



STOCK MARKET FORECASTING

a machine learning case study

Table of Contents

- Business Background
- Data Understanding and Exploration
- Modeling and Evaluation
- Deployment
- Conclusion and Recommendation



Background

Data Scientist working at a brokerage firm.

As the data scientist of the firm, we are given a task to create a model to help stock market forecasting .

Business Objectives

The objective of this project is to help stock investors gain margin or high advantage in the stock market



Expected Output

The output of this project is a model prediction on the direction of the market stock trend to help analyst in stock market forecasting.

Project Limitation

With the limitation of time and data we will be focusing only on a certain number of stock exchange indexes

Analytic Approach

- Machine Learning Technique

Supervised learning (classification) to predict market direction

- Performance Measures

Precision, Recall, F1, Accuracy

Data Collection

The dataset of stock exchange market is taken from Kaggle [Stock Exchange Data](#) by Mattiuzc

Data Shape

104,229 rows

9 columns

Data Description

Dataset statistics	Count	Variable Types	Count
Number of variables	9	Categorical	2
Number of observations	104224	Numeric	7
Missing cells	0		
Missing cells(%)	0.0%		
Duplicate	0		
Total size in memory	7.2MiB		
Average record size in memory	72.0 B		

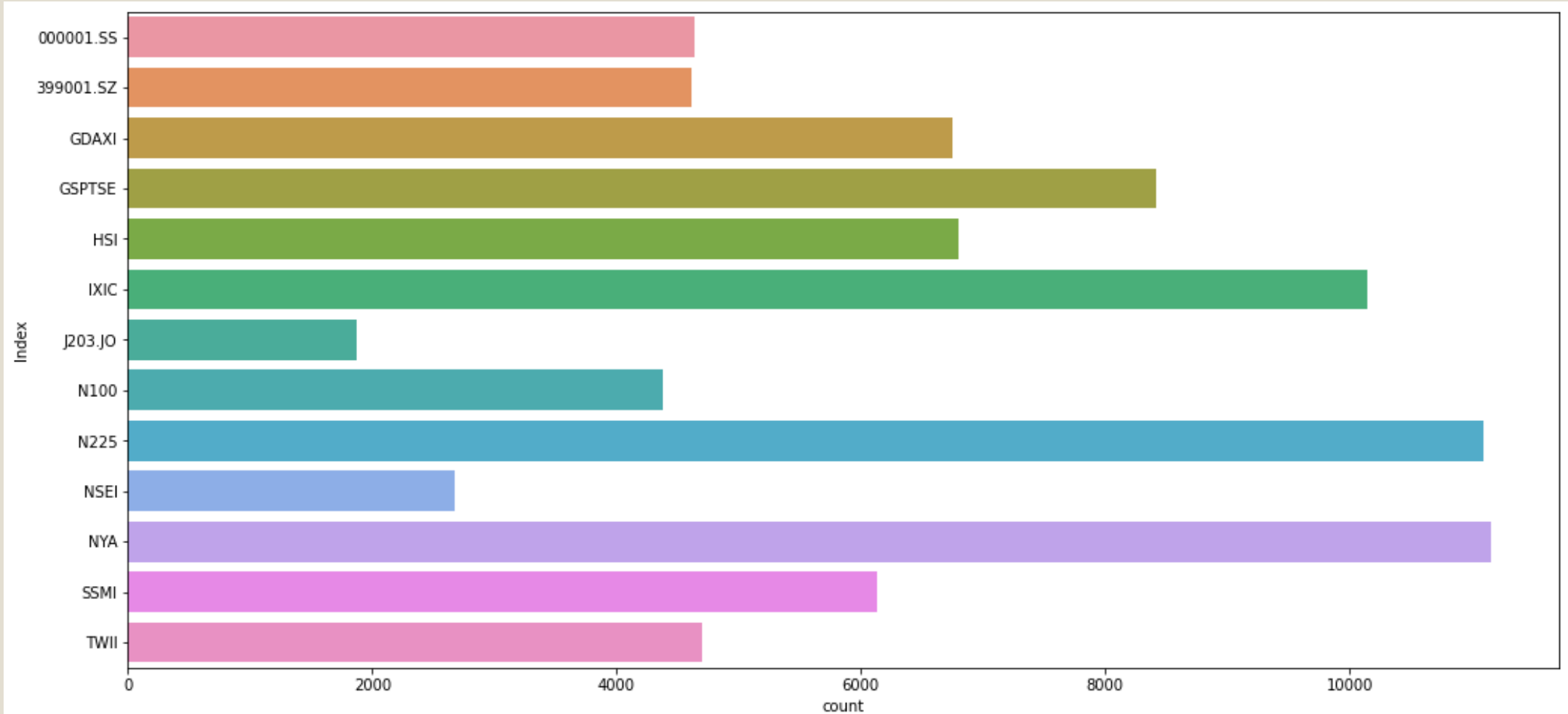
Column Explanation

No	Feature	Data Type	Description	Feature Type
1	Index	Object	Stock Indexes	Categorical
2	Date	Object	Time of Observation	Categorical
3	Open	Float64	Opening price	Numerical
4	Hight	Float64	Highest price during trading day	Numerical
5	Low	Float64	Lowest price during trading day	Numerical
6	Close	Float64	Closing price	Numerical
7	Adj Close	Float64	Closing price adjusted for dividends and stock splits	Numerical
8	Volume	Float64	Number of shares traded during trading day	Numerical
9	CloseUSD	Float64	Close price in terms of USD	Numerical

