# Principles Of Bigdata Management

**Phase-2 Report**

**By**

**Team-23**

Hemanth Kumar Reddy Dantu(16233525)

Mahesh Chowdary Jamallamudi(16234558)

Theyab Alharbi(16194252)

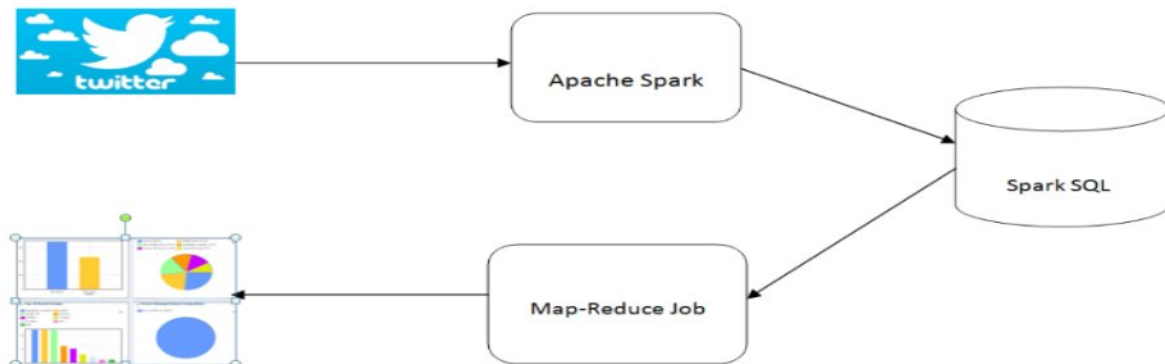Arun Kumar Anthati(16232818)

**Project Description:**

For the phase-2 we have collected to more than 1 GB tweets on Banking systems. We have used hashtags for banks and collected the tweets on Banking systems.

After the tweets collections, we have analyzed the banking tweets and visualized using "d3.js" and "high charts". Total we have implemented 6 queries, 3 queries using "Java RDD" and 3 queries using "DataFrames".

We have implemented the dynamic web project it first process the data based on the query and the result stored on .csv file and immediately it calls html page to visualize the collected results.

**Architecture Diagram:**

**Software Technologies and Tools:**
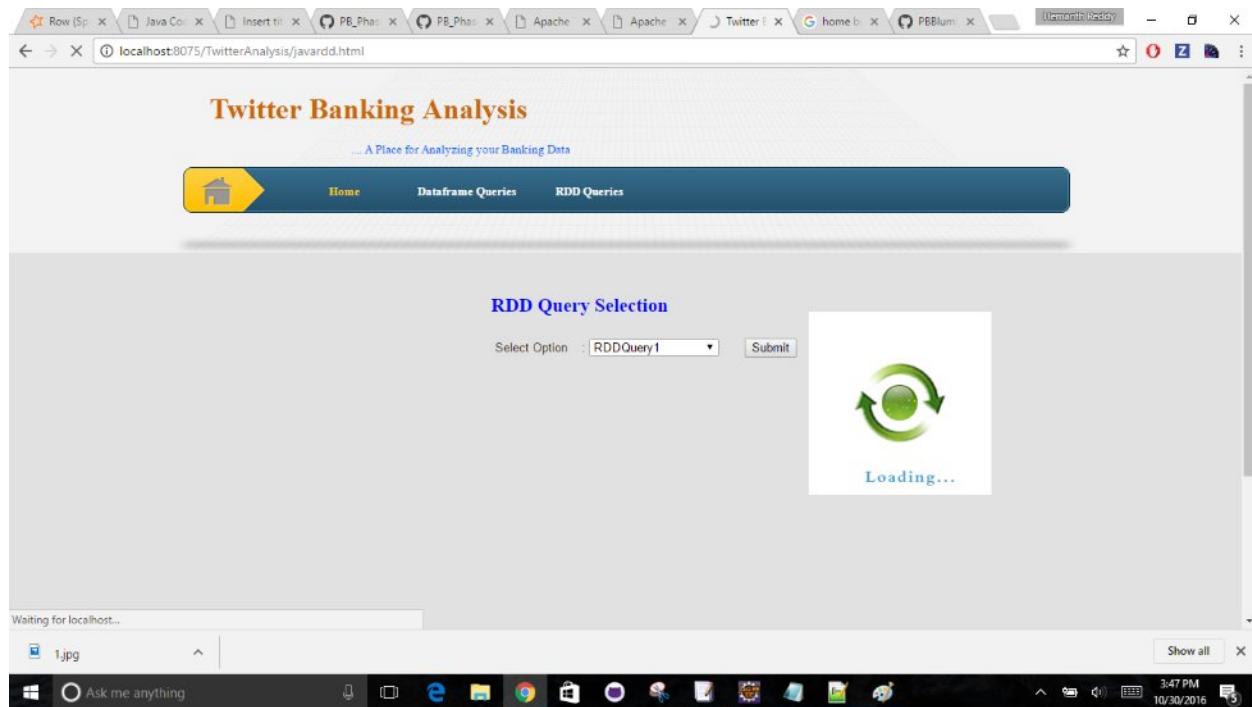
**User Interface:** HTML,CSS. D3.js and High Charts.

