# Software Design Document
# Fake Tweets Detection

Radwa Mostafa, Abdelrahman Tolba , Mariam Khaled, John Gerges

March 2, 2020

# 1 Introduction

## 1.1 Purpose

This software design document purpose is to fully describe the architecture of our web-based fake tweets detector.it will explain in details, the components of the system represented in the block diagram, the order of the project with sequence diagram, also the implementation of the project and its development will be shown in the class diagram . This software design document (SDD) is, therefore, intended for the stakeholders and developers of the our system

## 1.2 Scope

Detection of Fake news it has an enormous social impact as it can mislead citizens into believing misinformation about a specific product/ person or situation and the spread of economic and political fake news can directly affect the stock market and government economics in general. Also, Tourism income can be affected by fake news, as spreading fake news in countries can make tourists not interested in visiting certain places which costs losing much income.. This targets social media users that rely on social as their source of information.
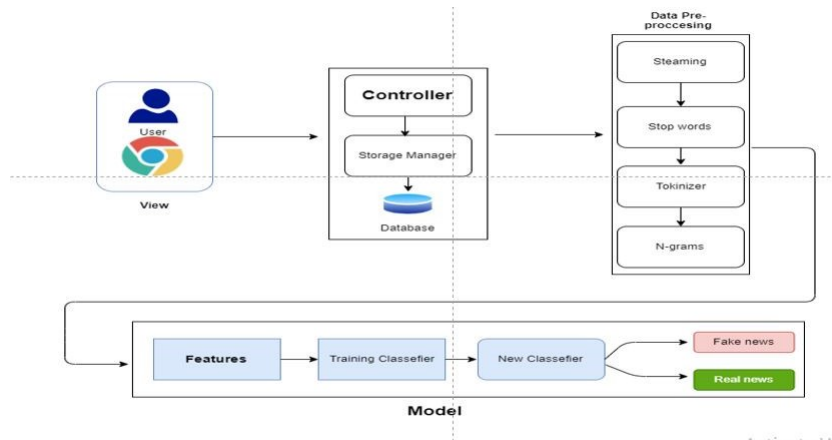
## 1.3 Overview

Nowadays Twitter is one of the main source of news for a huge number of people, and twitter has no restrictions on what people write so this produce a big percentage of fake news . Those fake news become one of the most critical problems in the 21 century as it can affect politics , society and economy . This document will explore the use of an artificially intelligent computer system to help enhance the human ability to differentiate between the fake and real news and will also highlights the features and algorithms helping figuring out the fake news with the best accuracy.
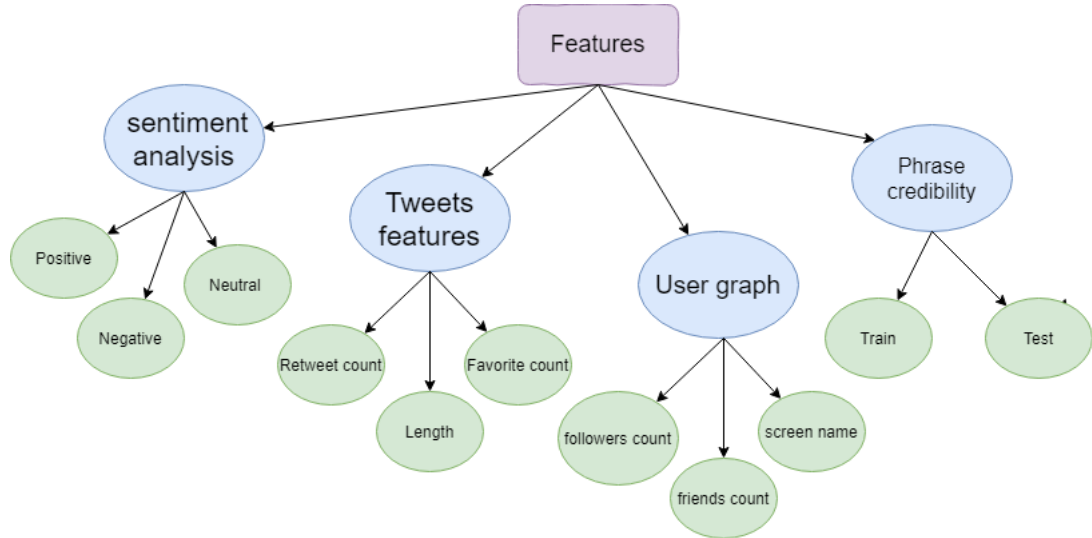
## 1.4  Definitions and Acronyms

- NLP :Natural language processing is a subfield of computer science, information

- API:An application programming interface (API) is a set of routines, protocols,and tools for building software applications

- GUI:Graphical user interface

- NLTK:Natural Language Toolkit [3] (NLTK) is a platform used for building Python programs that work with human language data for applying in statistical natural language processing (NLP)

- Python An interpreted high-level programming language for generalpurpose programming.

