

Stock Market Forecasting using Data Mining Technique

Presenter:

Marina Shah binti Muhammad Zabri Tan (WQD 180011)

Objective

The objective of this report is to present framework of analysis based on clustering algorithm, decision tree and time series analysis to forecast stock market price.

Framework



Objective

To analyze and forecast stock market data



Data acquisition

Crawl data from websites
store data in MySQL database



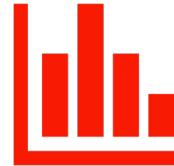
Data cleaning

Clean the data and convert it for further analysis



Analysis & Modeling

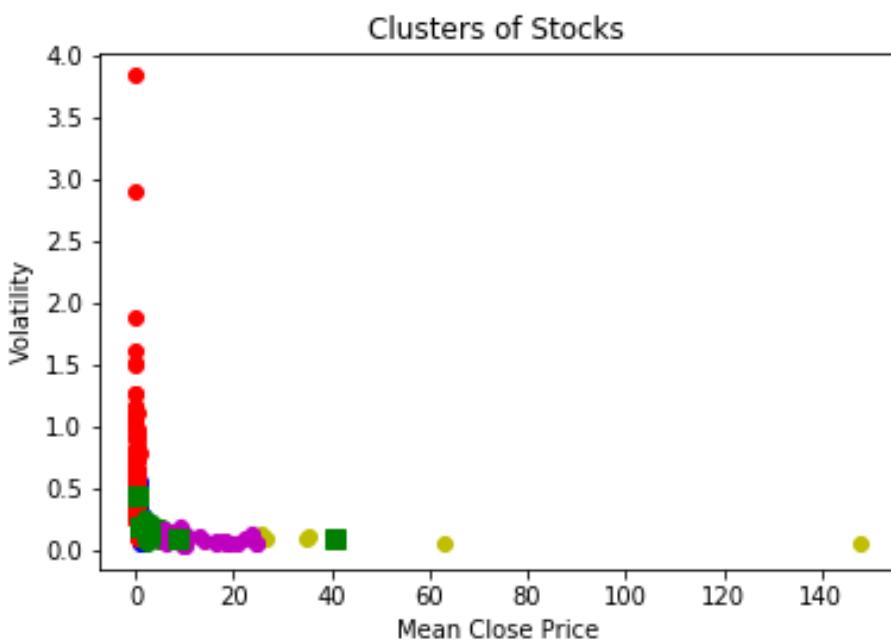
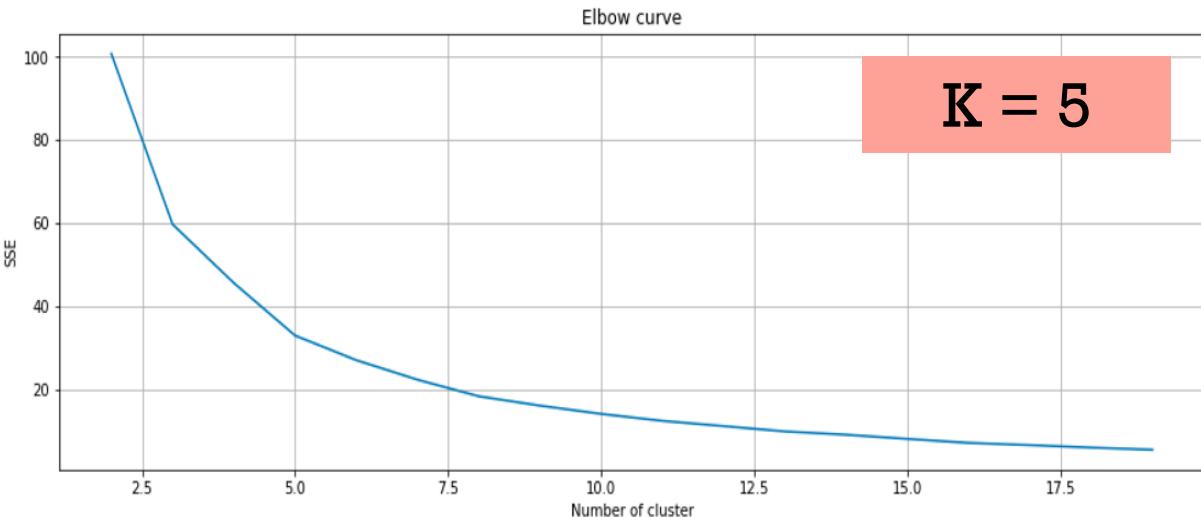
Descriptive
Explorative
Machine learning



