**COLLECTING TWEETS USING TWITTER STREAMING API'S**

**Sub: Principles of Big Data Management**

**Phase 2 of Project**

Sirisha Rella          (16268829)

Vineeth Dudipalli    (16274444)

Rajeshwari Cholleti (16272108)

Instructor

Dr. PRAVEEN RAO, Ph.D.

## Abstract:

Phase 2 of this project deals with the following requirements:

1. Writing interesting analytical queries on twitter data that we collected.

2. Developing interesting visualizations.

## Title:

Technology in different Domains.

## Technologies & Tools Used:

1. Spark

2. Python

3. Tableau

4. Spyder
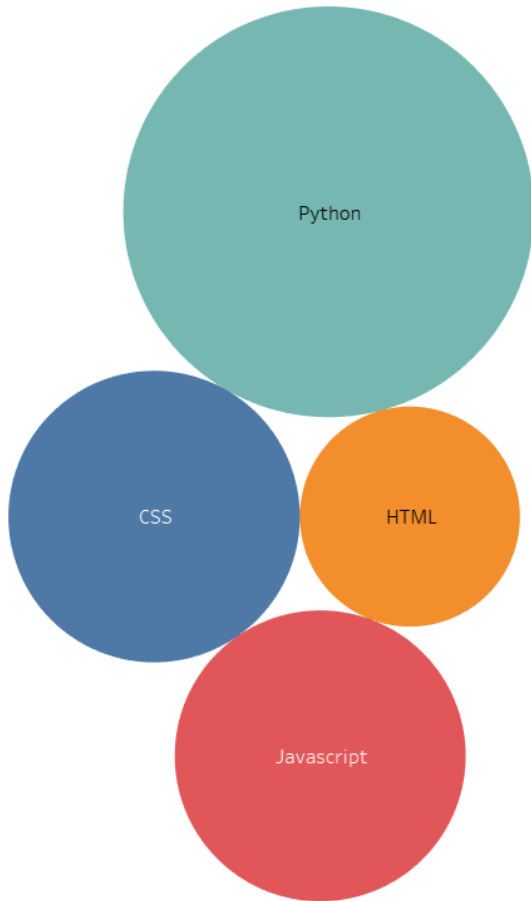
## Queries and Analysis:

## Query-1:

In this query we found the number of users tweeting about Python, JavaScript, HTML and CSS.

## Code:

```
sqlDF = spark.sql("SELECT COUNT(*) AS NumberOfTweets, 'Python' as Language FROM Technology where text LIKE '%Python%'\
    UNION\
    SELECT COUNT(*) AS NumberOfTweets,'Javascript' as Language FROM Technology where text LIKE '%javascript%'\
    UNION\
    SELECT COUNT(*) AS NumberOfTweets, 'HTML' as Language FROM Technology where text LIKE '%HTML%'\
    UNION\
    SELECT COUNT(*) AS NumberOfTweets, 'CSS' as Language FROM Technology where text LIKE '%CSS%'\
    ")
pd = sqlDF.toPandas()
```

**Visualization:**

**Query-2:**

In this query we found how Microsoft, IBM and Cerner associates/employees are actively tweeting about AI.

**Code:**

```
df.createOrReplaceTempView("Technology")
sqlDF = spark.sql("SELECT 'Microsoft' as Company, count(*) as Count from Technology where text like '%Microsoft%' and text like '%AI%'\
    UNION\
    SELECT 'IBM' as Company, count(*) as Count from Technology where text like '%IBM%' and text like '%AI%'\
    UNION\
    SELECT 'Cerner' as Company, count(*) as Count from Technology where text like '%Cerner%' and text like '%AI%'")
pd = sqlDF.toPandas()
```
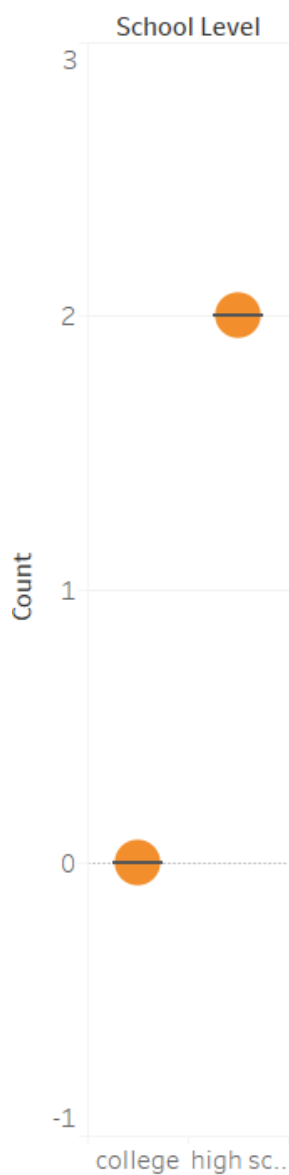
**Visualization:**

**Query-3:**

In this query we found count of active high school/college students'
achievements/participation in technology.