Splitting Data set

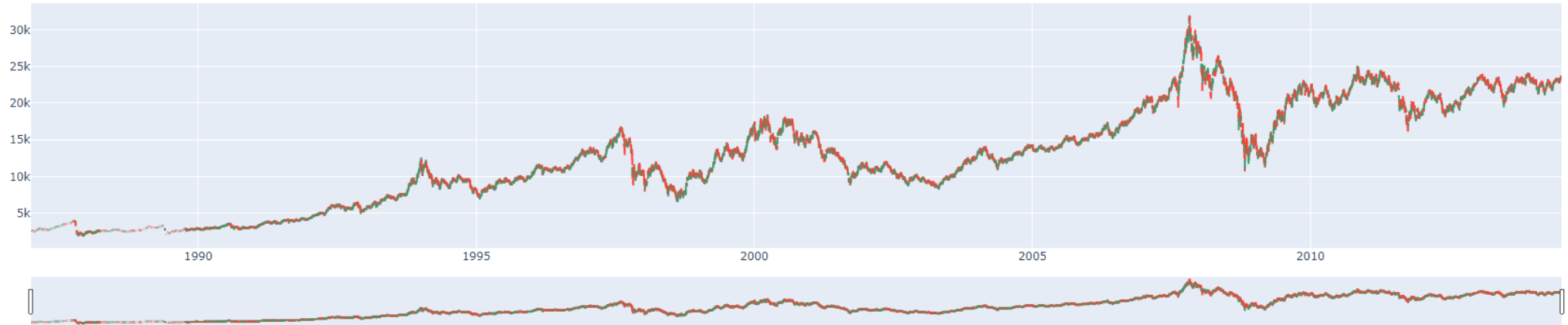
Before preprocessing we will be splitting each indexes price history based on the date of the day. The split will be a 80 %, 20% split.