# 2  System Overview

### 2.0.1   Features



Our System aims to detect fake tweets by using machine learning and a combination of algorithms The proposed diagram using MVC model ( Model View Controller),The user's view a web based application ,entering a topic to search on then passing tweets related to this topic to the controller , this tweets needs pre-processing to be well defined data to work on like ( stemming , stop words , toknizer, n-grams). After that the controller will send the organized tweets to work on additional features :

- The first feature is the sentiment analysis which works on determining the effect of the tweet wither it's positive, negative or neutral.

- The second feature is extracting the tweet features itself like the retweet count,length of tweet and the favourite count.

- The third feature is the user graph which is analyzing the user's information like followers count and screen name.

- The final feature is the phrase credibility which training and testing using classifier

,those features differentiate between real and fake accounts and tweets on twitter .

# 3 System Architecture
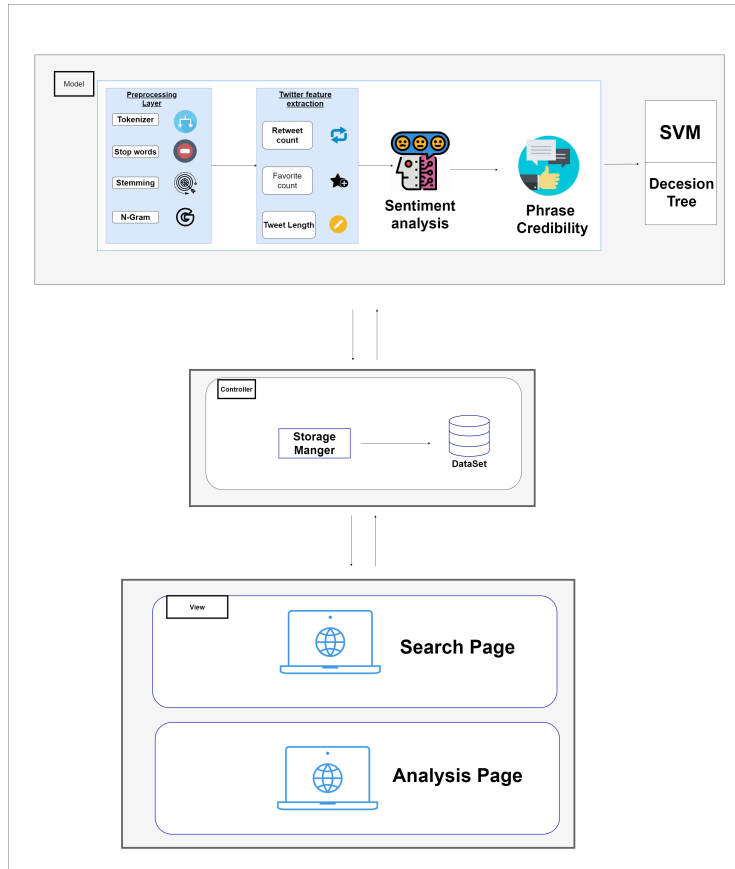
## 3.1 Architectural Design



Figure 1: Architectural Design

### 3.1.1 Model

The model part is responsible for the functionality of the system , which is first doing pre-processing on data to make it in the form that is suitable for the functions to use (Tokenizer, stop words, Stemming,N-Gram) and then extracting tweet features like: Retweet count, Favourite count and Tweet length , sentiment analysis and phrase credibility , the final step is entering those features is a hybrid classifier between SVM and Decision Tree those are the algorithms we will use and sends it to the controller to save it in the database.
1- Algorithms:

- SVM : SVM plot each data item as a point in n-dimensional space where n is the number of classes we have which is fake and real ,then we perform classification by finding the hyperplane that differentiate the classes very well.

