

## Patient Care analysis - HealthCare

### Problem statement

Health Care Industry desires to classify the Patients using their pathology data for their CARE (Self-management, Doctor-Advise, Further Diagnostic and Chronic Medication) management improvement that facilitates to build a multi-classification model to build **CARE Management Model (CMM)** with right classification of patient.

### Tools Used:

**Spark Cluster**  
**Spark MLLIB (Machine Learning Library)**  
**R & Python Language**

### Apache Spark

Apache Spark is a lightning-fast cluster computing technology, designed for fast computation. It is based on Hadoop MapReduce and it extends the MapReduce model to efficiently use it for more types of computations, which includes interactive queries and stream processing. The main feature of Spark is its in-memory cluster computing that increases the processing speed of an application. Spark is designed to cover a wide range of workloads such as batch applications, iterative algorithms, interactive queries and streaming. Apart from supporting all these workload in a respective system, it reduces the management burden of maintaining separate tools.

### Spark SQL

Spark SQL is a component on top of Spark Core that introduces a new data abstraction called Schema RDD, which provides support for structured and semi-structured data.

### Spark Streaming

Spark Streaming leverages Spark Core's fast scheduling capability to perform streaming analytics. It ingests data in mini-batches and performs RDD (Resilient Distributed Datasets) transformations on those mini-batches of data.

### MLlib (Machine Learning Library)

MLlib is a distributed machine learning framework above Spark because of the distributed memory based Spark architecture. It is, according to benchmarks, done by the MLlib developers against the Alternating Least Squares (ALS) implementations. Spark MLlib is nine times as fast as the Hadoop

disk-based version of Apache Mahout(before Mahout gained a Spark interface).

## **GraphX**

GraphX is a distributed graph-processing framework on top of Spark. It provides an API forexpressing graph computation that can model the user-defined graphs by using Pregel abstraction

## **Python**

Python is a dynamic, interpreted (bytecode-compiled) language. There are no type declarations of variables, parameters, functions, or methods in source code. This makes the code short and flexible, and you lose the compile-time type checking of the source code. Universally, Python has gained a reputation because of it's easy to learn. The syntax of Python programming language is designed to be easily readable. Python has significant popularity in scientific computing. The people working in this field are scientists first, and programmers second. Nowadays working on bulk amount of data, popularly known as big data. The more data you have to process, the more important it becomes to manage the memory you use. Here Python will work very efficiently.

## **Scala**

If you're the kind of data scientist who deals with large datasets, Scala will be invaluable. It's practically the de facto language for the current Big Data tools like Apache Spark, Finagle, Scalding, etc. Many of the high performance data science frameworks that are built on top of Hadoop usually are written and use Scala or Java. The reason Scala is used in these environments is because of its amazing concurrency support, which is key in parallelizing a lot of the processing needed for large data sets. It also runs on the JVM, which makes it almost a no-brainer when paired with Hadoop. Like Java, Scala is object-oriented, and uses a curly-brace syntax reminiscent of the C programming language. Unlike Java, Scala has many features of functional programming languages like Scheme, Standard ML and Haskell, including currying, type inference, immutability, lazy evaluation, and pattern matching.

## Problem Understanding

This study is commissioned with the following objectives:

The care management model system gives computer the ability to learn without being explicitly programmed. (Machine Learning) To classify patients based on different criteria. The Care management models will facilities for

- 1) Personalized Medicine
- 2) Predictive Analytics and Preventive Measures
- 3) Self-Motivated Care
- 4) Disease Modelling and Mapping