**Code:**

```
sqlDF = spark.sql("SELECT 'high school' as SchoolLevel, count(*) as Count from Technology where text like '%high school%' and text like '%technology%' \
    UNION\
    SELECT 'college' as SchoolLevel, count(*) as Count from Technology where text like '%college%' and text like '%technology'")
pd = sqlDF.toPandas()
```
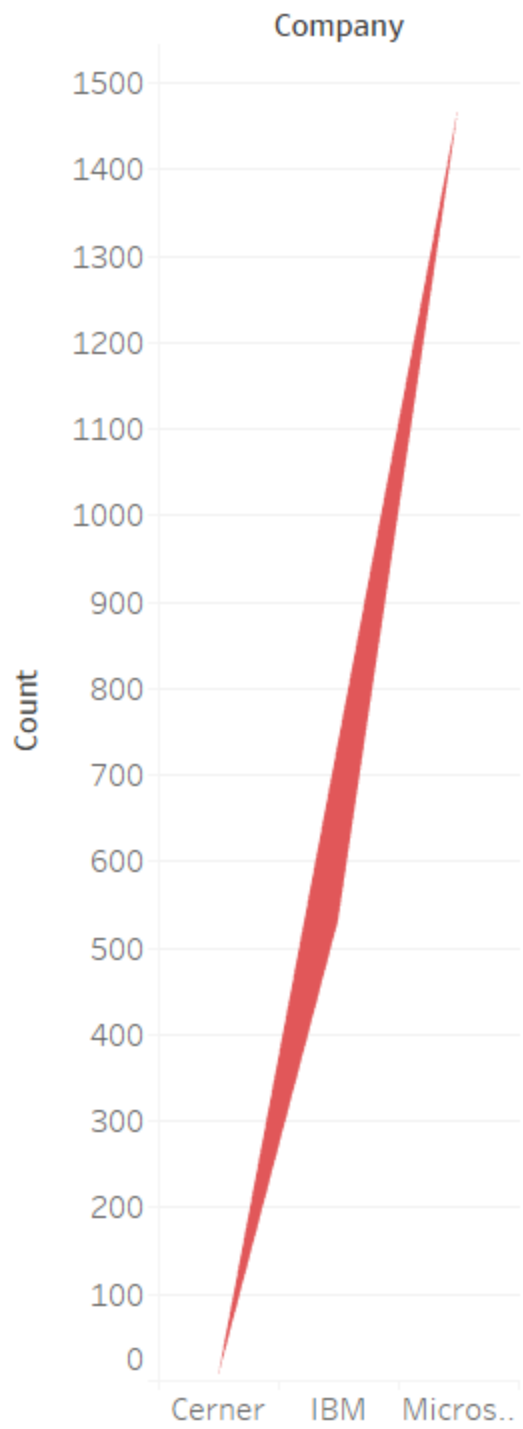
**Visualization:**

**Query-4:**

In this query we found the Count of people who are actively tweeting about Microsoft, Cerner and IBM