- Decision Tree: Decision tree uses the tree representation to solve the problem in which each leaf node corresponds to a class label and attributes are represented on the internal node of the tree.

2-libraries:

- Twitter API:its REST API allows you to read and write Twitter data; in other words, it can be used to create new tweets, read user profiles and the data of followers (among other data from each profile), since it identifies the various Twitter applications and the users who register.

- NLTK: he Natural Language Toolkit (NLTK) is a platform used for building Python programs that work with human language data for applying in statistical natural language processing (NLP). It contains text processing libraries for tokenization, parsing, classification, stemming, tagging and semantic reasoning.

- SKLEARN : Scikit-learn is a free machine learning library for Python. It features various algorithms like support vector machine, random forests, and decision tree and it also supports Python numerical and scientific libraries like NumPy.

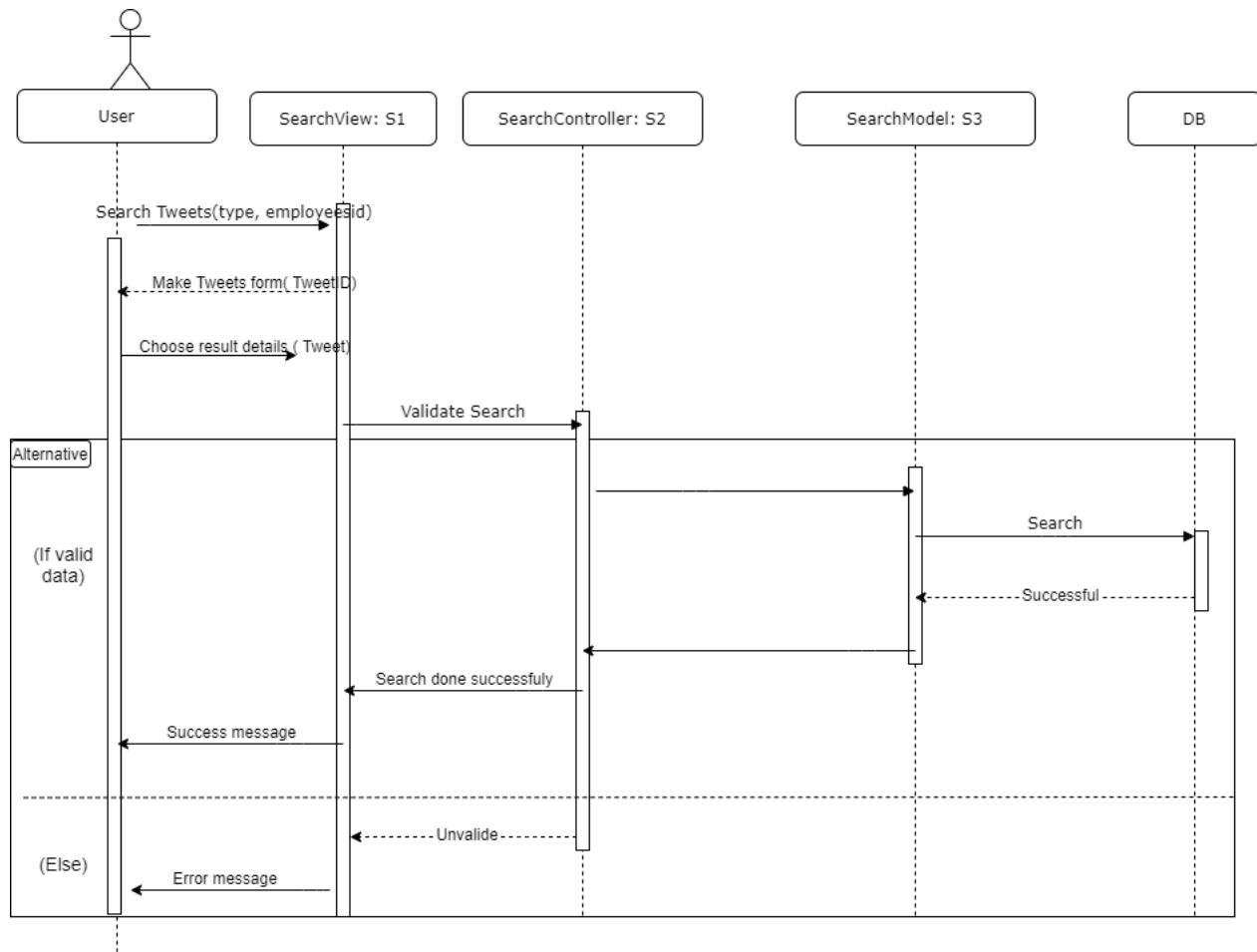- Numpy: This library is responsible for handling arrays.