## Sample Source:

### Naive Bayes Classification Algorithm For Data Accuracy

#### NB.py

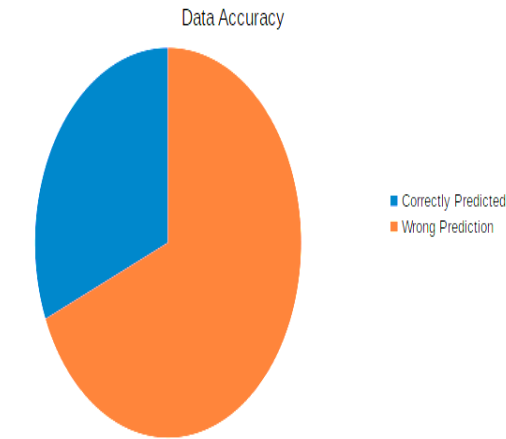
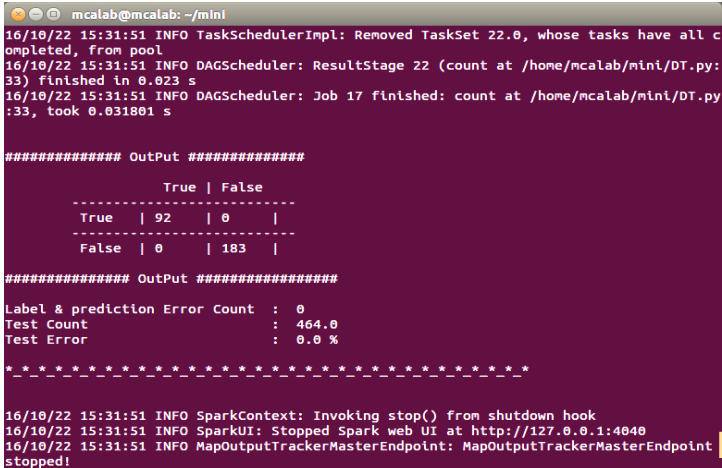
```
# Import required packages
from __future__ import print_function, division
from pyspark.mllib.classification import NaiveBayes, NaiveBayesModel
from pyspark.mllib.util import MLUtils
from pyspark import SparkContext
# Load and parse the data file.
sc = SparkContext("local", "classification")
data = MLUtils.loadLibSVMFile(sc, '/home/mcalab/mini/HealthSVM.txt')
# Split data approximately into training (50%) and test (50%)
training, test = data.randomSplit([0.5, 0.5])
# Train a naive Bayes model.
model = NaiveBayes.train(training, 1.0)
predictions = model.predict(test.map(lambda x: x.features))
predictions.foreach(print)
predictions.coalesce(1).saveAsTextFile("/home/mcalab/mini/output1")
# Make prediction and test accuracy.
predictionAndLabel = test.map(lambda p: (model.predict(p.features), p.label))
predictionAndLabel.foreach(print)
cp = predictionAndLabel.filter(lambda (x, v): x == v).count()
tc = test.count()
accuracy = (cp/tc)*100
print('\n')
print('*****Output*****')
```

SCREEN SHOTS

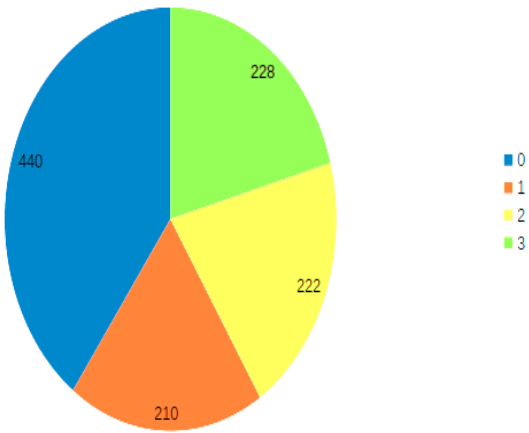
Naive Bayes Classification Algorithm For Data Accuracy

