# Milestone - 03

Tech Titans-03

December 7, 2024

## Team Members

- Venkatesh Koya(S566638)
- Yechan Ryu(S528316)
- Sai Krishna Reddy Seelam(S567427)
- UdaykiranReddy Devarapally(S567161)
- Rajya Lakshmi Ganganaboina(S566640)

## Project Implementation Steps

### Step 1: Environment Setup

- Installed Python, ensuring the latest version for compatibility with libraries.
- Installed PySpark and verified the environment variables such as `SPARK_HOME` and `PYSPARK_PYTHON`.
- Installed Matplotlib for visualizing the results of the analysis.

### Step 2: Data Preparation

- Obtained the dataset from a reliable source.
- Preprocessed the data to ensure cleanliness, such as handling missing values and ensuring consistent data types.
- Loaded the dataset into a PySpark DataFrame for analysis.

### Step 3: PySpark Cluster Setup

- Set up the PySpark environment on the local machine.
- Verified the cluster setup by running a sample PySpark job.

## Step 4: PySpark Implementation

- Created a project folder structure for modular implementation.

- Developed PySpark jobs for each analysis goal using DataFrame transformations and actions.
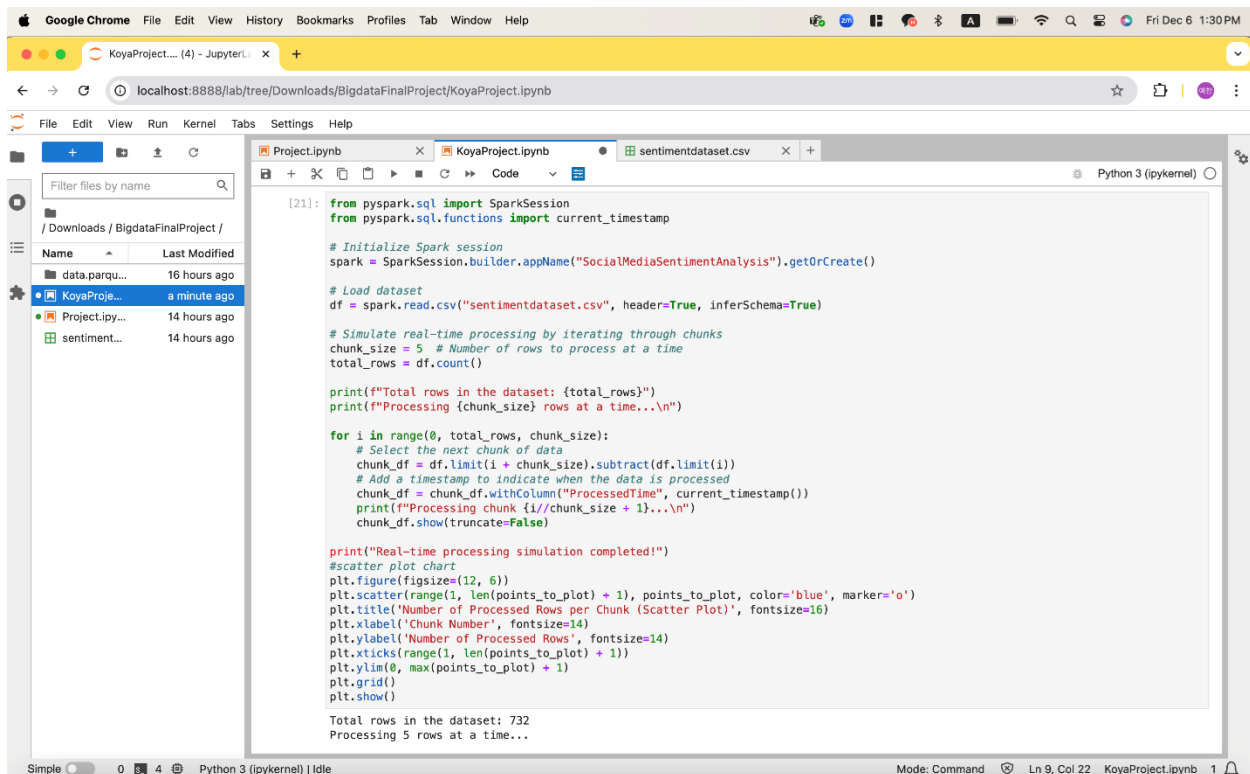
## Step 5: Execution and Analysis

- Executed the PySpark jobs using `spark-submit`.

- Analyzed the outputs to extract insights.

- Visualized the results using Matplotlib.