Index's record distribution

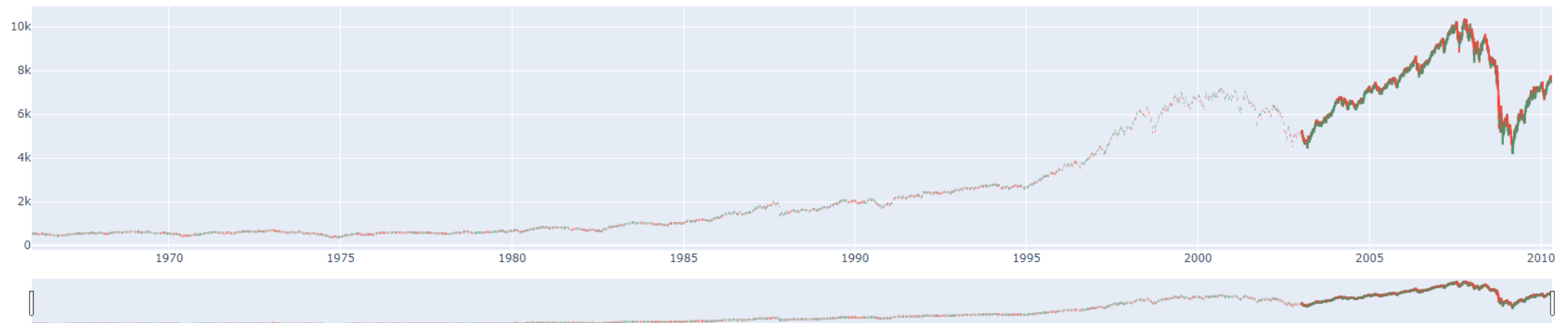


Opening and closing price based on their respective date

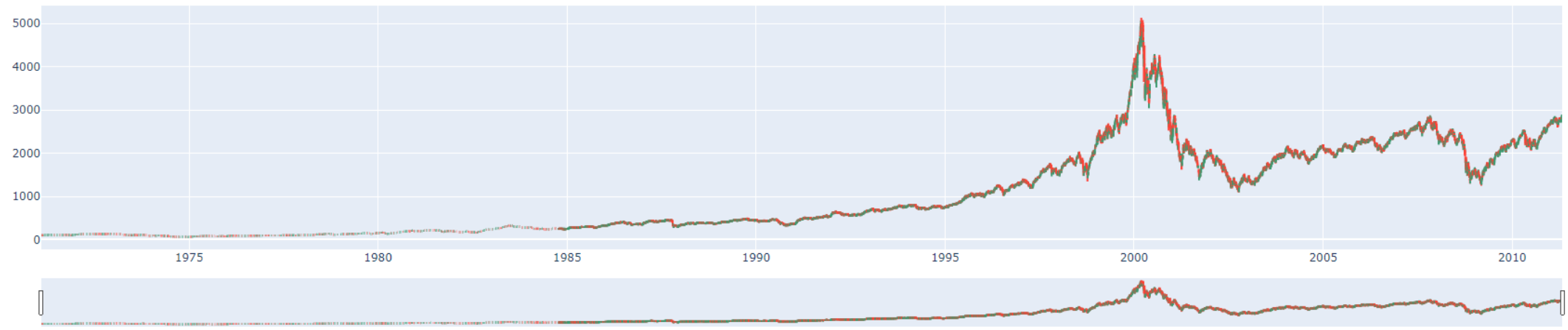
HSI



NYA



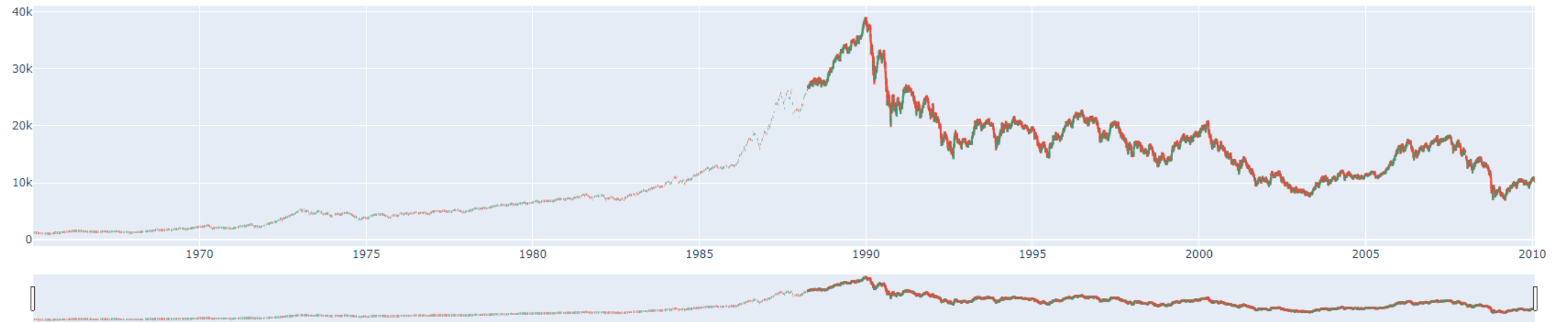
IXIC



000001.SS



N225



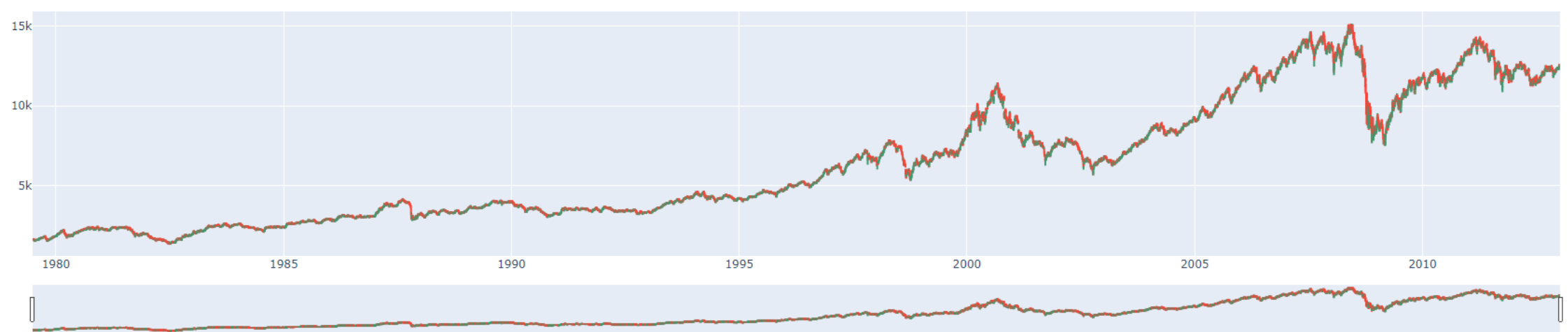
N100