Running the problem on Spark

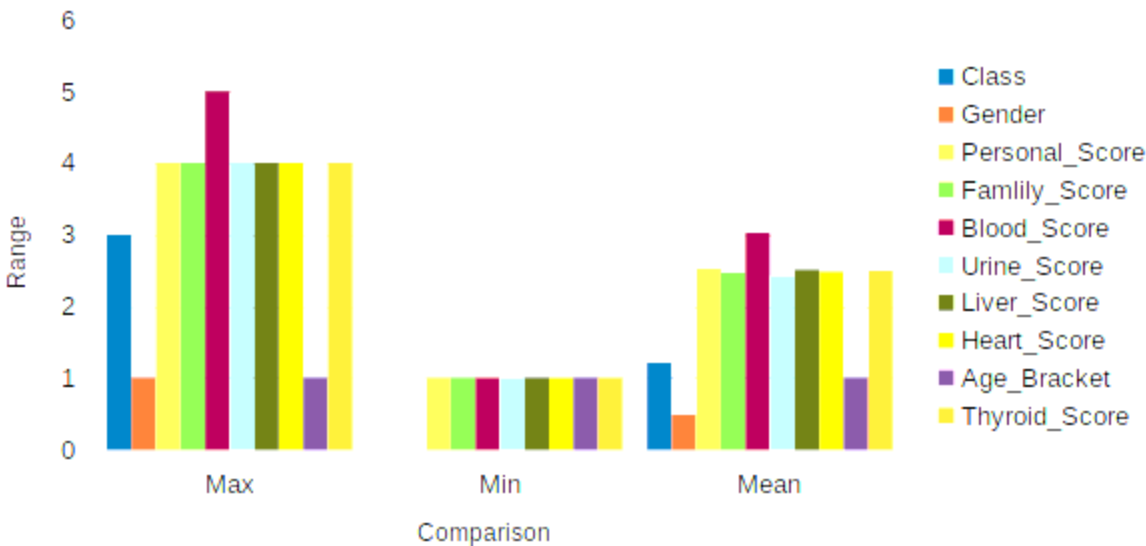
Output Confusion Matrix



Patient Break-ups on Blood and Heart Scores



Summary Report on data Scores



## **CONCLUSION**

Big data is the reality and is going to stay there for a long time. While have a care management model by any hospital they can predict the future condition of any patient without any wait. Time is a major problem for health care. Within the time anything may happens. So life is saved by care management model. So the model created by this project is so important in the everyday life. But the project model wants to more accurate. The only way to improve the accuracy of this model is collect more data and applies it for modelling. So in future hospital have to concentrate on big data preparation and management in order to create a good model for the system. The project suggest to gain success hospital need to change service manner and increase service on patients who are professionally unbalanced and try to efficiently utilise care management model for patients, serious and quick action needed when pathology values get changing abnormally.

## **Customer Segmentation - Banking**

### **Problem:**

BANK wants to sell the new product in insurance/loans so analyse and report customer's money spending trends. This segmentation/information can be used to provide solution to meet the requirement of the customer which might be very unique and product can be offered.

### **Tools Used:**

- Hadoop Cluster**
- Sqoop in Hadoop**
- Mysql database**
- PIG scripts in Hadoop**
- Hive Data warehouse in Hadoop**
- R Language**

### **Sqoop in Hadoop Framework**

The evaluated data from MYSQL is being imported using Sqoop into Hadoop environment for further tuning and analysis. The sqoop import command is used to bring the table from MYSQL to HDFS.

### **Pig**

The imported data from sqoop is given an input to the Pig grunt shell. The log data is loaded in to pig shell using load command. The loaded log data may contain some of the text data in an unstructured format. These unstructured data are cleansed using pig filter commands. This fine-tuned data are stored as comma separated csv(Comma Separated) file for further analysis using pig storage.

### **Hive**

The fine tuned data which is in csv format is loaded into a database inside hive. A table with suitable data type is created inside the database using HQL create query, after analysing the data. The data in the csv file is loaded in to the table created for taking adhoc reports using select query. The result values are stored in to HDFS as csv file for further analysis in R.

### **R LANGUAGE**

The resultant values stored as csv file from Hive based on the query is taken as an input to R using read command. The loaded data is used to create

summary of the data for analysis, and for creation of the Chart or Graphs for Visualization and easy understanding about the analysis.

