

Name: Purna Sunil Jadhav

SapId: 60004220127

Batch: C2-2

Course: Big Data Infrastructure laboratory

Course Code: DT19CEEL6011

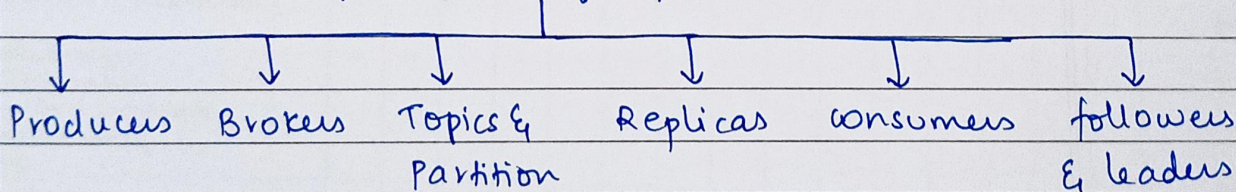
EXPERIMENT 10

AIM: Perform sentiment analysis using Kafka.

THEORY: Data reasoning is the process of transmitting a continuous data (also known as streams) typically fed into stream processing software to derive valuable insights.

Apache Kafka is an open source, distributed streaming platform that enables the development of real time, event driven application.

Components of Apache Kafka



- Producers: Producers in Kafka publish messages to one or more topics.

- Brokers: A Kafka cluster comprises one or more servers that are known as brokers. Broker works as a container that can hold multiple topics

- Topic: A stream of messages that are a part of specific category or feed name is referred to as a kafka topic.
- Partitions: Topics in kafka are divided into a configurable no. of parts, which are known as partitions.
- Replicas: Replicas are like backups for partition in kafka.
- Leaders & followers: Every partition will have one server that plays the role of a leader for that partition. Leader will perform read and write operations. Followers will replicate the data of the leader.

CONCLUSION: Thus, we performed sentiment analysis using kafka.



Name:	Prerna Sunil Jadhav
Sap Id:	60004220127
Class:	T. Y. B. Tech (Computer Engineering)
Course:	Big Data Infrastructure Laboratory
Course Code:	DJ19CEEL6011
Experiment No.:	10

AIM: Perform Sentiment Analysis using Kafka.

Step 1 : As the Log Data is unstructured, we parse and create a structure from each line, which will in turn become each row while analysis.

```
1 import re
2 from pyspark.sql import Row
3 # This is the regex which is specific to Apache Access Logs parsing, which can be modified according to
  different Log formats as per the need
4 # Example Apache log line:
5 # 127.0.0.1 - - [21/Jul/2014:9:55:27 -0800] "GET /home.html HTTP/1.1" 200 2048
6 # 1:IP 2:client 3:user 4:date time 5:method 6:req 7:proto 8:respcode 9:size
7 APACHE_ACCESS_LOG_PATTERN = '^(\S+) (\S+) (\S+) \[([^\w:/]+\s[+-]\d{4})\] \"(\S+) (\S+) (\S+)\" (\d{3})
  (\d+)\"
8
9 # The below function is modelled specific to Apache Access Logs Model, which can be modified as per
  needs to different Logs format
10 # Returns a dictionary containing the parts of the Apache Access Log.
11 def parse_apache_log_line(logline):
12     match = re.search(APACHE_ACCESS_LOG_PATTERN, logline)
13     if match is None:
14         raise Error("Invalid logline: %s" % logline)
15     return Row(
16         ip_address = match.group(1),
17         client_idend = match.group(2),
18         user_id = match.group(3),
19         date = (match.group(4)[:6]).split(":", 1)[0],
20         time = (match.group(4)[:6]).split(":", 1)[1],
21         method = match.group(5),
22         endpoint = match.group(6),
23         protocol = match.group(7),
24         response_code = int(match.group(8)),
25         content_size = int(match.group(9))
26     )
```

