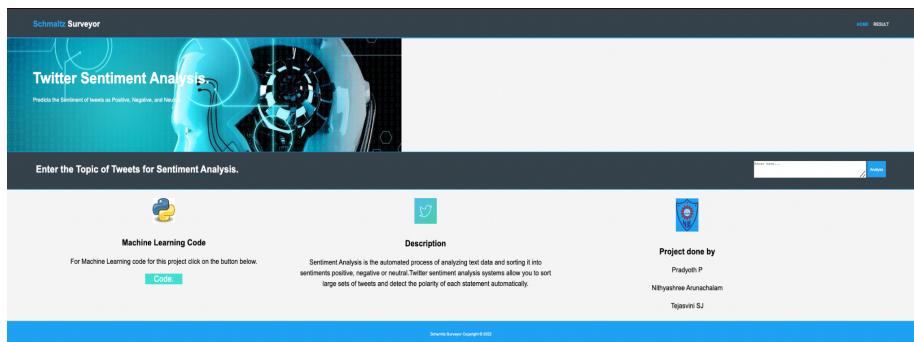
The next was to test the data for different classifiers and the same has been illustrated in-depth in the System Implementation part.

The data were tested with 4 classifiers, namely

1. Random Forest
2. Support Vector Machines
3. k Nearest Neighbours
4. Logistic Regression

Finally, based on the same, we chose the Random Forest which gave the output we needed with the highest accuracy.

The front end has been tested with various inputs of different keywords.



## **Chapter 7:**

### **Future Enhancements**

Firstly, English is definitely one of the most largely used languages for communication on social media platforms. But there is a lot of content that is produced in regional languages as well.

Detecting sentiment is definitely hard for regional languages as most of the models built are for English. But this is not an impossible task. This software can be used as a starting point for creating models to detect sentiment in regional languages.

Secondly, such software can be used in a widespread manner and this functionality can be added by users who want to avoid unnecessary negative comments on their page and want to filter out negative content.

---

## **Chapter 8:**

### **Conclusion**

In conclusion, as said in the introduction and the abstract, we've clearly demonstrated the Sentiment Analysis being done LIVE, along with the Comparison between various ML Models, and choosing the right one for our test case.

<b>The Model used as Classifier</b>	<b>Accuracy Score (in %)</b>
Support Vector Machines	95.38
Logistic Regression	94.50
k Nearest Neighbors	93.69
<b>Random Forest</b>	<b>95.60</b>

From the above data comparison of the classifiers used, it is clear that Random Forest Model is the way to go with the highest accuracy of approx. 95.60%

Classifying the tweet into positive and negative would help us understand and analyze the extent of positivity and negativity on Twitter. Sentiment Analysis is a great way to analyze the response to a particular tweet. The model at a backend to a web application helps us analyze the sentiments, as and when needed, and helps us with a dashboard of the same as well.

## **Chapter 9:**

### **References**

- [1]. Fang, Xing, and Justin Zhan. "Sentiment analysis using product review data." Journal of Big Data 2.1.
  - [2] Faizan. "Twitter Sentiment Analysis" International Journal of Innovative Science and Research Technology (2019)
  - [3] Chirag Kariya and Priti Khodke. "Twitter Sentiment Analysis" 2020 International Conference for Emerging Technology (INCET) Belgaum.
  - [4] Amolik, Akshay, et al. "Twitter sentiment analysis of movie reviews using machine learning techniques." International Journal of Engineering and Technology 7.6.
  - [5] Scikit Learn user guide [https://scikit-learn.org/stable/user\\_guide.html](https://scikit-learn.org/stable/user_guide.html)
  - [6] Numpy Documentation <https://numpy.org/doc/>
  - [7] Pandas Documentation <https://pandas.pydata.org/docs/>
  - [8] Dataset  
<https://www.kaggle.com/datasets/dv1453/twitter-sentiment-analysis-analytics-vidya>
  - [9] Flask Documentation <https://flask.palletsprojects.com/en/2.1.x/>
  - [10] NLTK Documentation <https://www.nltk.org/>
-