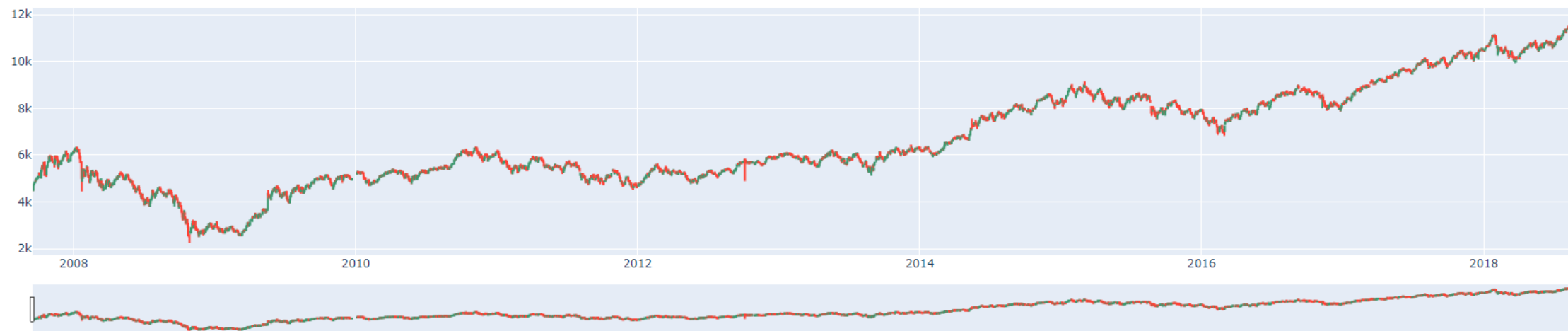
399001.SZ



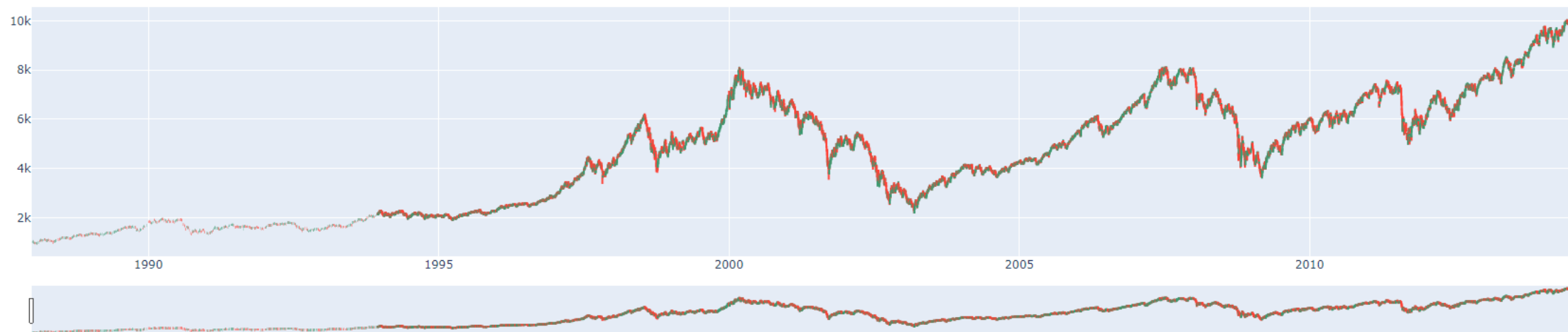
GSPTSE



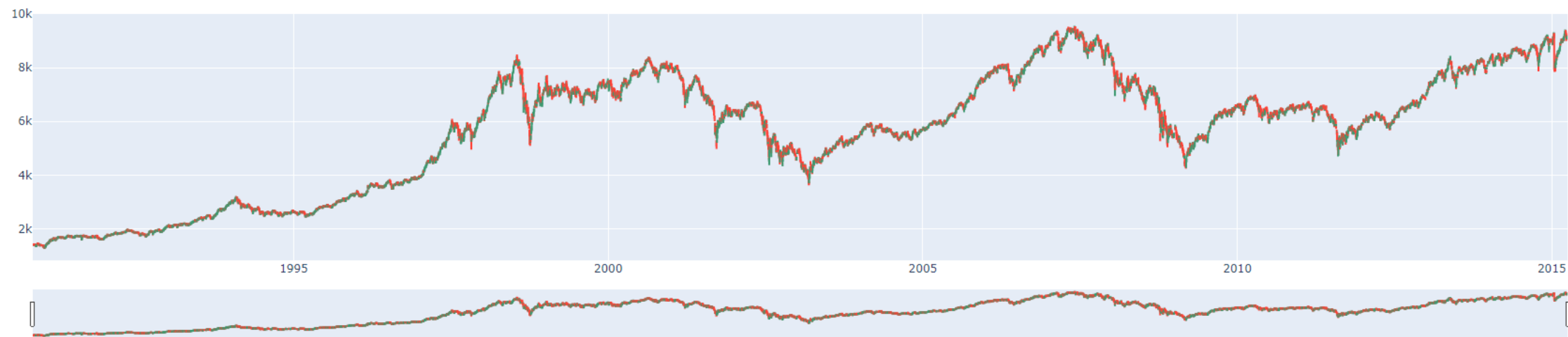
NSEI



GDAXI



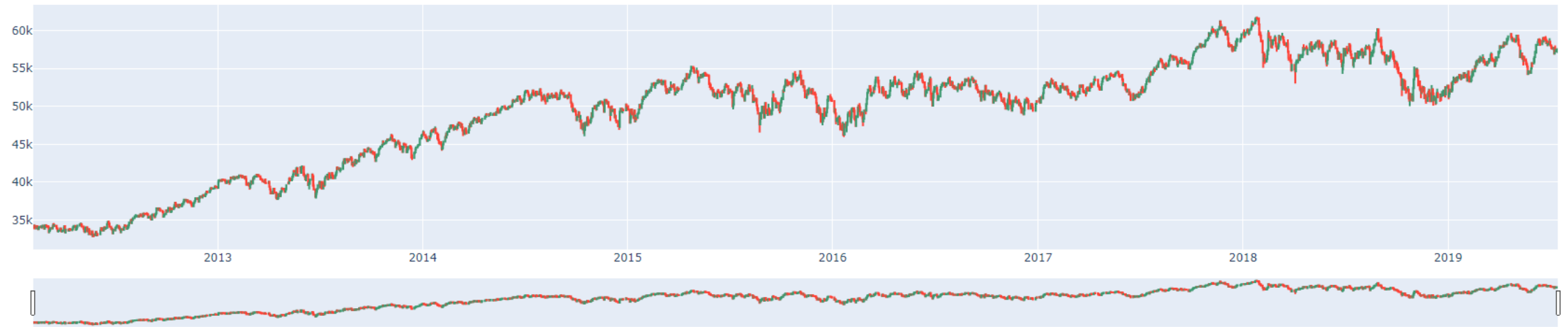
SSMI



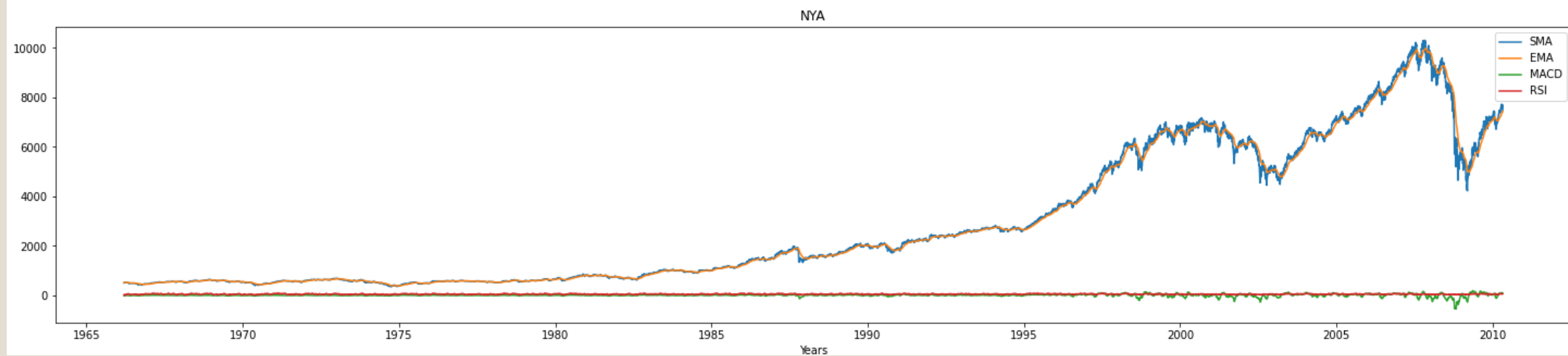
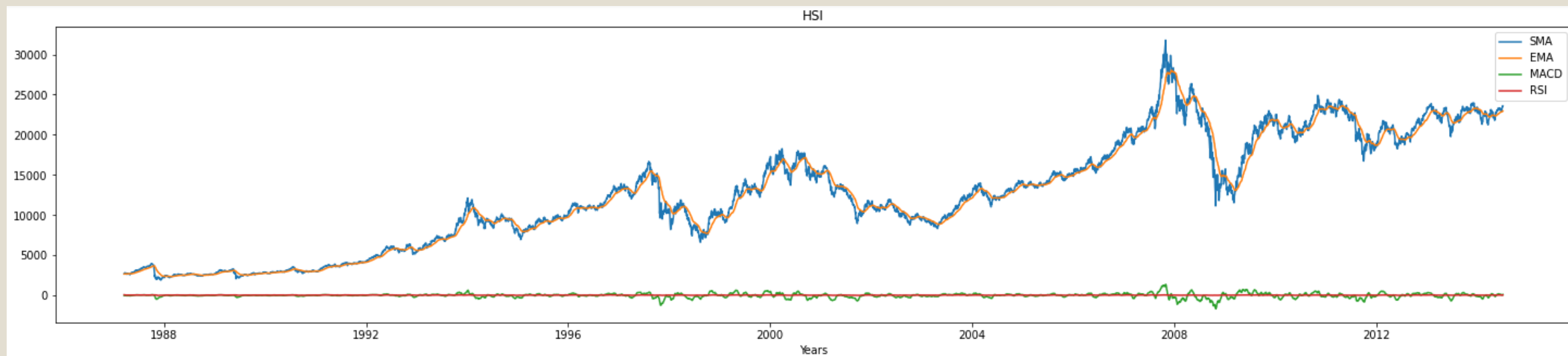
TWII

