# Stock Price Prediction

COURSE: BIG DATA SYSTEM ENGINEERING USING SCALA (CSYE 7200)

PROFESSOR: ROBIN C. HILLYARD

TEAM: ASWATHNARAYAN KIRUBAKARAN, MEENAKSHI MUTHIAH

# Goals



▶ To determine the future value of a company stock or other financial instrument traded on an exchange

▶ Perform Time series analysis on Stock data using Scala and Spark

# Data Sources

- Huge stock market Dataset: https://www.kaggle.com/borismarjanovic/price-volume-data-for-all-us-stocks-etfs/version/2/home

- S&P 500 Stock Data: https://www.kaggle.com/camnugent/sandp500

- Yahoo Finance

- The dataset is taken is from Kaggle and has data for around 500 companies.

- Each data file has 8 columns and 2000 rows per file.

# Use cases

- ▶ Provide the user with future stock price for each company

- ▶ Provide the sentiment analysis for individual stock which the user wants to invest

# Methodology

data cleaning → Exploratory data analysis → Time series analysis → Twitter sentiment analysis using spark → Machine learning models

# Data Cleaning- Databricks

- Handle missing values

  removing nulls

  substituting nulls with the median
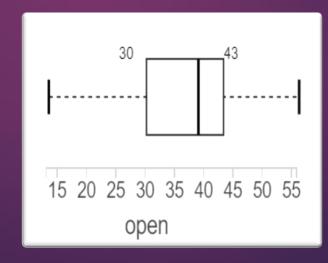
## Handling missing values

```
//count the missing values by summing the boolean output of the isNull() method with spark sql
import org.apache.spark.sql.functions.{sum, col}
df.select(df.columns.map(c => sum(col(c).isNull.cast("int")).alias(c)): _*).show


+----+----+----+---+-----+------+----+
|date|open|high|low|close|volume|Name|
+----+----+----+---+-----+------+----+
|   0|  11|   8|  8|    0|     0|   0|
+----+----+----+---+-----+------+----+

import org.apache.spark.sql.functions.{sum, col}
```
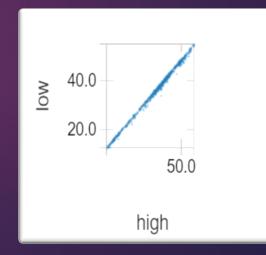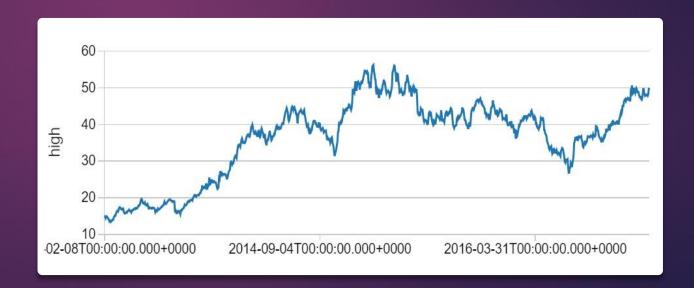
# EDA with Databricks

- Summary statistics
- Correlation
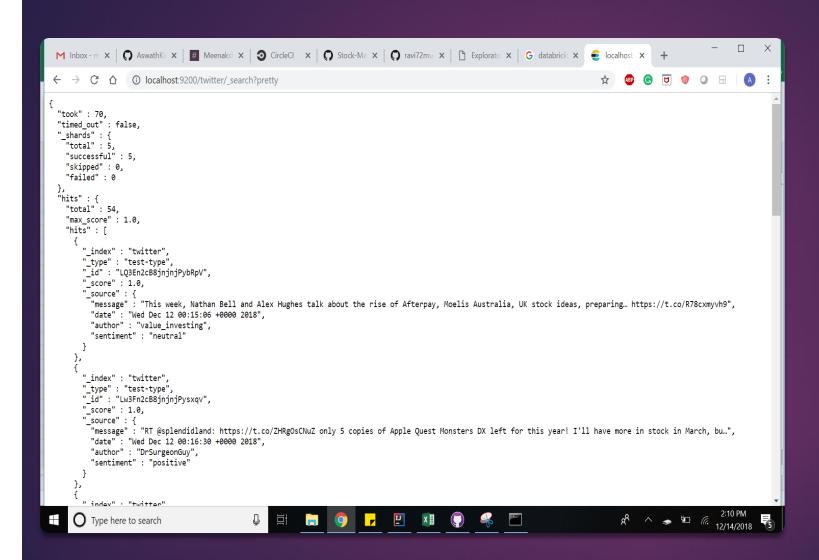- Trend analysis
- Outlier detection