**Tools:** Eclipse

**Environment:** Apache Spark

**No SQL Database:** Spark SQL.

**Backend:** Java Spark and Servlets.

# RDD Queries:
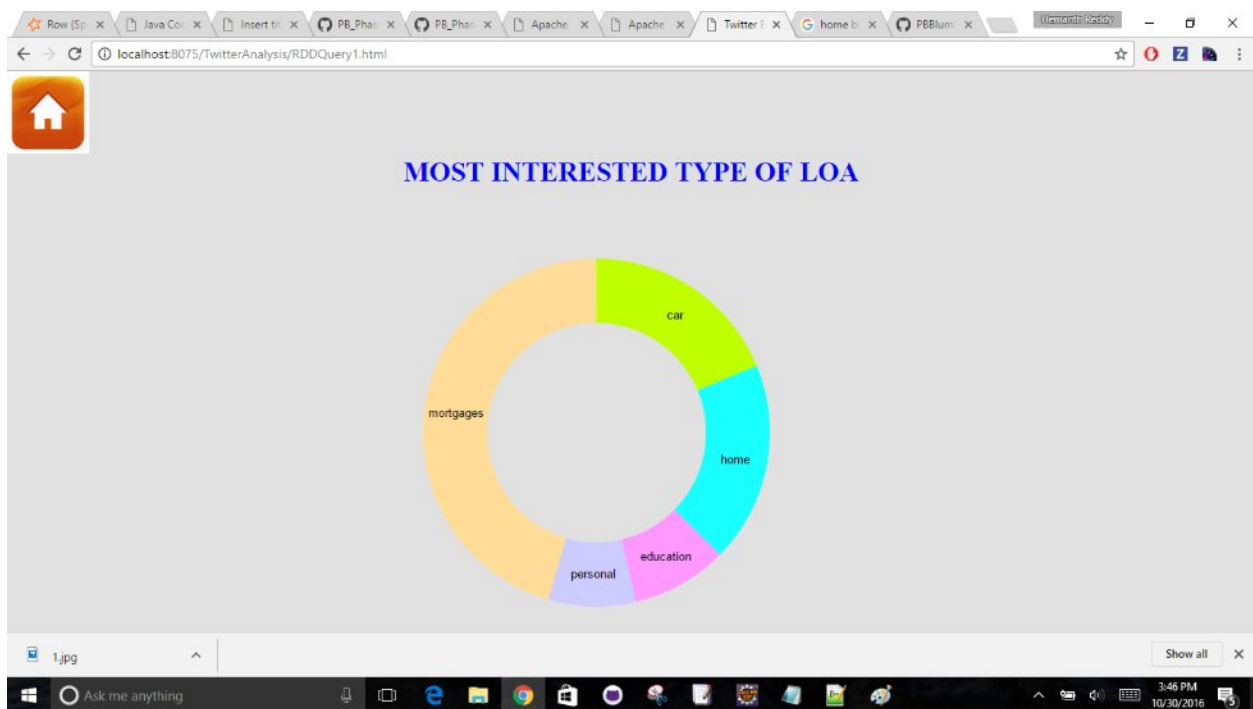
**Home Page Visualization:**

**RDD Query-1:** The most interested loan types that the users discussing on twitter.

**.CSV file Results:**

| LoanType | Count |
|----------|-------|
| car | 23 |
| home | 23 |
| education | 11 |
| personal | 10 |
| mortgages | 56 |

**Visualization Results:**

**RDD Query-2:** This query analyzes people's opinion on bank security. We have analyzed the positive(Secure) and negative (insecure) tweets.

**.CSV file Results:**

| Words | Count | | |
|-------|-------|---|---|
| Secure | 27335 | | |
| Insecure | -15739 | | |
| | | | |

**Visualization Results:**

**RDD Query-3:** This query analyses the top banks on service wise.

**.CSV file Results:**

| A | B | C | D | E | F |
|---|---|---|---|---|---|
| Name | Count | | | | |
| Lloyds Bank | 164350 | | | | |
| Guaranty Trust | 128312 | | | | |
| HDFC Bank | 92439 | | | | |
| TD Bank | 63593 | | | | |
| Next Bank | 43652 | | | | |
| | | | | | |

**Visualization Results:**

# DataFrame Queries:

**DataFrame Query-1:** The most popular banks on twitter. We have analyzed the most followers of the each banks and visualized the top 6 banks.

## .CSV file Results:

| Name | Count |
|------|-------|
| World Bank | 1721682 |
| Guaranty Trust Bank | 527570 |
| Jodrell Bank | 230042 |
| Kotak Mahindra Bank | 113469 |
| Diamond Bank | 87193 |
| ICICI Bank | 84555 |
| | |