Step 2: Create Spark Context, SQL Context, DataFrame (is a distributed collection of data organized into named columns. It is conceptually equivalent to a table in a relational database)



```

1 from pyspark import SparkContext, SparkConf
2 from pyspark.sql import SQLContext
3 import apache_access_log # This is the first file name , in which we created Data Structure of Log
4 import sys
5
6 # Set up The Spark App
7 conf = SparkConf().setAppName("Log Analyzer")
8 # Create Spark Context
9 sc = SparkContext(conf=conf)
10 #Create SQL Context
11 sqlContext = SQLContext(sc)
12
13 #Input File Path
14 logFile = 'Give Your Input File Path Here'
15
16 # .cache() - Persists the RDD in memory, which will be re-used again
17 access_logs = (sc.textFile(logFile)
18               .map(apache_access_log.parse_apache_log_line)
19               .cache())
20
21 schema_access_logs = sqlContext.createDataFrame(access_logs)
22 #Creates a table on which SQL like queries can be fired for analysis
23 schema_access_logs.registerTempTable("logs")

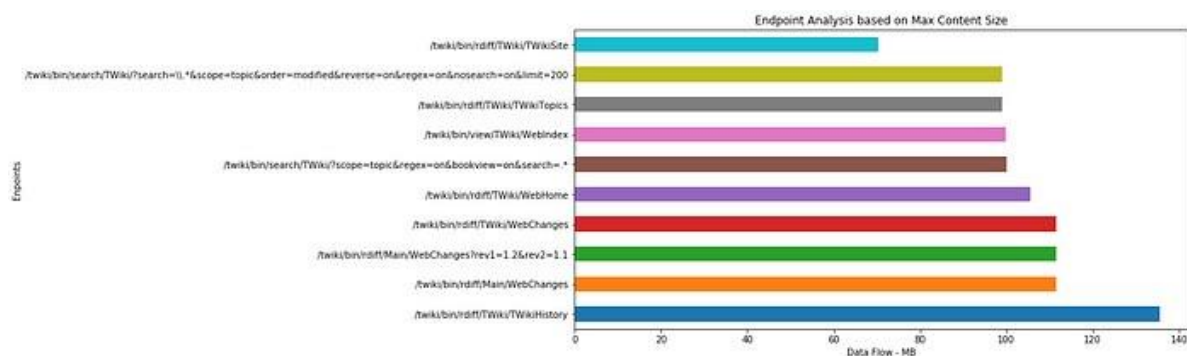
```

Step 3 : Analyze Top 10 Endpoints which Transfer Maximum Content in MB

```

1 #Top 10 Endpoints which Transfer Maximum Content
2 #.rdd.map() - Will convert the resulted rows from SQL query into a map
3 # .collect() - actually executes the DAG to get the overall results
4 topEndpointsMaxSize = (sqlContext
5                       .sql("SELECT endpoint,content_size/1024 FROM logs ORDER BY content_size DESC LIMIT 10")
6                       .rdd.map(lambda row: (row[0], row[1]))
7                       .collect())
8 # Plot Analysis Code
9 bar_plot_list_of_tuples_horizontal(topEndpointsMaxSize,'Data Flow - MB','Endpoints','Endpoint Analysis
   based on Max Content Size')

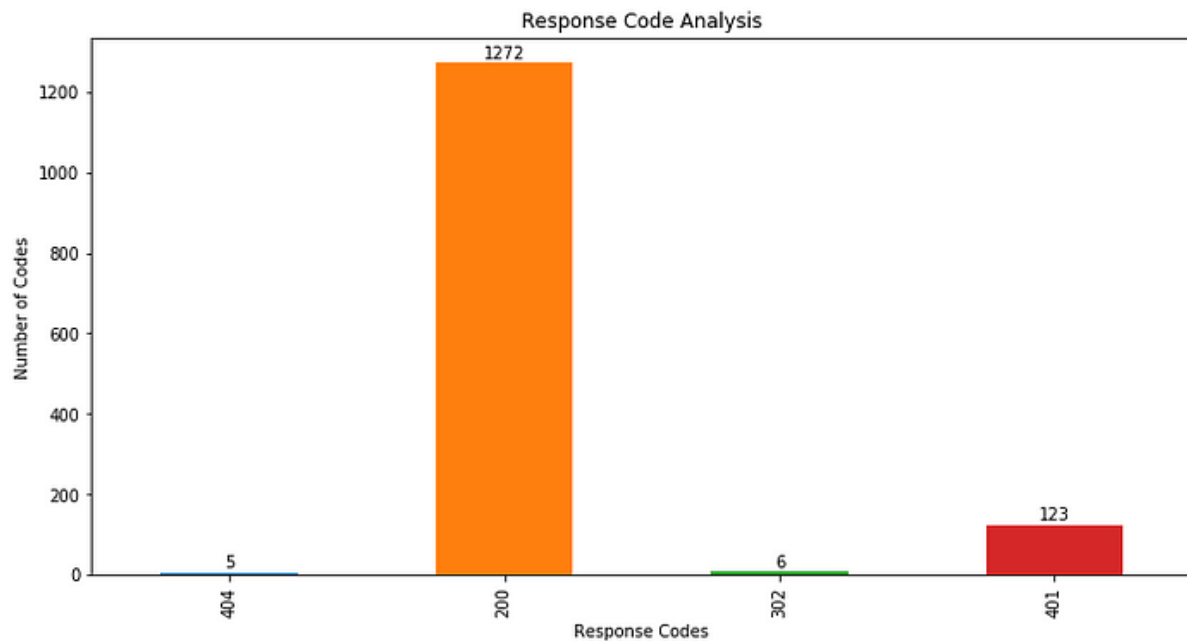
```