### 3.1.2 view

It is responsible for the presentation of data and representing the User Interface(UI). We have two different interfaces one is responsible for retrieving the data from the user and the second one is responsible for displaying the output data for the user and the analysis related to the data.
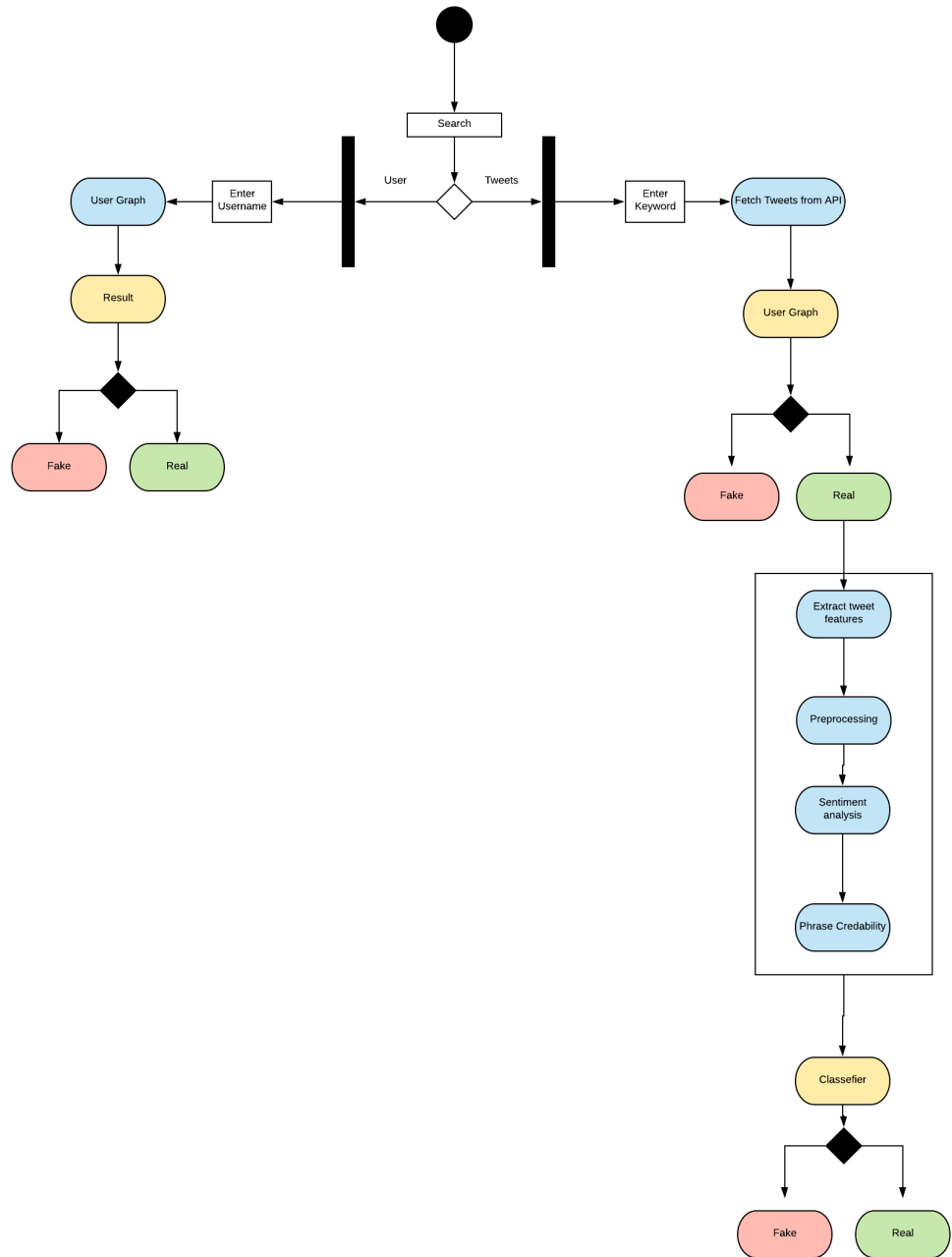
### 3.1.3   Controller

It is responsible for binding the view and model. The interactions and requests made within the view are taken and sent to the database to fetch data with the use of models then it forward data to the view again to be shown. The controller we have is: the user controller that deals with the user input that will be stored in the database and after applying our functionality on it, it then sends data to the view model to display the result for the user.