## MySQL

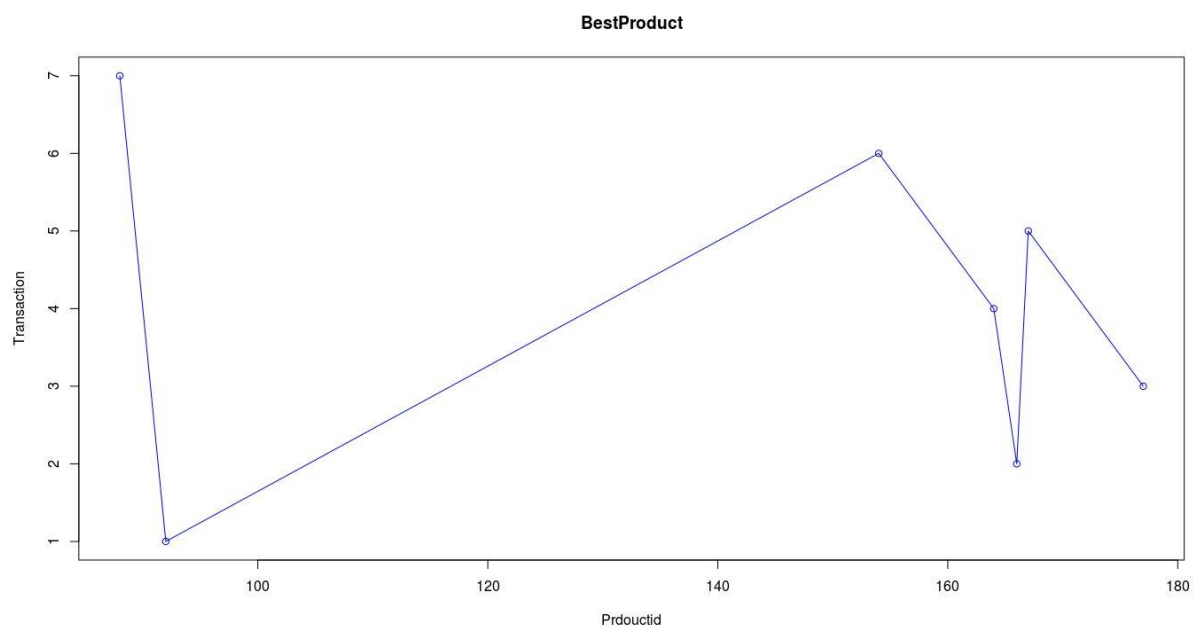
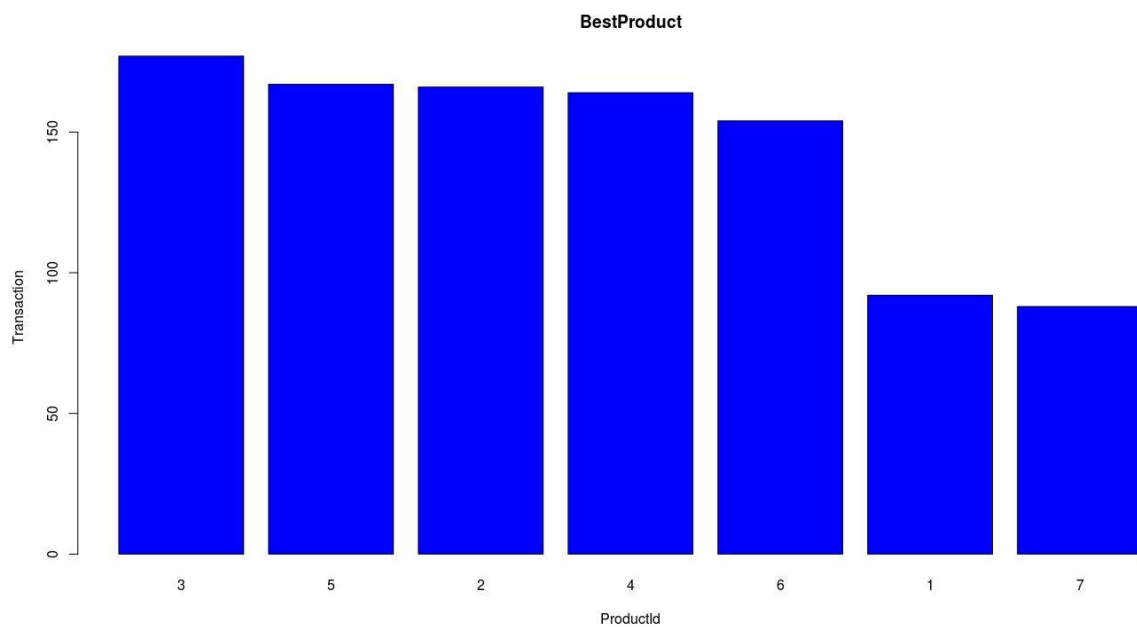
Un-structured tower data is being captured from tower data management system. The log file containing the tower utilization details of different towers is copied into mySql. A single linear table is created in mySql

### DATA PROFILING

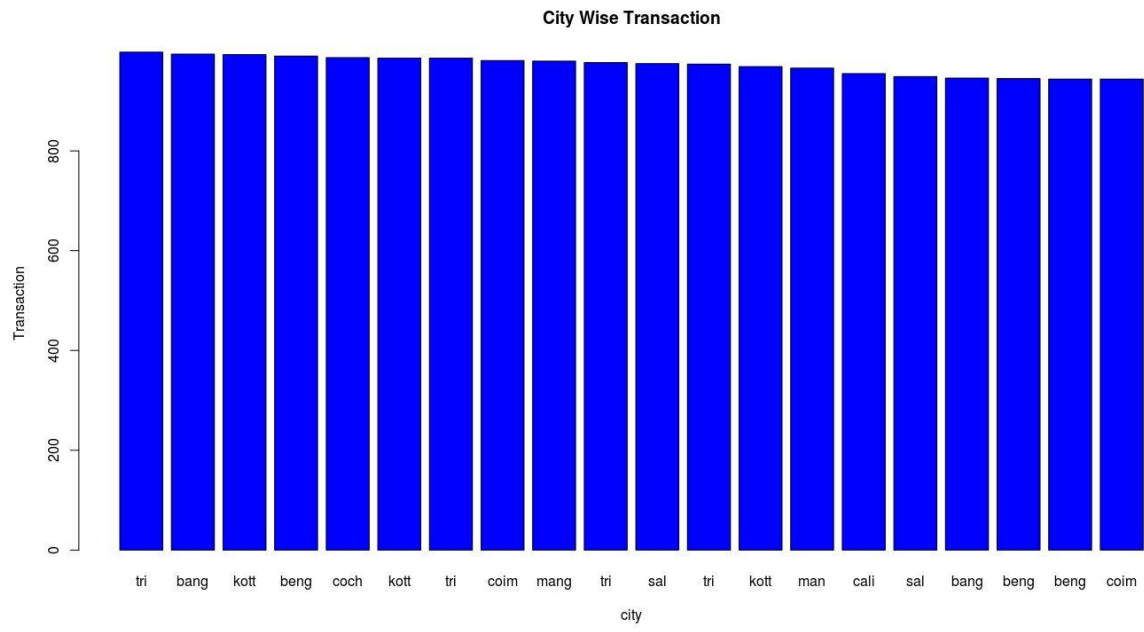
Product Item	No of Records	Transaction Amount
Savings Account	183	92609400
Current Account	321	166077747
Term Deposits	350	177386416
Loans	338	164462435
Credit Cards	331	167896232
Gold	312	154077082
Insurance	165	88702632
<b>Total</b>	<b>2000</b>	<b>1011211944</b>

