# Final Project (24.3)

## Stock Predictor – Price & Media

# Agenda and Team Members

- Brian Young
- Carl Coffman
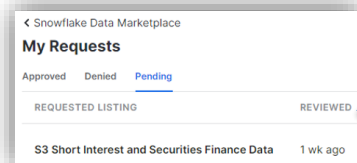- Push Palani
- Steven Crutcher

# Final Project - Stocks

## Project Scope

- Predict future stock price and how it may be influenced by social media noise – such as what had happened recently to GameStop and AMC.

## Site Hosts

- Stock Predictor (ccc-gh.github.io)
  https://ccc-gh.github.io/FinalProject-Stock/

- GitHub Code Repository
  https://github.com/CCC-GH/FinalProject-Stock

## Key Data Sources

Used:
- yfinance - # Yahoo! Finance market data
  $ pip install yfinance
  import yfinance

- pandas-datareader – Tweet data feed
  $pip install pandas-datareader
  import pandas_datareader

- Reddit – user-generated bulletins
  https://www.reddit.com/r/datasets/
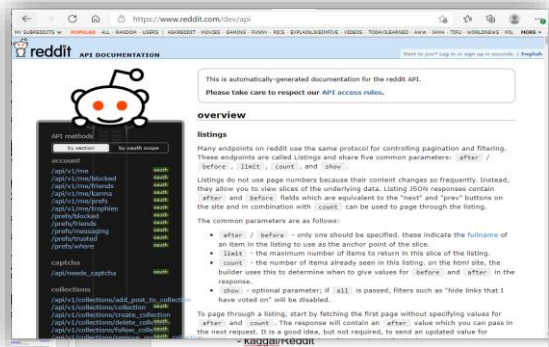
- Kaggle – WallStreetBest posts
  https://www.kaggle.com/datasets

Reviewed:
- snowflake – great data cloud connectivity store/retrieve, but no access given to stock-shorting data.



< Snowflake Data Marketplace
**My Requests**
Approved    Denied    **Pending**
REQUESTED LISTING                         REVIEWED
S3 Short Interest and Securities Finance Data    1 wk ago
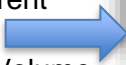
# Social Media and Cloud Data

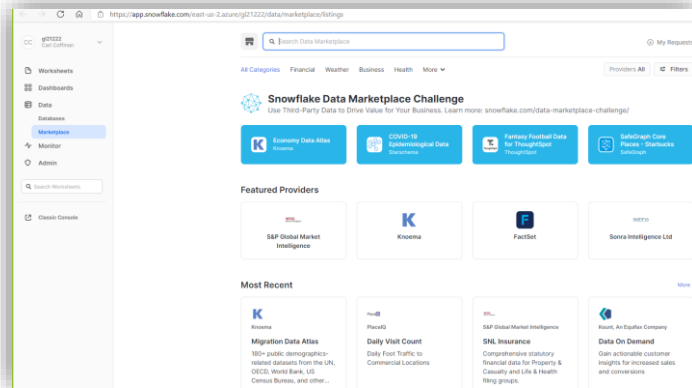Reddit – https://www.reddit.com/r/datasets/



## Data Challenges

- No significant, historical warehouse for social media data.

- Finding "free" historical data source for a particular stock's short activity – current (.info) and basic-historical is made available; Open, High, Low, Close, Volume



snowflake – Cloud Data/warehouse & Market Data



Kaggle –
WallStreetBest posts



Northwestern

# Libraries Used

**from sklearn:**
- linear_model import LinearRegression
- tree import DecisionTreeRegressor
- model_selection import train_test_split
- metrics import mean_squared_error

**\*ARIMA** - A popular and widely used statistical method for **time series forecasting**.
- from statsmodels.tsa.arima_model import ARIMA

*Other:*
- import pandas as pd
- import numpy as np
- import matplotlib.pyplot as plt
- from pandas.plotting import lag_plot
- from pandas.tseries.holiday import USFederalHolidayCalendar
- from pandas.tseries.offsets import CustomBusinessDay
- from pandas_datareader.data import DataReader
- import yfinance as yf
- import datetime