Academic Year: 2022-2023

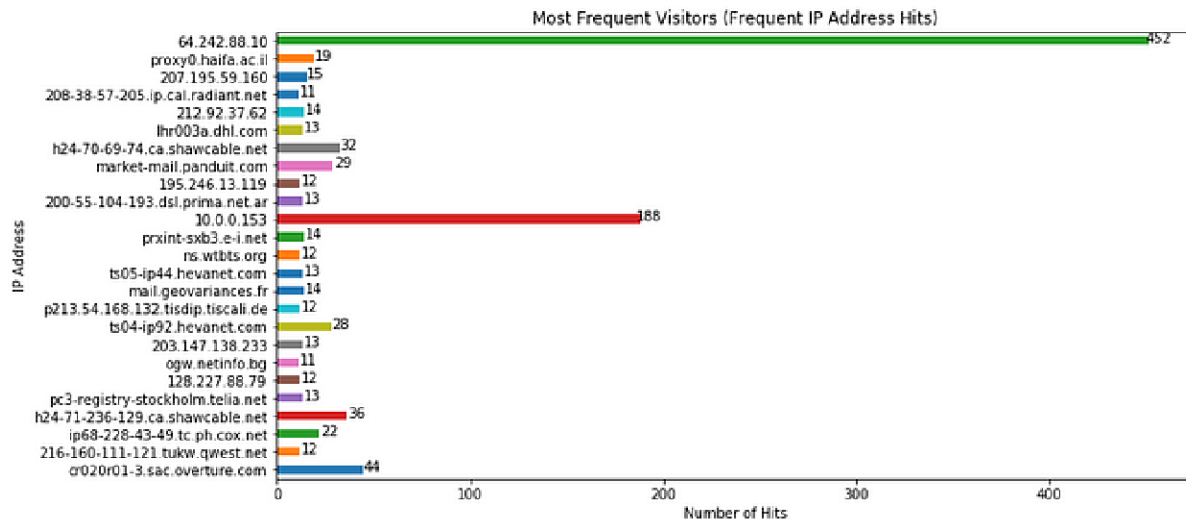
```
1 # Response Code Analysis
2 responseCodeToCount = (sqlContext
3     .sql("SELECT response_code, COUNT(*) AS theCount FROM logs GROUP BY
4         response_code")
5     .rdd.map(lambda row: (row[0], row[1]))
6     .collect())
7
8 # Code to Plot the results
9 def bar_plot_list_of_tuples(input_list,x_label,y_label,plot_title):
10     x_labels = [val[0] for val in input_list]
11     y_labels = [val[1] for val in input_list]
12     plt.figure(figsize=(12, 6))
13     plt.xlabel(x_label)
14     plt.ylabel(y_label)
15     plt.title(plot_title)
16     ax = pd.Series(y_labels).plot(kind='bar')
17     ax.set_xticklabels(x_labels)
18     rects = ax.patches
19     for rect, label in zip(rects, y_labels):
20         height = rect.get_height()
21         ax.text(rect.get_x() + rect.get_width()/2, height + 5, label, ha='center', va='bottom')
```



```
1 # Most Frequent Visitors (Most Frequent IP Address visits).
2 frequentIpAddressesHits = (sqlContext
3     .sql("SELECT ip_address, COUNT(*) AS total FROM logs GROUP BY ip_address HAVING total >
4         10")
5     .rdd.map(lambda row: (row[0], row[1]))
6     .collect())
7
8 bar_plot_list_of_tuples_horizontal(frequentIpAddressesHits,'Number of Hits','IP Address','Most Frequent
9     Visitors (Frequent IP Address Hits)')
```