## 3.2   Decomposition Description

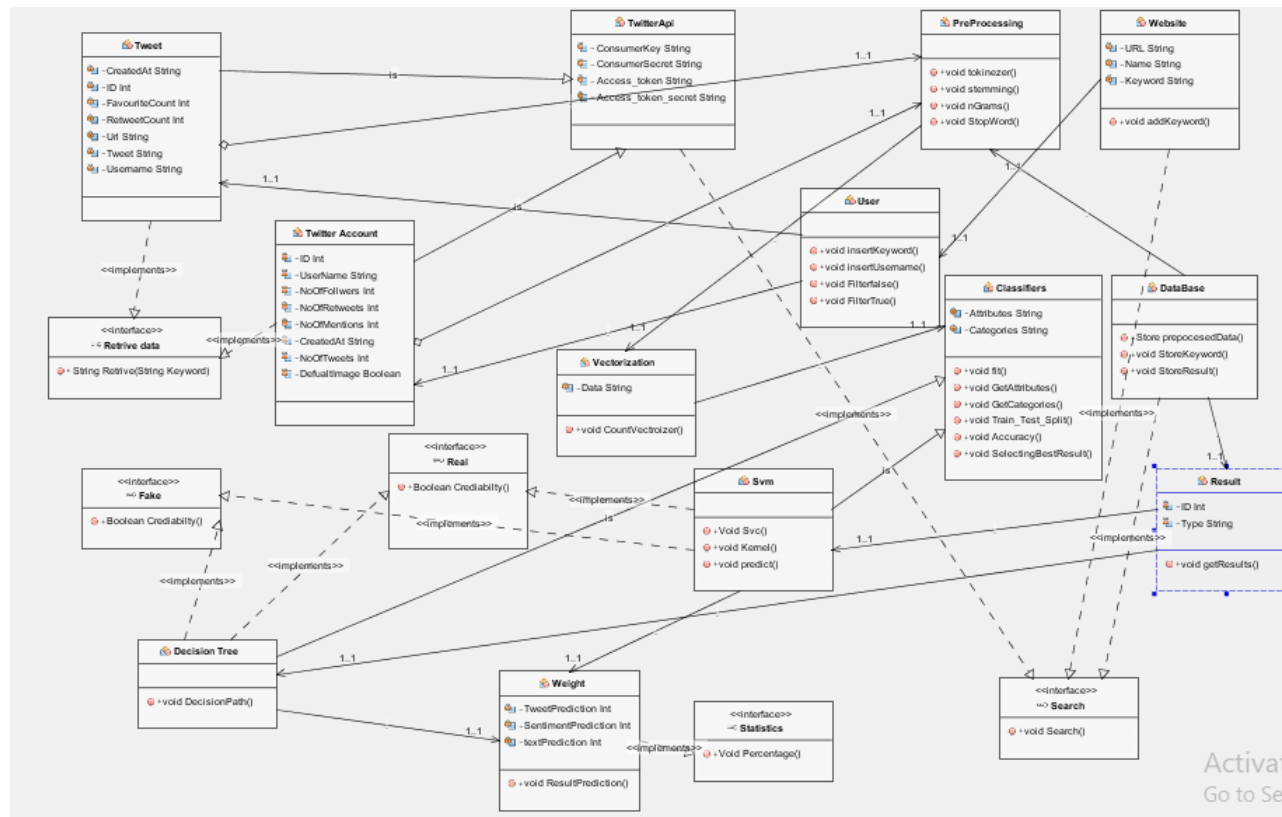### 3.2.1   System Sequence Diagram

## 3.3     Activity Diagram

## 3.4 User side

User first enters a username to search one and then we perform a user graph to show the result which is the user is credible or not

## 3.5 tweets side

user enters a tweet to search on and then the system fetch from twitter API returning the related tweets , after that we perform user graph and neglect the fake ones , the real user's tweets entered in the pre-processing stage , sentiment analysis , tweet features and the phrase credibility , all those stage's output entered in a classifier helping in differentiation between real and fake tweets

### 3.5.1 Class Diagram



class name:

## 3.6 Design Rationale

As mentioned previously, we have used Model-View-Controller (MVC) as our architecture as it helped us separate the functionality and data of our system

8

from the presentation. So, we can easily make modifications, re-use and optimize functionality part as it is our core. Also the software we are developing efficiency and accuracy is a very important aspect of it so it will be very sensitive with data so this should be developed in a very accurate and efficient way.