copied for quick process and data evaluation using my-sql capabilities.

```
Terminal
(27444,1,2,2005-02-01,136838,Tamilnadu,Trichy)
(27445,1,4,2005-02-01,568581,Kerala,Kottayam)
(27446,1,5,2005-02-01,816446,Karnataka,Banglore)
(27447,1,1,2005-02-01,18485,Kerala,Calicut)
(27448,2,5,2005-02-01,181528,Kerala,Kottayam)
(27449,1,5,2005-02-01,987271,Kerala,Kottayam)
(27450,2,3,2005-02-01,615058,Karnataka,Banglore)
(27451,3,2,2005-02-01,809993,Kerala,Calicut)
(27452,3,4,2005-02-01,968640,Kerala,Trivandrum)
(27453,1,5,2005-02-01,302360,Tamilnadu,Trichy)
(27454,1,4,2005-02-01,372003,Karnataka,Banglore)
(27455,1,3,2005-02-01,622840,Kerala,Kannur)
(27456,2,4,2005-02-01,425988,Karnataka,Banglore)
(27457,2,3,2005-02-01,790622,Kerala,Kottayam)
(27458,3,6,2005-02-01,974075,Kerala,Kannur)
(27459,1,4,2005-02-01,666589,Tamilnadu,Trichy)
(27460,3,6,2005-02-01,251213,Kerala,Calicut)
(27461,1,4,2005-02-01,42318,Karnataka,Banglore)
(27462,1,7,2005-02-01,260082,Kerala,Calicut)
(27463,1,3,2005-02-01,193793,Kerala,Kottayam)
(27464,2,7,2005-02-01,357830,Karnataka,Banglore)
(27465,3,6,2005-02-01,708180,Karnataka,Banglore)
(27466,2,3,2005-02-01,249834,Kerala,Trivandrum)
(27467,2,2,2005-02-01,619721,Karnataka,Banglore)
(27468,2,6,2005-02-01,220852,Kerala,Kottayam)
(27469,3,1,2005-02-01,601814,Tamilnadu,Trichy)
(27470,1,3,2005-02-01,907790,Kerala,Trivandrum)
(27471,1,2,2005-02-01,121255,Karnataka,Banglore)
(27472,1,5,2005-02-01,871979,Kerala,Calicut)
(27473,3,4,2005-02-01,997915,Karnataka,Banglore)
(27474,2,5,2005-02-01,528934,Tamilnadu,Trichy)
(27475,1,2,2005-02-01,846494,Kerala,Kottayam)
(27476,3,2,2005-02-01,134360,Karnataka,Banglore)
(27477,1,2,2005-02-01,818272,Tamilnadu,Trichy)
(27478,3,3,2005-02-01,132362,Kerala,cochin)
(27479,1,5,2005-02-01,912515,Karnataka,Banglore)
(27480,1,7,2005-02-01,254013,Tamilnadu,Trichy)
(27481,2,4,2005-02-01,861060,Kerala,Kannur)
(27482,2,5,2005-02-01,998011,Tamilnadu,Trichy)
(27483,2,2,2005-02-01,603360,Kerala,Kannur)
grunt> dump finetuneddata;
```







## **Sentiment Analysis – Text mining**

### **Problem statement**

Text mining is gaining importance due to problem of discovering useful information from the data deluge that the organizations are facing today. This project intends to present a broad overview of text mining and its components and techniques and their use in various business applications. This project gives a description about text mining and the reasons for its increased importance over the years. This is followed with presenting a generic process framework for text mining and describes its different components and sub-components, business applications, and brief description of text mining tools available in the market. Text mining involves the pre-processing of document collections (text categorization, feature/term extraction, etc.), the storage of the intermediate representations, the techniques to analyze these intermediate representations (such as distribution, analysis, clustering, trend analysis, and association rules), and visualization of the results.