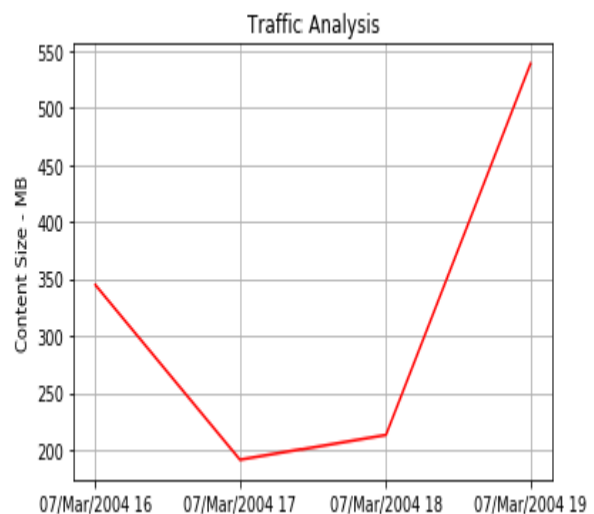
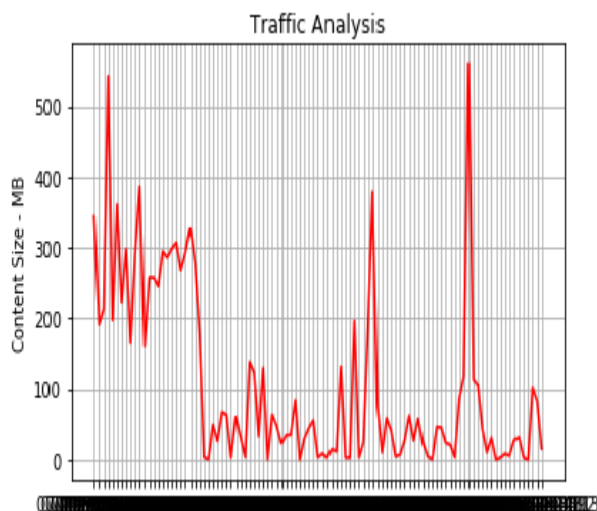
Academic Year: 2022-2023



```

1 # Traffic Analysis for past One Week
2 trafficWithTime = (sqlContext
3                     .sql("SELECT date, content_size/1024 FROM logs")
4                     .rdd.map(lambda row: (row[0], row[1]))
5                     .collect())
6 time_series_plot(trafficWithTime)
7

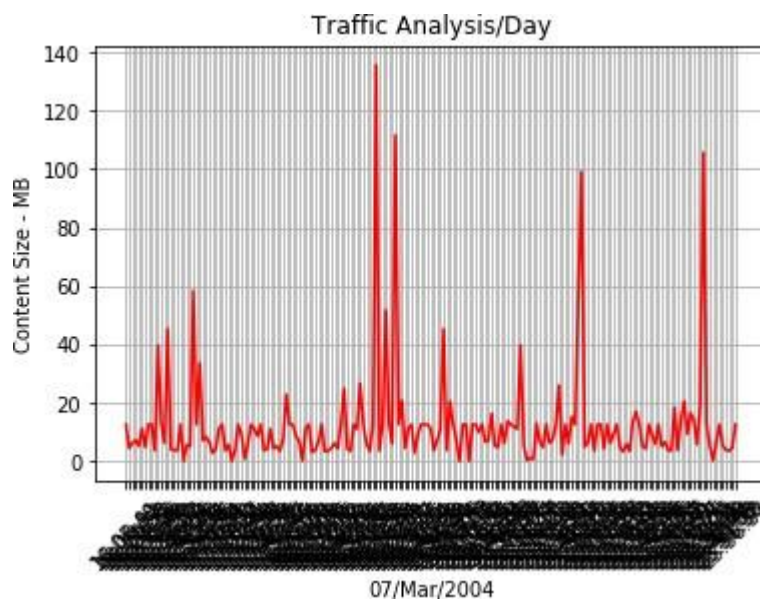
```