### 3.7 Possible Algorithms

There were so many alternative algorithms that we could have used like decision trees and support vector machine that are generally very good for text classification.

Decision Trees: It is a decision support tool that uses a tree-like model of 11 decisions and their possible consequences, including chance event outcomes, resource costs, and utility. It is one way to display an algorithm that only contains conditional control statements.

SVM: Linear SVM is given a set of train data which belong to a certain class to find an optimal separating line. It tries to maximize the distance between each class to avoid mis-classification. Then a test data are given to be classified to one of the classes formed before[5]. We have chooses the linear SVM as our classifier as after several experiments it was found out that the SVM gave the best result and most accurate with the highest f-measure which is technically the mean between the precision and the recall score.

## 4 Data Design

### 4.1 Data Description

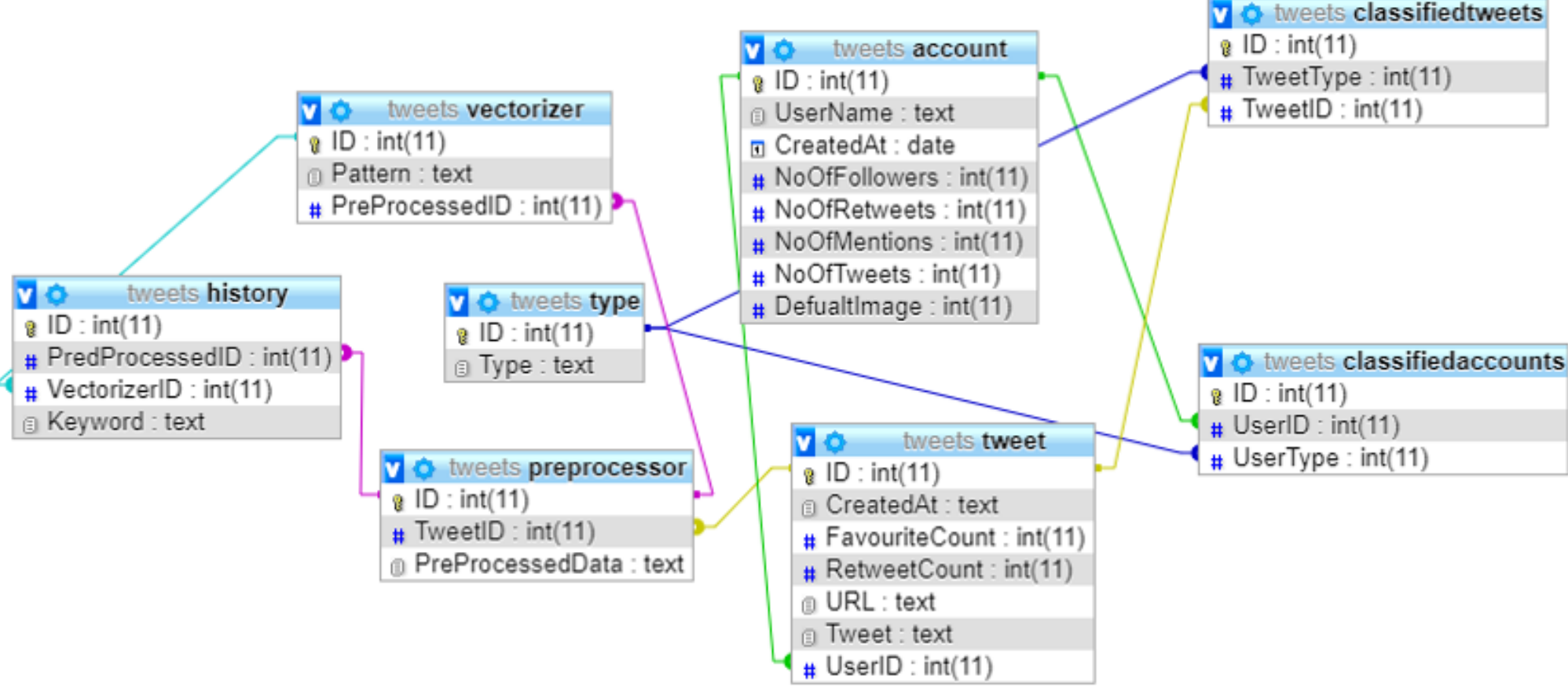

## 4.2   Data Dictionary

- UserAccount: this table saves all information about the user like number of people following him , number of retweets on his tweets .also saves the search result about any user being searched on.

- classifiedTweets: this tables store a the tweets Id and it's type Id weither it's fake or real. It has relations with tweet and type.

- classifiedAccounts:this tables store all the account Id and it's type Id weither it's fake or real.It has relation with account and type.

- Pre-processor: In this table we saves all the data after pre-processing done on it. it has relation with history, vectorizer and tweets.

- Tweets: this table saves all information about the tweets like number of favourites , number of retweets on this tweets. also saves the searched on tweets extracted from twitter API.