### **DATA PREPARATION**

#### **Data Set Description (Select, Clean, Construct, Integrate and Format Data)**

This data's are collected from twitter and doing some cleansing and analysis. After it going to be text mining by using python or R to do some process and finally we get exact output from the result.

#### **Text Mining Framework Components**

The different stages in the text mining framework are described below:

#### **Textual Data Sources**

The textual data is available in numerous internal and external data source like electronic text, call centre logs, social media, corporate documents, research papers, application forms, service notes, emails, etc.

#### **Pre-processing**

Pre-processing tasks include methods to collect data from the disparate data sources. This is the preliminary step of identifying the textual information for mining and analysis. Pre-processing tasks apply various feature extraction methods against the data. Pre-processing tasks include different types of techniques to transform the raw, unstructured, original format data into structured, intermediate data format. Knowledge discovery operations are conducted against the structured intermediate data.

For the preparation of unstructured data into a structured data format, different techniques are needed than those of traditional data mining systems where the knowledge discovery is done against the structured data sources. Various preprocessing techniques exist and can be used in combination to create structured data representation from raw textual data. Therefore different combinations of techniques can be used based on the type of the raw textual data.

## **Text Cleansing**

Text cleansing is the process of cleansing noisy text from the textual sources. Noisy textual data can be found in SMS, email, online chat, news articles, blogs and web pages. Such text may have spelling errors, abbreviations, non-standard terminology, missing punctuation, misleading case information, as well as false starts, repetitions, and special characters.

Noise can be defined as any kind of difference in the surface form of an electronic text from the original, intended or actual text. The text used in the in short message service (SMS) and on-line forums like twitter, chat and discussion boards and social networking sites is often distorted mainly because the recipients can very well understand the shorter form of the longer words and also reduces the time and effort of the sender. Most of the text is created and stored so that humans can understand it, and it is not always easy for a computer to process that text.

## **Removing stop words**

Stop words are words which are filtered before or after processing of textual data. There is not one definite list of stop words which all tools use, if even used. Some tools specifically avoid removing them to support phrase search. The most common stop words found in the text are “the”, “is”, “at”, “which” and “on”. These kinds of stop words can sometimes cause problems when looking for the phrases that include them. Some search engines remove some of the most common words from the query on order to improve performance.

## **Tools Used:**

### **Python Language**

#### **Python**

Python is a dynamic, interpreted (bytecode-compiled) language. There are no type declarations of variables, parameters, functions, or methods in source code. This makes the code short and flexible, and you lose the compile-time type checking of the source code. Universally, Python has gained a reputation because of it's easy to learn. The syntax of Python programming language is designed to be easily readable. Python has significant popularity in