# Times Series Analysis

- Stationary check
- Durbin-Watson Test for Auto Correlation
- Smoothing

# Twitter Sentiment Analysis

- ▶ Spark Twitter Streaming
- ▶ Sentiment Analysis using Stanford NLP
- ▶ Stored the results in ElasticSearch
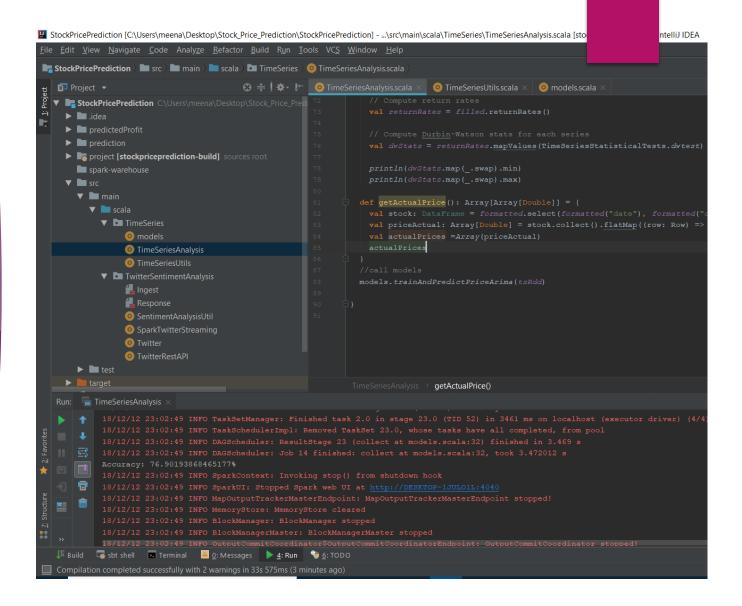- ▶ Visualization using Kibana

Elastic Search

# Machine Learning Models-Time Series Forecasting

- *Feature Engineering* : Converted the stock prices to vectors and applied smoothing

- *Models*: Utilized ARIMA model for stock price forecasting

- *Evaluation Metrics*: Accuracy and RMSE

# Acceptance Criteria

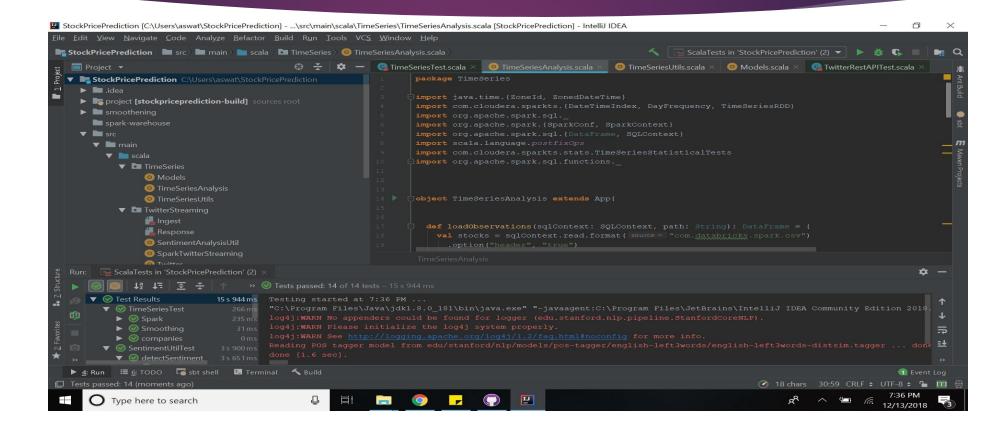▶ Stock price prediction accuracy should provide accuracy above 72%

# MileStones/Sprints

- ▶ Sprint 1: Perform data wrangling, integration and exploratory data analysis

- ▶ Sprint 2: Build machine learning models and perform training

- ▶ Sprint 3: Perform twitter streaming. Combine the work

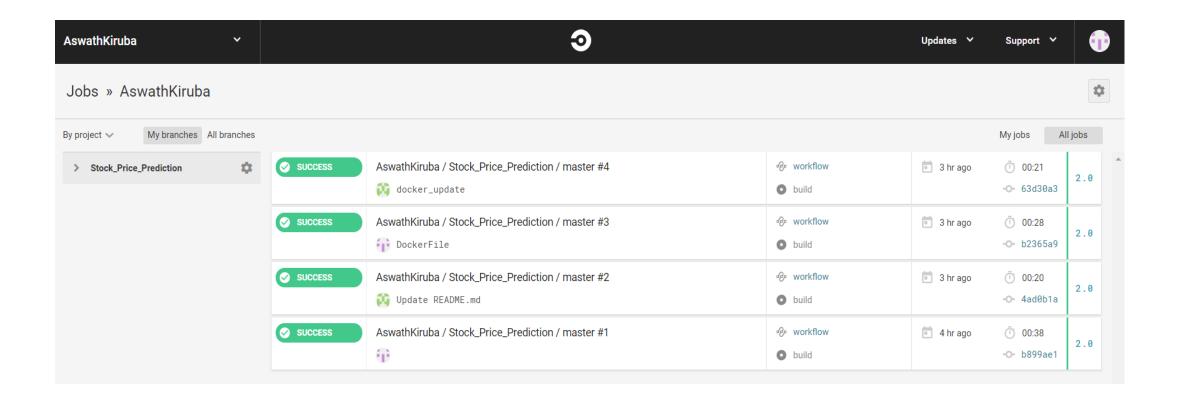- ▶ Sprint 4: Testing and Documentation

# Code

- ▶ Everything in Scala
  - ▶ Spark
  - ▶ Spark Mllib
  - ▶ Spark Streaming

# Unit Test

# Continuous Integration and Docker

# Thank you