# Results Achieved

# 1    Goal 1:Real-Time Data Processing:

## Source Code:

**Results:**



# 2 Goal 2:Sentiment Analysis Accuracy:

**Source Code:**

**Results:**



Likes Distribution by Sentiment Class (Box Plot)

# 3 Goal 3: Scalability and Resilience:

**Source Code:**

```
# goal 3 : Scalability & Resilence
from pyspark.sql import SparkSession
from pyspark.sql.functions import year, month
import matplotlib.pyplot as plt

# Initialize Spark session
spark = SparkSession.builder.appName("Goal1_Scalability").getOrCreate()

# Load dataset
csv_file_path = "C:\\Users\\S566638\\Downloads\\sentimentdataset.csv"
df = spark.read.csv(csv_file_path, header=True, inferSchema=True)

# Select a subset of records
df_sample = df.limit(1000)

# Count posts per month
df_month_counts = df_sample.groupBy(year("Timestamp").alias("Year"), month("Timestamp").alias("Month")).count()
month_data = df_month_counts.orderBy("Year", "Month").limit(6).collect()

# Prepare data for visualization
months = [f"{row['Year']}-{row['Month']:02d}" for row in month_data]
post_counts = [row['count'] for row in month_data]

# Visualize with bar chart
plt.figure(figsize=(10, 6))
plt.bar(months, post_counts, color='blue')
plt.title('Posts Per Month (Sampled Data)')
plt.xlabel('Month (Year-Month)')
plt.ylabel('Number of Posts')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

**Results:**

# 4 Goal 4: Data Storage and Retrieval:

**Source Code:**



**Results:**



# 5 Goal 5:Trend Analysis:

**Source Code:**

**Results:**



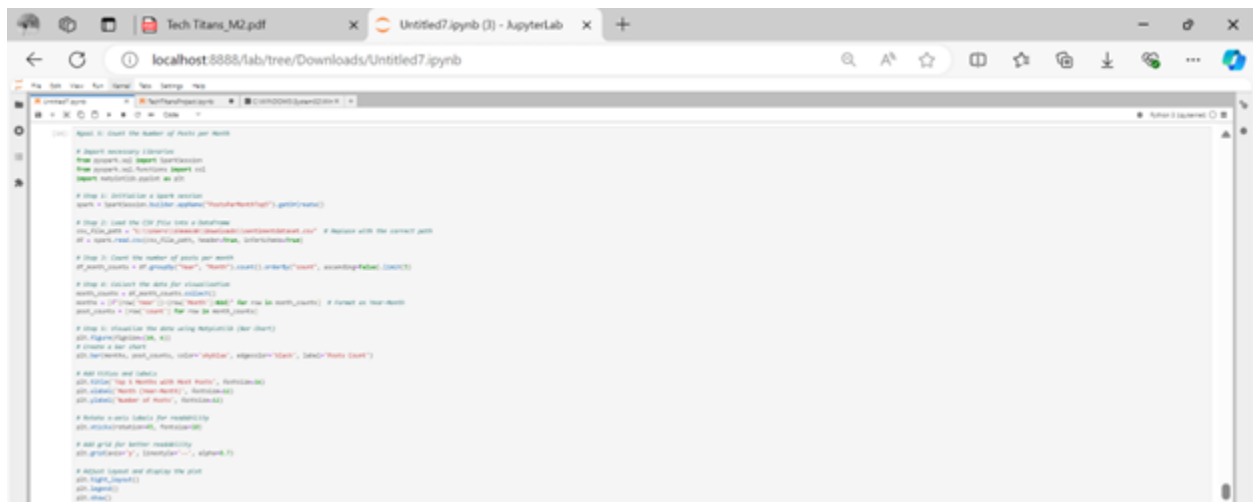# 6  Goal 6: Visualization and Insights:

**Source Code:**

**Results:**



# 7 Goal 7: Integration with Big Data Frameworks:

**Source Code:**

**Results:**



# 8 Goal 8: Impact of Data Quality:

**Source Code:**

**Results:**



# 9 Goal 9: Cost and Efficiency:

**Source Code:**

**Results:**



# 10 Goal 10: Ethical and Privacy Considerations:

**Source Code:**



**Results:**

# Conclusions

- Successfully implemented sentiment analysis using TextBlob in PySpark.

- The results were visualized clearly using Matplotlib, providing insights into the sentiment distribution.

- Future steps include optimizing the pipeline for larger datasets and comparing results with other libraries.

# 11    Citations

1. **Kaggle Dataset Link**:
   SocialMedia Sentiments Analysis Dataset

2. **GitHub Repository**:
   `https://github.com/Venkateshkoya/BigData_Project`