	A	B	C	D	E	F	G	H
1	text		favorite	favoriteCount	replyToSN	created	truncated	replyToSID
2	1 RT @dr3amscometrue: #ilic #cdallas was lit ♥️👏		FALSE	0 NA		2016-09-27 05:19:52	FALSE	NA
3	2 @Dylan_Butler I think he saving his seat for the #LIC		FALSE	0 Dylan_Butler		2016-09-27 02:31:11	FALSE	7805785080
4	3 These raves this month!n read: 1 #GTA #LIC #SAFEINSOUND https://t.co/2EvaG76rWv		FALSE	4 NA		2016-09-27 02:25:15	FALSE	NA
5	4 RT @ZorlicRivera: After the rain but before the point at LIC #festibestie #LIC #feincolor #rave #EDM #PLUR #kandikid https://t.co/		FALSE	0 NA		2016-09-27 02:11:01	FALSE	NA
6	5 RT @shylyinynteshy: friends who headbong together stay together @ #LIC @BrittaniBeasley @chick_11 @syndreysms4 http://f		FALSE	0 NA		2016-09-27 01:42:24	FALSE	NA
7	6 After the rain but before the point at LIC #festibestie #LIC #feincolor #rave #EDM #PLUR #kandikid https://t.co/Qp0w9C		FALSE	13 NA		2016-09-27 01:25:13	FALSE	NA
8	7 RT @JimmyVanBramer: Great joining Borough President @MelindaKatz for a fundraiser at @S_Lewandowski's lovely home! #		FALSE	0 NA		2016-09-27 01:20:13	FALSE	NA
9	8 RT @JimmyVanBramer: Great to be at the @creekandcave hosting a debate watch party. Btw - #imWithHer #debateight #ide#		FALSE	0 NA		2016-09-27 01:08:05	FALSE	NA
10	9 RT @JimmyVanBramer: Great to be at the @creekandcave hosting a debate watch party. Btw - #imWithHer #debateight #ide#		FALSE	0 NA		2016-09-27 01:02:14	FALSE	NA
11	10 RT @JimmyVanBramer: Great to be at the @creekandcave hosting a debate watch party. Btw - #imWithHer #debateight #ide#		FALSE	0 NA		2016-09-27 01:02:00	FALSE	NA
12	11 Great joining Borough President @MelindaKatz for a fundraiser at @S_Lewandowski's lovely home! #LIC https://t.co/w		FALSE	0 NA		2016-09-27 01:01:35	FALSE	NA
13	12 RT @shylyinynteshy: friends who headbong together stay together @ #LIC @BrittaniBeasley @chick_11 @syndreysms4 http://f		FALSE	0 NA		2016-09-27 00:45:13	FALSE	NA
14	13 friends who headbong together stay together @ #LIC @BrittaniBeasley @chick_11 @syndreysms4 https://t.co/uyTE79oW		FALSE	7 NA		2016-09-27 00:28:24	FALSE	NA
15	14 RT @JimmyVanBramer: Great joining Borough President @MelindaKatz for a fundraiser at @S_Lewandowski's lovely home! #		FALSE	0 NA		2016-09-27 00:26:46	FALSE	NA
16	15 View of Long Island City from Manhattan tonight #qantigray #ilic #longislandcity #y: queens... https://t.co/SGLjF6wPk		FALSE	0 NA		2016-09-27 00:06:06	FALSE	NA
17	16 RT @NiallOfficial: Oh some Craic last night at my house with the LIC I Few drinks, father ted box set! What else would ya wan		FALSE	0 NA		2016-09-26 23:50:35	FALSE	NA
18	17 RT @JimmyVanBramer: Great joining Borough President @MelindaKatz for a fundraiser at @S_Lewandowski's lovely home! #		FALSE	0 NA		2016-09-26 23:49:32	FALSE	NA
19	18 Great joining Borough President @MelindaKatz for a fundraiser at @S_Lewandowski's lovely home! #LIC https://t.co/w		FALSE	0 NA		2016-09-26 23:45:47	FALSE	NA
20	19 RT @Wichitschick: Please help me find the loml @LICDallas #icdallas @TexasRaves @DallasEDMfamily #lic https://t.co/2P		FALSE	0 NA		2016-09-26 23:18:48	FALSE	NA
21	20 Please help me find the loml @LICDallas #icdallas @TexasRaves @DallasEDMfamily #lic https://t.co/2P224TFYG		FALSE	3 NA		2016-09-26 23:16:38	FALSE	NA
22	21 #ilic #cdallas was lit ♥️👏		FALSE	1 NA		2016-09-26 22:01:18	FALSE	NA
23	22 #ilicdallas #ilic #feincolor https://t.co/aU6f1RXi5p		FALSE	1 NA		2016-09-26 22:00:47	FALSE	NA
24	23 We thought they were going to cancel on us! #icdallas #ilic #ifeincolor https://t.co/vV8Tvm1tkS		FALSE	1 NA		2016-09-26 21:59:56	FALSE	NA
25	24 RT @mainimppals: I felt pretty in pink last night @LIC https://t.co/KBkfw6XqQr		FALSE	0 NA		2016-09-26 21:25:48	FALSE	NA
26	25 The new office is pretty dope. #YasQueens #LIC #NYC #Madewell https://t.co/iecmUj0D9j		FALSE	0 NA		2016-09-26 21:16:47	FALSE	NA
27	26 RT @AlewfeNYC: Our #mealtime Monday special Eggplant Parmigiana #alewfeNYC #ilic @ Alewife NYC https://t.co/3nzUoFte		FALSE	0 NA		2016-09-26 20:34:52	FALSE	NA
28	27 #ilic #cdallas was lit ♥️👏		FALSE	0 NA		2016-09-26 20:28:29	FALSE	NA
29	28 I've worked here for three years. Never knew I had roof top access. #ilic #longislandcity https://t.co/EJL4sH0b		FALSE	0 NA		2016-09-26 20:13:56	FALSE	NA
30	29 RT @LICMarket: Zucchini with lemon, basil, grana padano and a light dusting on breadcrumb #divine #ilic #licmarket https://t.c		FALSE	0 NA		2016-09-26 20:13:55	FALSE	NA
31	30 Our #licmarket Monday special Eggplant Parmigiana #alewfeNYC #ilic @ Alewife NYC https://t.co/oeHx453Hee		FALSE	0 NA		2016-09-26 19:58:38	FALSE	NA
32	31 Boom! #LIC's tallest tower in 2017 will be The Harrison, a 120 unit, 27 story tall luxury condo development!... https://t.co/ABAP7		FALSE	0 NA		2016-09-26 19:08:35	TRUE	NA
33	32 AFTERNOON !! dog walk #LIC https://t.co/KJL2j0LgR		FALSE	0 NA		2016-09-26 18:45:05	FALSE	NA
34	33 RT @jeanmarieEvilly: "We knew that we didn't want this completely sanitized neighborhood." Court Square, #LIC has new civi		FALSE	0 NA		2016-09-26 18:32:13	FALSE	NA

```
from tweepy import Stream
```

```
from tweepy import OAuthHandler
```

```
from tweepy.streaming import StreamListener
```