## Visualization Results:

**DataFrame Query-2:** The most tweeted timings on banks. We have analyzed the top 8 twitting times on banks.

**.CSV file Results:**

| Time | Count |
|------|-------|
| Sat Oct 26 20:57:29 +0000 2016 | 51 |
| Thu Oct 17 20:40:00 +0000 2016 | 50 |
| Wed Oct 16 19:58:00 +0000 2016 | 47 |
| Sat Oct 26 20:31:51 +0000 2016 | 46 |
| Mon Oct 14 20:10:50 +0000 2016 | 46 |
| Wed Oct 09 19:43:52 +0000 2016 | 44 |
| Sat Oct 26 20:31:45 +0000 2016 | 43 |
| Wed Oct 16 17:52:30 +0000 2016 | 43 |

**Visualization Results:**

**DataFrame Query-3:** This query analyzes the Public and Private and other banks tweets percentage.

**.CSV file Results:**

| TweetStatus | Percentage |
|---|---|
| CommercialBanks% | 10 |
| InvestmentBanks% | 6 |
| others% | 82 |

**Visualization Results:**

# Log files:

**Source Code:**

Source code is available on GitHub. Please find the URL below.

https://github.com/HemanthDantu/PB_Phase-2_Project/tree/master/Source%20Code

**References:**

https://github.com/nivdul/spark-in-practice/blob/master/src/main/java/com/handson/spark/dataframe/DataFrameOnTweets.java

https://github.com/AgilData/spark-rdd-dataframe-dataset/blob/master/src/main/java/example_java/dataframe/JavaDataFrameExample.java

https://spark.apache.org/docs/1.6.0/sql-programming-guide.html

http://www.programcreek.com/java-api-examples/index.php?api=org.apache.spark.sql.Row

https://github.com/stefani75/workspace/blob/b90a63f2f3028fb358b28a77ac416b223d37a52a/projet1/Hands-On-Spark-java-solution/src/main/java/com/duchessfr/spark/part3/sparksql/FunWithSparkSQL.java

https://github.com/AshokYaganti/PB_Phase2_TwitterAnalysis