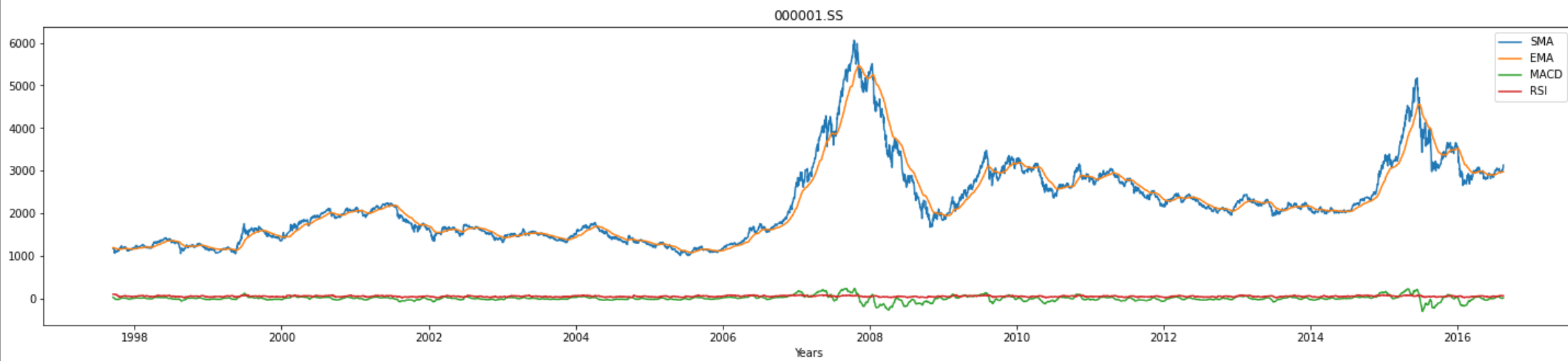
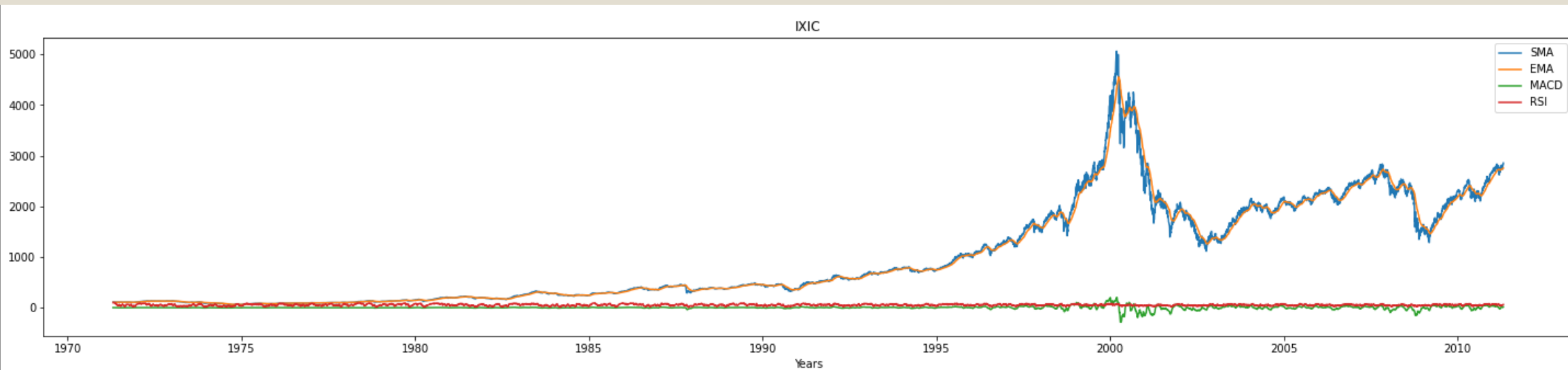