#use consumer key, consumer secret, access token, access secret as follows.

```
ckey="28befderyHlepe25CA"
```

```
csecret="u3xE1SdgQQxzOd97BzgssrygXZ3ezih4MbirBQO6yzJ9SVQGm"
```

```
atoken="468010305-KDYJJvdbey6pxrX9ssdgteI26YmhBxHjGwuknEmsuGBwJ"
```

```
asecret="0R8NEyFPasssPKfcP2sdstbDoQ6OsstejbUXHTpGCxg9pAh1i"
```

```
class listener(StreamListener):
```

```
def on_data(self, data):
```

```
print(data)
```

```
return(True)
```

```
def on_error(self, status):
```

```
print status
```

```
auth = OAuthHandler(ckey, csecret)
```

```
auth.set_access_token(accessToken, secret)
```

```
twitterStream = Stream(auth, listener())
```

```
tweet = twitterStream.filter(track=["#SBI"])
```

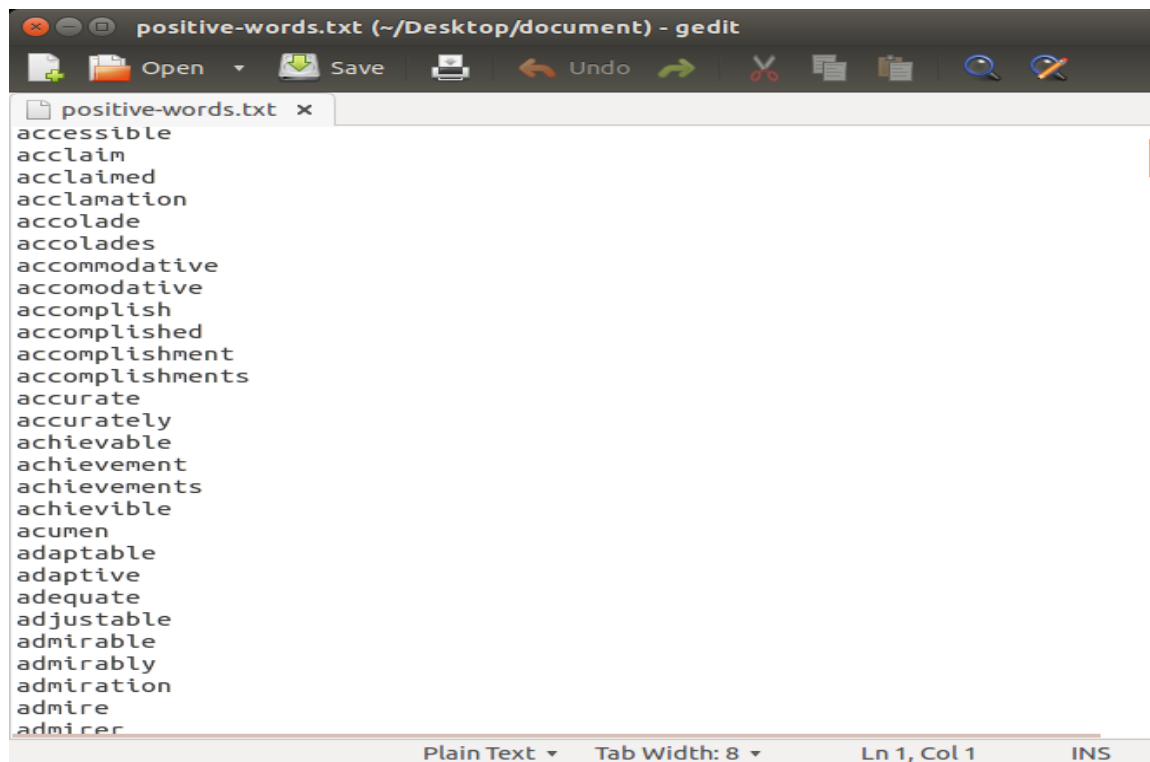


Figure 2