**Code:**

```
df.createOrReplaceTempView("Technology")
sqlDF = spark.sql("SELECT 'Microsoft' as Company, count(*) as Count from Technology where text like'%Microsoft%'\
    UNION\
    SELECT 'IBM' as Company, count(*) as Count from Technology where text like '%IBM%'\
    UNION\
    SELECT 'Cerner' as Company, count(*) as Count from Technology where text like '%Cerner%'")
pd = sqlDF.toPandas()
```
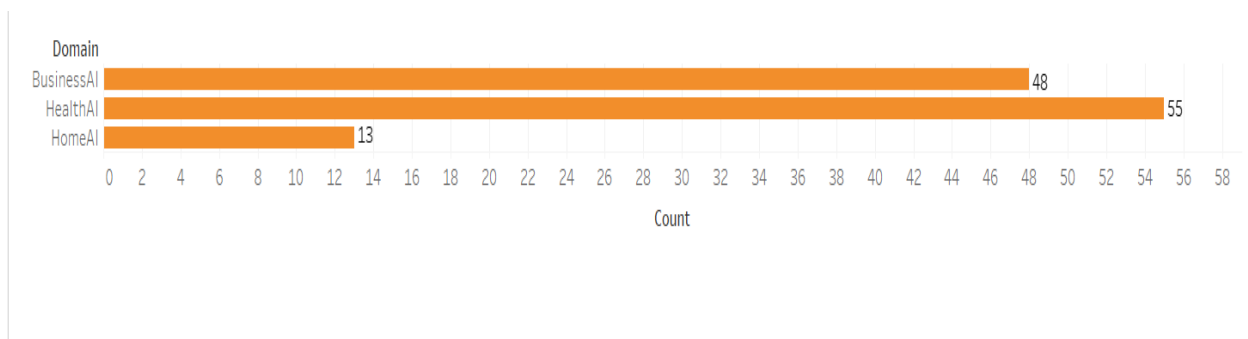
**Visualization:**

Company

Count

| Cerner | IBM | Micros.. |

**Query-5:**

In this query we found tweets about how AI is impacting in health, home and business AI.

**Code:**

```
sqlDF = spark.sql("SELECT 'HomeAI' as Technology, count(*) as Count from Technology where text like '%home%' and text like '%AI%'\
    UNION\
    SELECT 'HealthAI' as Technology, count(*) as Count from Technology where text like '%health%' and text like '%AI%'\
    UNION\
    SELECT 'BusinessAI' as Technology, count(*) as Count from Technology where text like '%business%' and text like '%AI%'")
pd = sqlDF.toPandas()
pd.plot(kind="bar", x="Technology", y="Count")
```
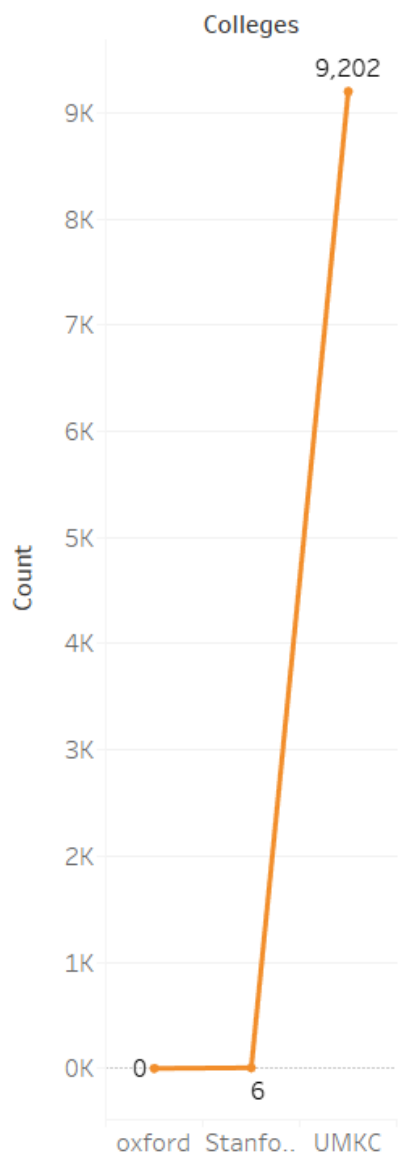
**Visualization:**

**Query-6:**

In this query we found which students who are actively participating/tweeting about python programming based on college name.