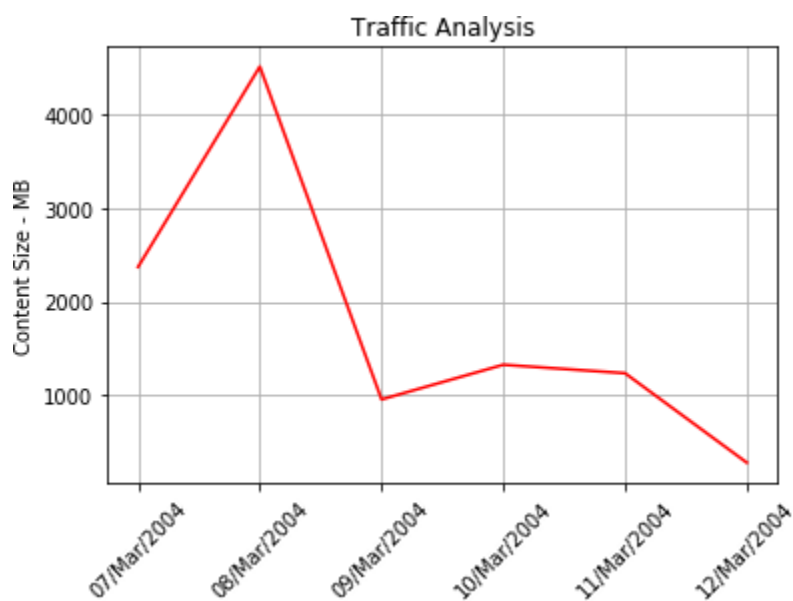


Academic Year: 2022-2023

```
1 # Overall Traffic Analysis for a Day
2 Day = '07/Mar/2004'
3 trafficperDay = (sqlContext
4                   .sql("SELECT time,content_size/1024 FROM logs where date='07/Mar/2004'")
5                   .rdd.map(lambda row: (row[0], row[1]))
6                   .collect())
7 time_series_plot(trafficperDay,Day,'Content Size - MB','Traffic Analysis/Day')
```



Outliers can be clearly detected by analysis the spikes and which end points were been hit at time by what IP Addresses.





Here, we can see an unusual spike on 8th March, which can be analyzed further for identifying discrepancy.

Code for Plot Analysis:



```
1 def time_series_plot(input_list,x_label,y_label,plot_title):
2     x_labels = [val[0] for val in input_list]
3     y_labels = [val[1] for val in input_list]
4     dict_plot = OrderedDict()
5     for x,y in zip(x_labels,y_labels):
6         # cur_val = x.split(":", 1)[0]
7         cur_val = x.split(" ")[0]
8         #print(cur_val)
9         dict_plot[cur_val] = dict_plot.get(cur_val, 0) + y
10    input_list = list(dict_plot.items())
11    x_labels = [val[0] for val in input_list]
12    y_labels = [val[1] for val in input_list]
13    plt.plot_date(x=x_labels, y=y_labels, fmt="r-")
14    plt.xticks(rotation=45)
15    plt.title(plot_title)
16    plt.xlabel(x_label)
17    plt.ylabel(y_label)
18    plt.grid(True)
19    plt.show()
```




```
20
21 def bar_plot_list_of_tuples_horizontal(input_list,x_label,y_label,plot_title):
22     y_labels = [val[0] for val in input_list]
23     x_labels = [val[1] for val in input_list]
24     plt.figure(figsize=(12, 6))
25     plt.xlabel(x_label)
26     plt.ylabel(y_label)
27     plt.title(plot_title)
28     ax = pd.Series(x_labels).plot(kind='barh')
29     ax.set_yticklabels(y_labels)
30 for i, v in enumerate(x_labels):
31     ax.text(int(v) + 0.5, i - 0.25, str(v),ha='center', va='bottom')
```

```
32
33 # Frequent End Points
34 topEndpoints = (sqlContext
35                 .sql("SELECT endpoint, COUNT(*) AS total FROM logs GROUP BY endpoint ORDER BY total
36                       DESC LIMIT 10")
37                 .rdd.map(lambda row: (row[0], row[1]))
38                 .collect())
39 bar_plot_list_of_tuples_horizontal(topEndpoints,'Number of Times Accessed','End Points','Most
40 Frequent Endpoints')
```