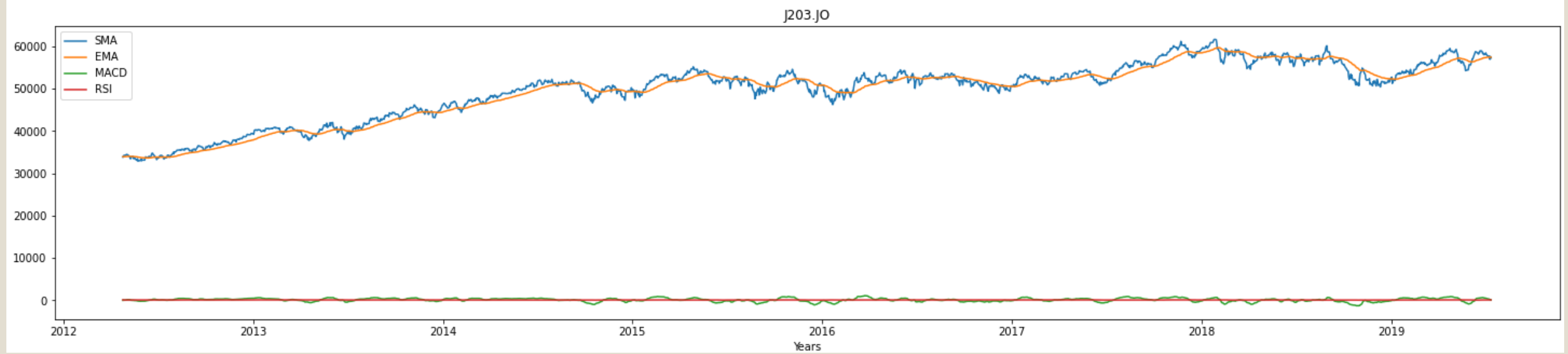
J203.JO



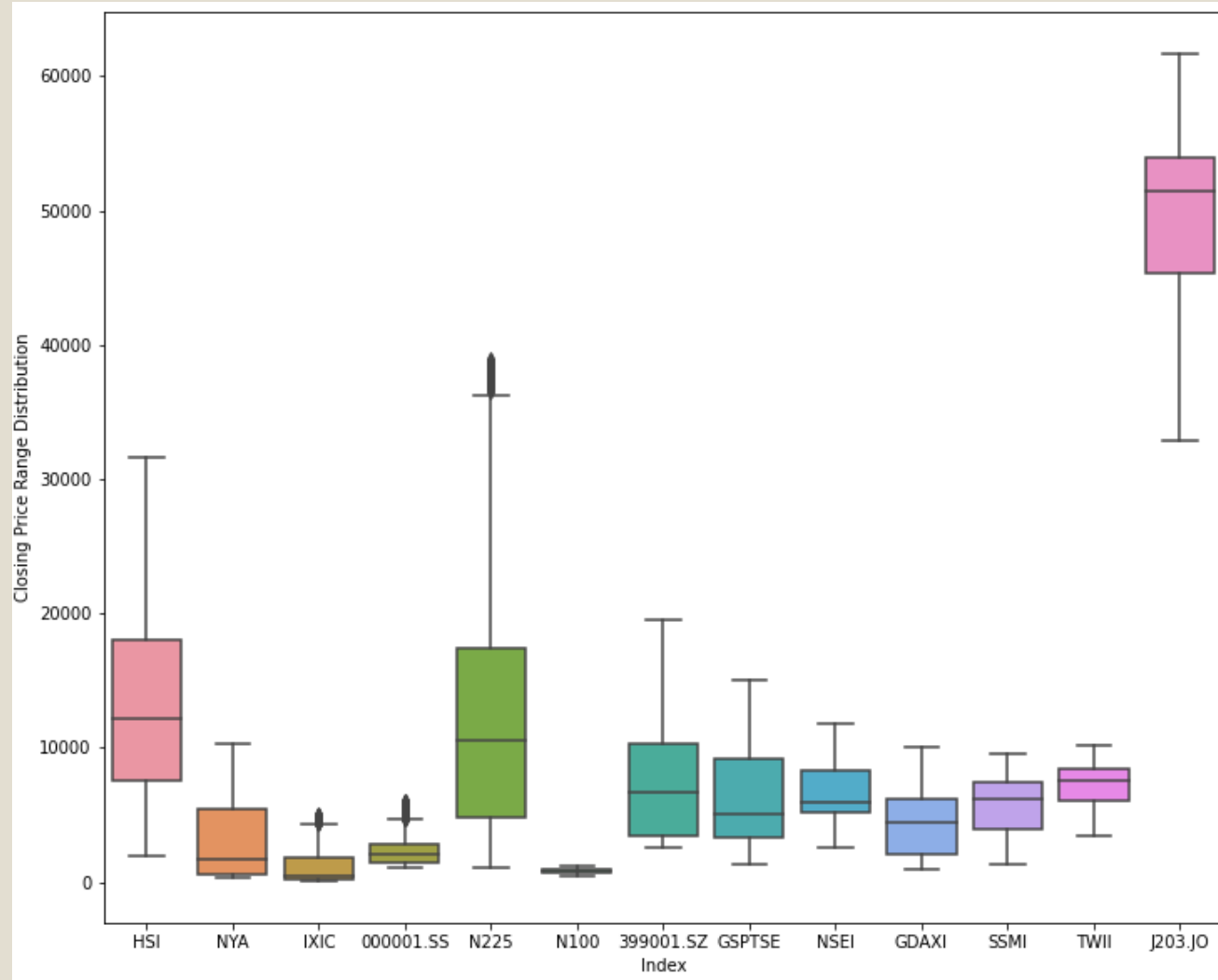
There are patterns occurring in the graph throughout the year we will try to generalize this pattern by smoothing the graph using moving averages.