**Code:**

```
sqlDF = spark.sql("SELECT 'Stanford' as Colleges, count(*) as Count from Colleges where text like '%Stanford%' and text like '%AI%'\
    UNION\
    SELECT 'oxford' as Colleges, count(*) as Count from Colleges where text like '%oxford%' and text like '%AI%'\
    UNION\
    SELECT 'UMKC' as Colleges, count(*) as Count from Colleges where text like '%UMKC%' or text like '%AI%'")
plt.show()
```
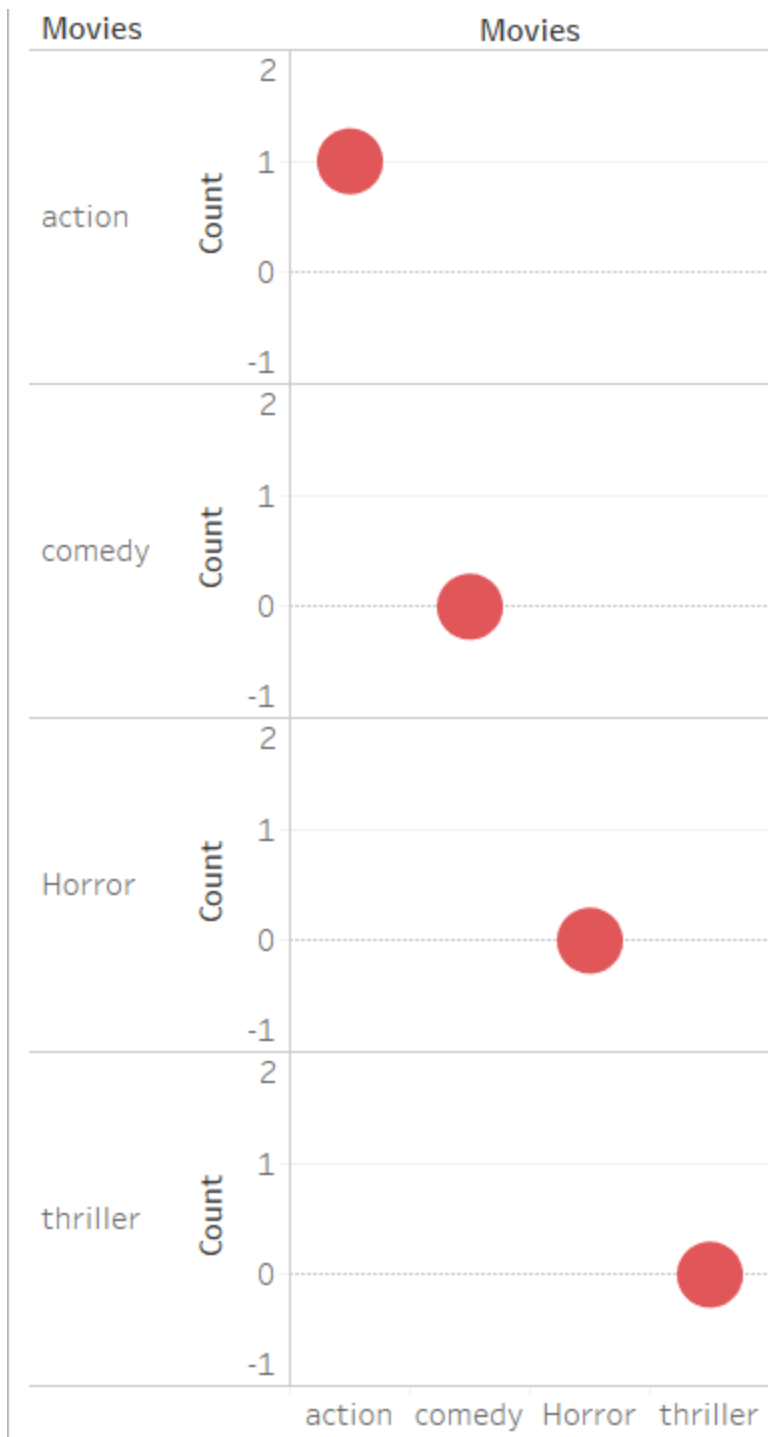
**Visualization:**

**Query-7:**

In this query we found tweets that talk about different kind of films that are using vfx effects.

**Code:**

```
df.createOrReplaceTempView("Movies")
sqlDF = spark.sql("SELECT 'Horror' as Movies, count(*) as Count from Movies where text like '%horror%' and text like '%vfx%'\
    UNION\
    SELECT 'comedy' as Movies, count(*) as Count from Movies where text like '%comedy%' and text like '%vfx%'\
    UNION\
    SELECT 'thriller' as Movies, count(*) as count from Movies where text like '%thriller' and text like '%vfx%'\
    UNION\
    SELECT 'action' as Movies, count(*) as Count from Movies where text like '%action%' and text like '%vfx%'")
```