- vectroizerr:Saves all data after changing it to patterns which is a sequence of numbers referring to each word that's made to be recognized by a classifier .it has relation with history and preprocessor.

- Output:Saves the statistics of each tweet and the percentage of the trueness of the tweet so it has relation with tweets table.

- Sentiment:Saves the result of the sentiment analysis feature which is the tweet having positive , negative or neutral effect so it has relation with tweets table.

# 5 Component Design

## 5.1 Algorithms

1-SVM : "SVM" plot each data item as a point in n-dimensional space where n is the number of classes we have which is fake and real ,then we perform classification by finding the hyper-plane that differentiate the classes very well by maximize the distance between each class for decreasing error percentage.

2- Decision Tree : Decision tree uses the tree representation to solve the problem in which each leaf node corresponds to a class label and attributes are represented on the internal node of the tree.

## 5.2 Features

1- Sentiment Analysis: Sentiment Analysis is the automated process of analyzing text data and sorting it into sentiments positive, negative or neutral. Performing Sentiment Analysis on data from Twitter using machine learning can help in differentiating between tweets.

2-Tweets features: we gather information about each tweet and actions made upon them like the retweet count , length of tweet and the favourite count helping us determining the credibility of the tweet.

3-User graph:we gather information about each user and actions made by them like the friends count , screen name and the favourite count and representing it using a graph structure helping us determining the credibility of the person is he trusted or not.

4-Phrase credibility : Training and testing multiple data sets to help use in the classification.
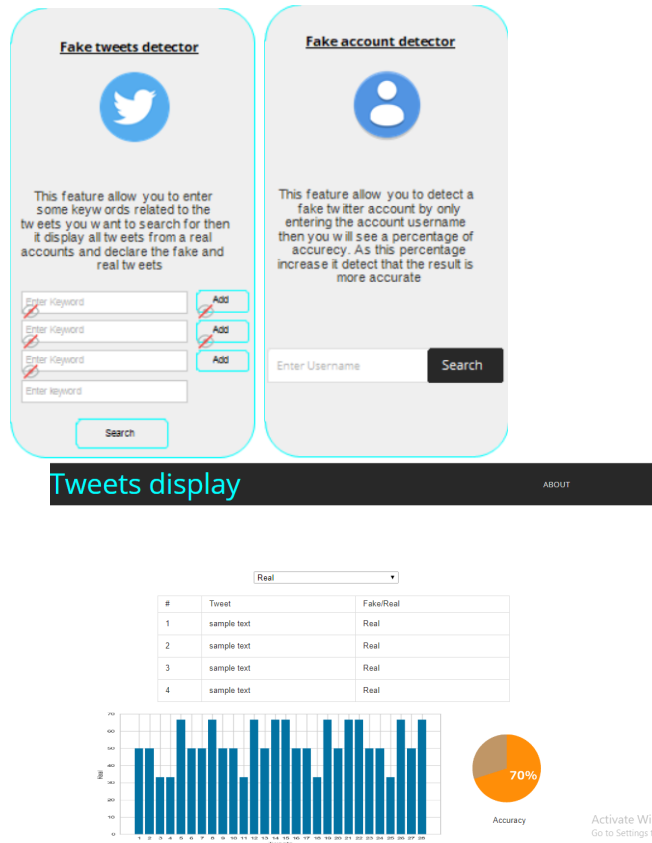
## 5.3    Data-set