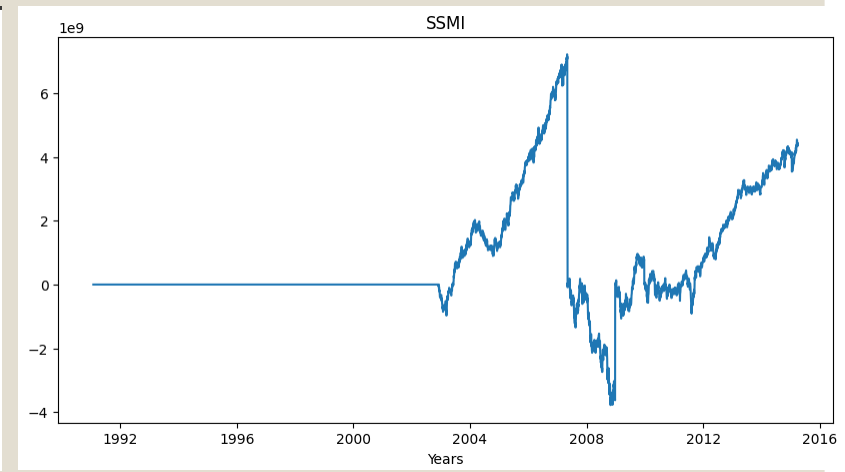
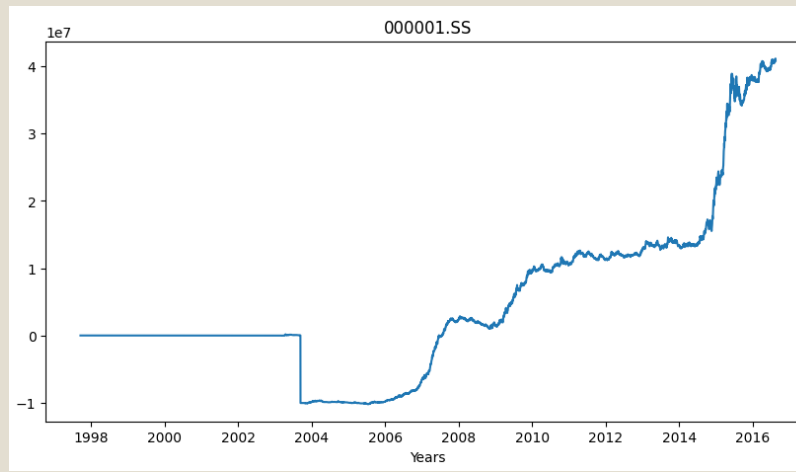
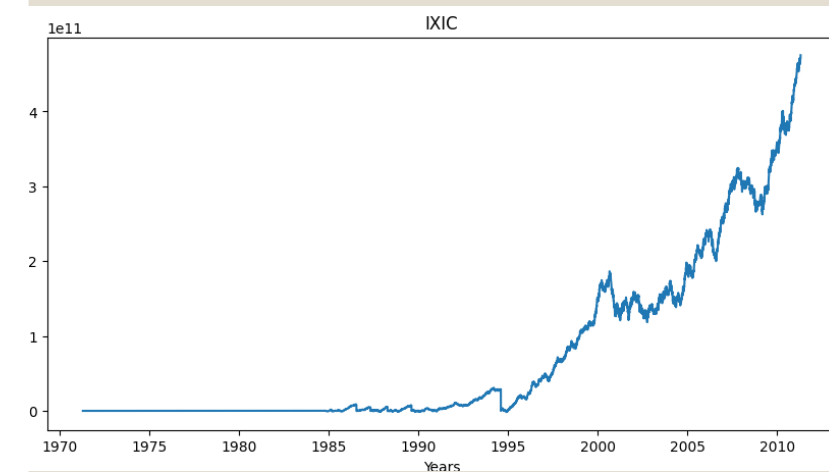
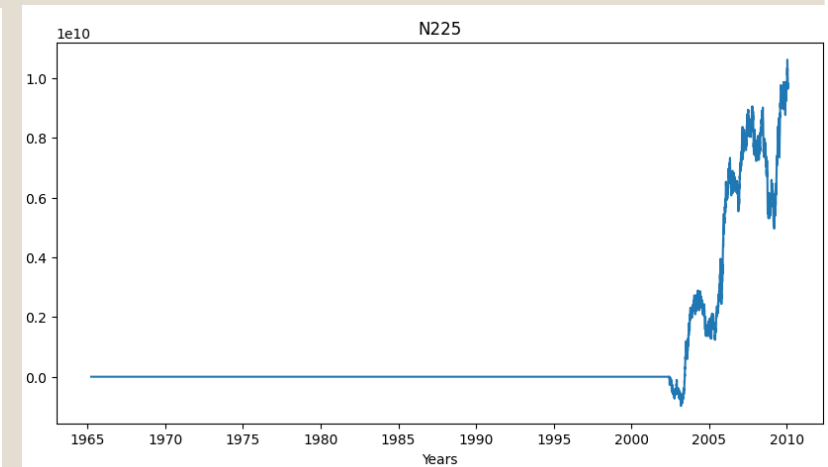
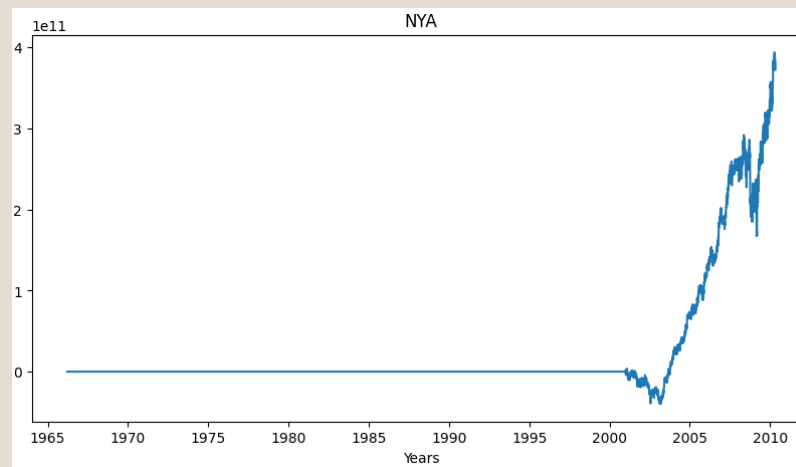
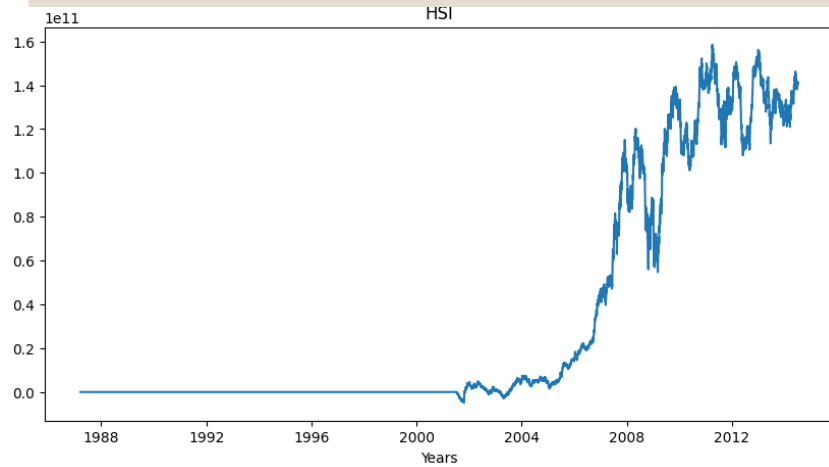