**Visualization:**

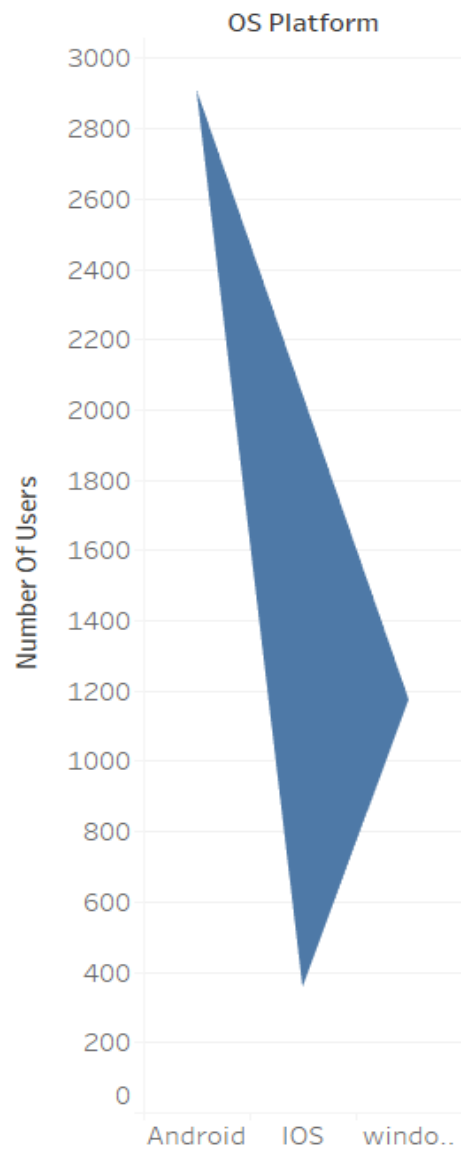| Movies | Movies |
|--------|--------|

## Query-8:

In this query we found number of people tweeting about IOS, Android and Windows operating systems.

**Code:**

```
sqlDF = spark.sql("SELECT COUNT(*) AS NumberOfUsers, 'IOS' as OSPlatform FROM Technology where text LIKE '%IOS%'\
        UNION\
        SELECT COUNT(*) AS NumberOfUsers,'Android' as OSPlatform FROM Technology where text LIKE '%Android%'\
        UNION\
        SELECT COUNT(*) AS NumberOfUsers, 'windows' as OSPlatform FROM Technology where text LIKE '%windows%'")
sqlDF.show()
pd = sqlDF.toPandas()
pd.to_csv('eight.csv', index=False)
pd.plot(kind="bar",x="OSPlatform",y="NumberOfUsers")
plt.show()
```

**Visualization:**

# Sheet 1

**OS Platform**



Number Of Users

Android    IOS    windo..

windows

## Query-9:

In this query we found number of people tweeting about WhatsApp, Instagram and snapchat

## Code:

```
df.createOrReplaceTempView("Technology")
sqlDF = spark.sql("SELECT COUNT(*) AS NumberOfUsers, 'Whatsapp' as OSPlatform FROM Technology where text LIKE '%whatsapp%'
        UNION\
        SELECT COUNT(*) AS NumberOfUsers,'Instagram' as OSPlatform FROM Technology where text LIKE '%instgram%'\
        UNION\
        SELECT COUNT(*) AS NumberOfUsers, 'Snapchat' as OSPlatform FROM Technology where text LIKE '%snapchat%'")
sqlDF.show()
pd = sqlDF.toPandas()
print(pd)
pd.to_csv('nine.csv', index=False)
pd.plot.pie(y='NumberOfUsers',labels=['whatsapp','instagram','snapchat'],figsize=(3,3))
plt.show()
```

## Visualization:

## Social Site

**Query-10:**

In this query we found number of people tweeting about SQL and NoSQL

**Code:**

```
sqlDF = spark.sql("SELECT COUNT(*) AS NumberOfTweets, 'SQL' as Database FROM Technology where text LIKE '%SQL%'\
        UNION\
        SELECT COUNT(*) AS NumberOfTweets,'NoSQL' as Database FROM Technology where text LIKE '%NoSQL%'")
pd = sqlDF.toPandas()
pd.to_csv('ten.csv', index=False)
pd.plot(kind="bar",x="Database",y="NumberOfTweets")
```

**Visualization:**

# Testing:

## Manual testing:

We tried to test the results manually by taking the keywords used for the query and finding the tweets for the same keywords using twitter search engine in google and tried to match the tweets with the twitter data we retrieved.

For instance, consider the VFX query for which we retrieved the profile user name of the particular tweet and searched it manually in twitter to find the tweet and check it whether he tweeted or not.

**Testing Based on tools:**

**Running in SPYDER:**



**Running in LINUX Terminal:**



As we can observe from the above two outputs both the outputs match each other therefore we confirm that the result is correct.

# Code Link:

https://github.com/VineethDvv/Principles-of-Bigdata-