**Linear Regression - Time Series Forecasting**

**\*ARIMA -** AutoRegressive Integrated Moving Average:

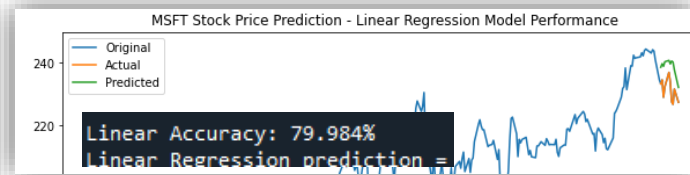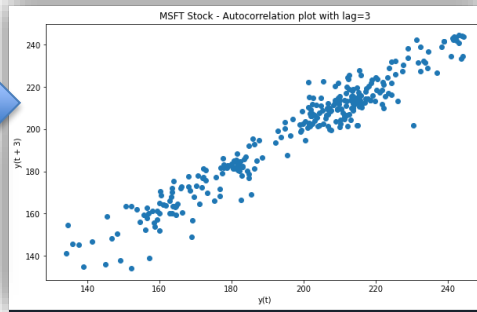- **AR:<Auto Regressive>** uses the dependent relationship between an observation and some predefined number of lagged observations (also known as "time lag" or "lag")

- **I:< Integrated >** model employs differencing of raw observations (e.g. it subtracts an observation from an observation at the previous time step) in order to make the time-series stationary

- *MA: < Moving Average >* model exploits the relationship between the residual error and the observations

**Model parameters** - expect as input parameters 3 arguments (p,d,q)
- **p** is the number of lag observations *(using 4 lags)*
- **d** is the degree of differencing *(not-stationary, days, so placed "1" in model)*
- **q** is the size/width of the moving average window

# ARIMA / Linear Regression

**ARIMA** (or Box-Jenkins method) is a good model to be applied to this type of data (there is auto-correlation in the data).



MSFT Stock - Autocorrelation plot with lag=3



MSFT Stock Price - 5/10/20/50/100 Day Moving Avg



MSFT Stock Price Prediction - Linear Regression Model Performance

Linear Accuracy: 79.984%
Linear Regression prediction =

**ARIMA Test -** 75% data split test model to train/fit model *10.2 MSE

```
train_data, test_data = df[0:int(len(df)*0.75)], df[int(len(df)*0.75):]
training_data = train_data['Close'].values
```

**ARIMA MSE -** *10.2 MSE
*note: avg squared value across all the test set predictions

```
Testing Mean Squared Error is 10.164232398775123
```

**order: p=4 lags, d=1 (not stationary)**

```
model = ARIMA(history, order=(4,1,0))
```

```
futureDF['Predict']=model_fit.forecast(steps=futureDays)[0]
```

## ARIMA Model Results

| Dep. Variable: | | D.y | No. Observations: | | 295 |
|---|---|---|---|---|---|
| Model: | | ARIMA(4, 1, 0) | Log Likelihood | | -851.658 |
| Method: | | css-mle | S.D. of innovations | | 4.340 |
| Date: | | Tue, 09 Mar 2021 | AIC | | 1715.316 |
| Time: | | 17:48:27 | BIC | | 1737.438 |
| Sample: | | 1 | HQIC | | 1724.174 |

| | coef | std err | z | P>\|z\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| const | 0.2450 | 0.210 | 1.168 | 0.243 | -0.166 | 0.656 |
| ar.L1.D.y | -0.3073 | 0.058 | -5.275 | 0.000 | -0.421 | -0.193 |
| ar.L2.D.y | 0.0593 | 0.061 | 0.976 | 0.329 | -0.060 | 0.178 |
| ar.L3.D.y | 0.0632 | 0.061 | 1.036 | 0.300 | -0.056 | 0.183 |
| ar.L4.D.y | -0.0200 | 0.058 | -0.342 | 0.733 | -0.134 | 0.095 |

### Roots

| | Real | Imaginary | Modulus | Frequency |
|---|---|---|---|---|
| AR.1 | -1.5637 | -1.2121j | 1.9785 | -0.3951 |
| AR.2 | -1.5637 | +1.2121j | 1.9785 | 0.3951 |
| AR.3 | 3.1479 | -1.7014j | 3.5783 | -0.0789 |
| AR.4 | 3.1479 | +1.7014j | 3.5783 | 0.0789 |

Linear Accuracy: 81.415%

# Reddit Data



Reddit Data Change

MSE: 0.15319122649933242, R2: 0.8150623337916859

```
X = post_data[['post_count_change', 'avg_score_change', 'Volume_change', 'Count_Score_300']]
y = post_data['Adj Close_change'].values.reshape(-1, 1)
```

Northwestern

# Stock Predictor Web Page

# Potential Next Steps

- Additional historical stock data such as short activity to focus on potential GameStop(GME) or AMC activity.

- If data can't be found, warehouse daily to build missing data – screen-scaping(beautiful-soup) is an option from sites such as [shortsqueeze](#) and/or [Daily Short Sale Volume](#).

- Further investigation of ARIMA model – test/track various P/D/Q configs

- Reduce time grain to hourly, minute for automated trading/prescription reporting (take automatic action on report).

- Build additional charts and dashboards to track multiple stocks.

# QUESTIONS?