Closing price distribution



- J203.JO shows a different price range than the other indexes

On Balance Volume of each indexes



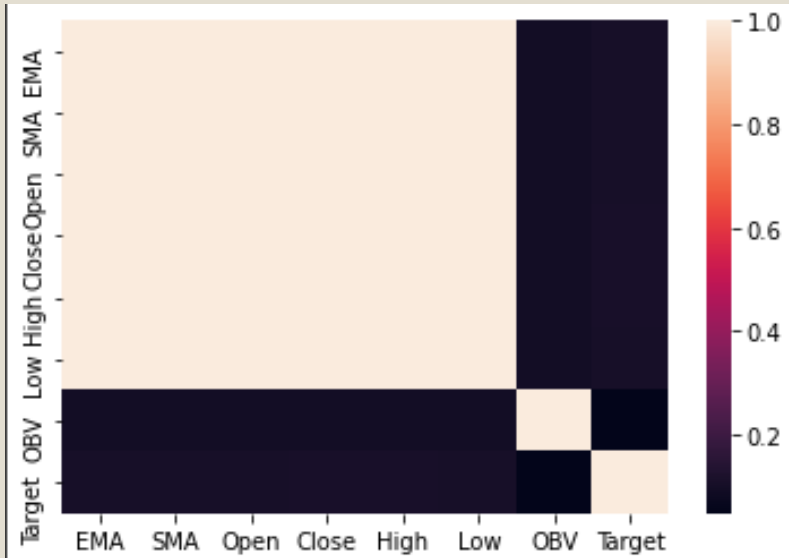
Data Preprocessing

- Feature engineering
- Feature selection
- Encoding

Feature engineering

- SMA (Simple Moving Averages)
$$= \frac{(A_1 + A_2 + \dots + A_n)}{n}$$
- EMA (Exponential Moving Averages)
$$= \text{Price Today} \times \left(\frac{\text{Smoothing}}{1 + \text{Days}} \right) + \text{EMA}$$
$$\text{Yesterday} \left(1 - \left(\frac{\text{Smoothing}}{1 + \text{Days}} \right) \right)$$
- CDMA (Convergence Divergence Moving Averages) = 12-EMA – 26-EMA
- RSI (Relative Strength Index)
$$= 100 - \left[\frac{100}{1 + \frac{\text{Ave Gain}}{\text{Ave Loss}}} \right]$$
- OBV (On Balance Volume)

Feature Selection & Encoding



- Since SMA have a high linear correlation with price columns we will use the SMA and drop the other three since SMA generalize the trend of the graph better.

- Next step we hot encode the indexes column

Modeling Baseline

Baseline model using
DummyClassifier

Precision scores:
0.501115055190106

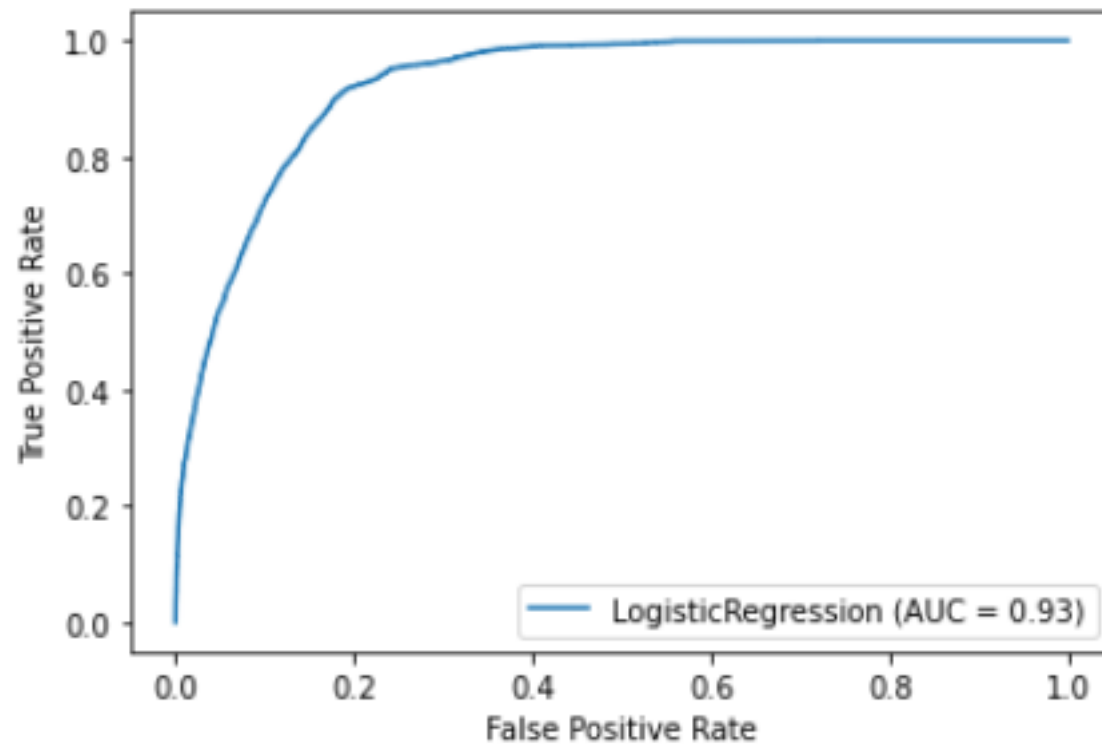
Recall scores:
0.501114973732381

F1 scores:
0.5011050082000984

Model Selection

	Model	Precision	Recall	F1	ROC AUC
0	LogisticRegression	0.776489	0.780509	0.778000	0.928984
1	LGBM	0.820723	0.750410	0.778058	0.937964
2	CatBoost	0.810967	0.749652	0.774454	0.936750
3	Decision Tree	0.745111	0.746629	0.745793	0.746629
4	XGBoost	0.826467	0.744556	0.775617	0.937831
5	Random Forest	0.782515	0.740514	0.758434	0.929760

ROC AUC Curve



- Since we want to improve on recall we have to move the threshold to get a better score

Model Evaluation on unseen data

	precision	recall	f1-score	support
0	0.82	0.98	0.89	14391
1	0.89	0.48	0.62	5758
accuracy			0.83	20149
macro avg	0.86	0.73	0.76	20149
weighted avg	0.84	0.83	0.82	20149

Conclusion

- Logistic regression still prove to be one of the best model to be able to predict the stock market direction trend
- The model could help analyst predict the market stock price with promising result in given more time and resource.