| text | user_id | tweet_id | retweet_c | favorite_c | length | Sentimen | predictior | FinalLabel |
|---|---|---|---|---|---|---|---|---|
| Hurricane | 8.80E+17 | 242 | 10 | 1 | 90 | 1 | 1 | 1 |
| Students 1 | 18813355 | 122 | 6 | 0 | 132 | 1 | 1 | 1 |
| Harvey ca | 2.72E+08 | 215 | 8 | 5 | 112 | 2 | 1 | 1 |
| I added a | 33600664 | 212 | 3 | 2 | 89 | 2 | 0 | 1 |
| DFW busir | 18440701 | 237 | 4 | 0 | 97 | 2 | 1 | 1 |
| Stories of | 4.11E+08 | 66 | 0 | 0 | 138 | 2 | 0 | 0 |
| 9 Weeks A | 2.57E+09 | 50 | 5 | 3 | 81 | 1 | 0 | 1 |
| @BestFrie | 3.09E+08 | 204 | 3 | 0 | 144 | 2 | 1 | 1 |
| Companie | 5.37E+08 | 7 | 2 | 1 | 139 | 2 | 1 | 1 |
| The Energ | 74014041 | 102 | 0 | 0 | 107 | 2 | 1 | 0 |
| Hurricane | 28606058 | 35 | 0 | 0 | 87 | 2 | 1 | 0 |
| @sethme | 9.23E+17 | 164 | 0 | 0 | 127 | 0 | 1 | 0 |
| How Otis | 2.03E+08 | 226 | 1 | 2 | 108 | 2 | 1 | 1 |
| Jobs repo | 47673131 | 114 | 0 | 0 | 121 | 1 | 0 | 0 |
| @CharityI | 7.97E+17 | 122 | 0 | 0 | 140 | 1 | 1 | 1 |
| 12 nonpro | 36369382 | 16 | 0 | 0 | 117 | 2 | 1 | 0 |
| Hurricane | 85842228 | 241 | 0 | 0 | 78 | 2 | 1 | 0 |
| Bonfire W | 5.66E+08 | 232 | 0 | 1 | 140 | 1 | 1 | 0 |
| Hurricane | 8.84E+17 | 200 | 1 | 0 | 100 | 2 | 0 | 0 |
| super bov | 3.31E+09 | 34 | 0 | 0 | 88 | 1 | 1 | 0 |
| Congrats I | 2.54E+08 | 211 | 0 | 0 | 123 | 1 | 0 | 0 |
| Hurricane | 18073211 | 80 | 0 | 0 | 87 | 2 | 1 | 0 |

Our data set contain the data used for training and testing the classifiers,also when the user search for a tweet or a news the system fetch the data from Twitter API and store it in the data-set.Our data set contain 9 column :

- text: This column contain the tweets itself.

- user-id: The twitter user ID fetched from twitter API.

- tweet-id: Every Tweet has its unique ID.

- retweet-count: Number of retweets made on this tweet.

- favorite-count :How many times people favorite this tweet.

- Length: The length of tweet in characters.

- Sentiment: "1" means positive ,"0" means negative and "2" means neutral.

- predictor: "1" means tweet is credible and "0" means tweet is not credible.

- FinalLabel: "1" means tweet is Real and "0" means tweet is Fake.

# 6  Humnan Interface Design

## 6.1  Screen Images



## 6.2  Screen Objects and Actions

Our web-based system is constructed in a friendly way to all type of users , as it have a small number of instructions that can be easily memorized which give the user a privilege to navigate through our system pages smoothly . The user starts by entering a keyword to search on and can specify the max number of tweets he wants to search for , the way to show it and also the interval of time of the tweets , Also the user can enter a username to know wither this user is trust-worthy or not.

# 7  Requirements Matrix

Provide a cross reference that traces components and data structures to the requirements in your SRS document. Use a tabular format to show which system

components satisfy each of the functional requirements from the SRS. Refer to the functional requirements by the numbers/codes that you gave them in the SRS.

# 8    References