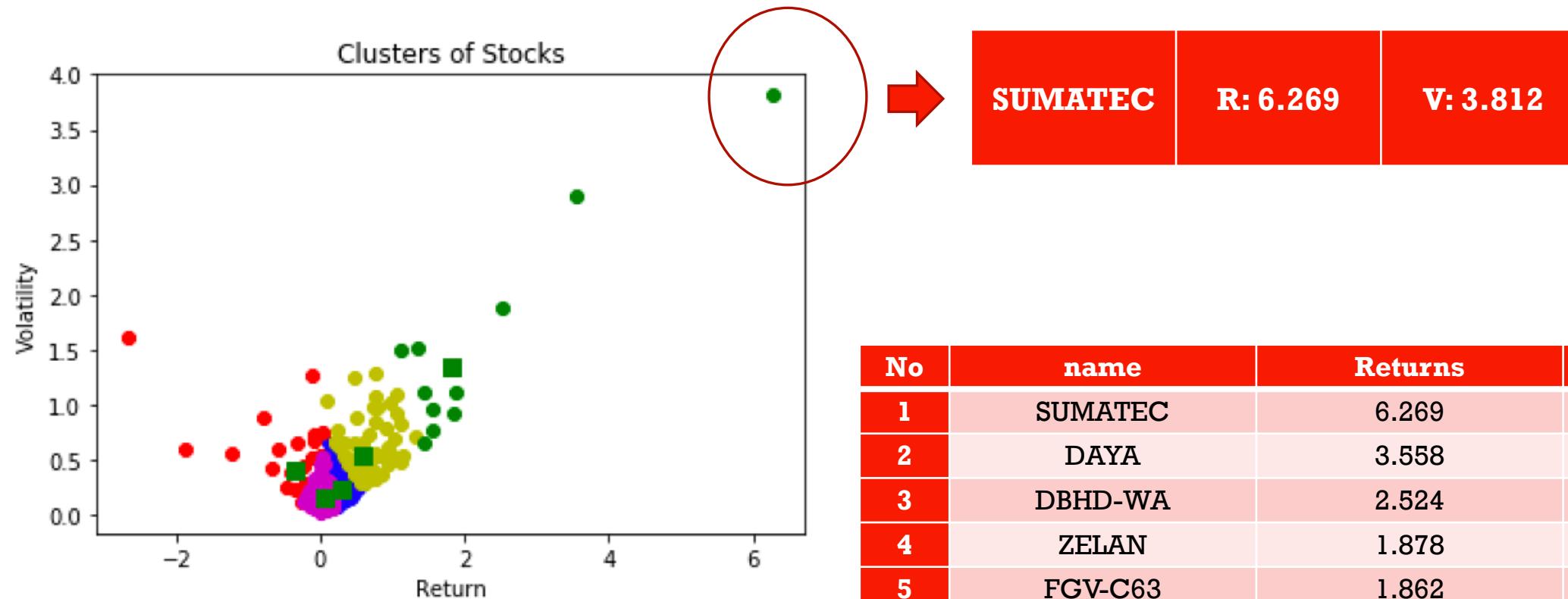
Visualization

Tables, graphs

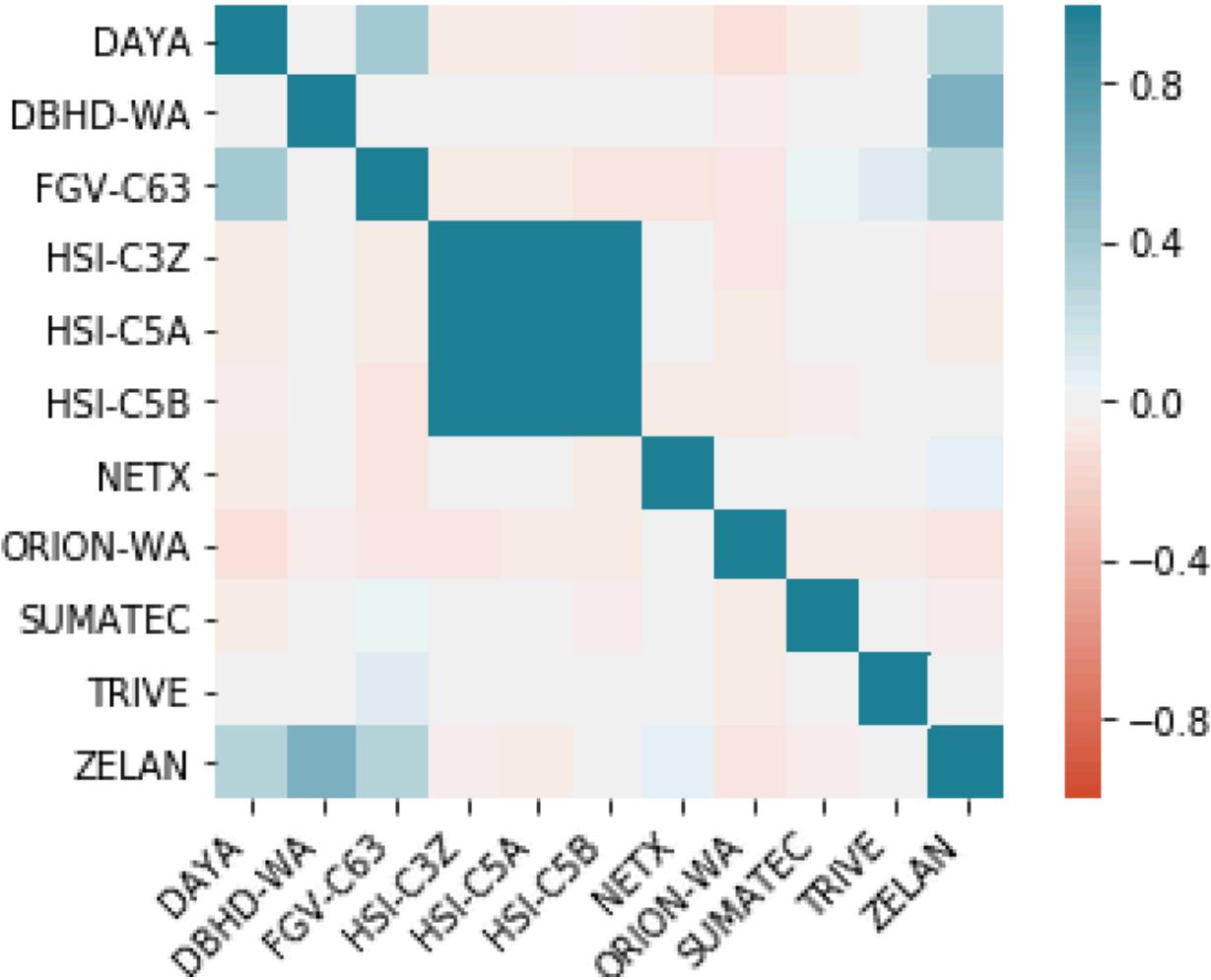
Clustering



| Name | Close Price | Volatility |
|---------|-------------|------------|
| NESTLE | 147.667 | 0.045 |
| SUMATEC | 0.0092 | 3.8428 |



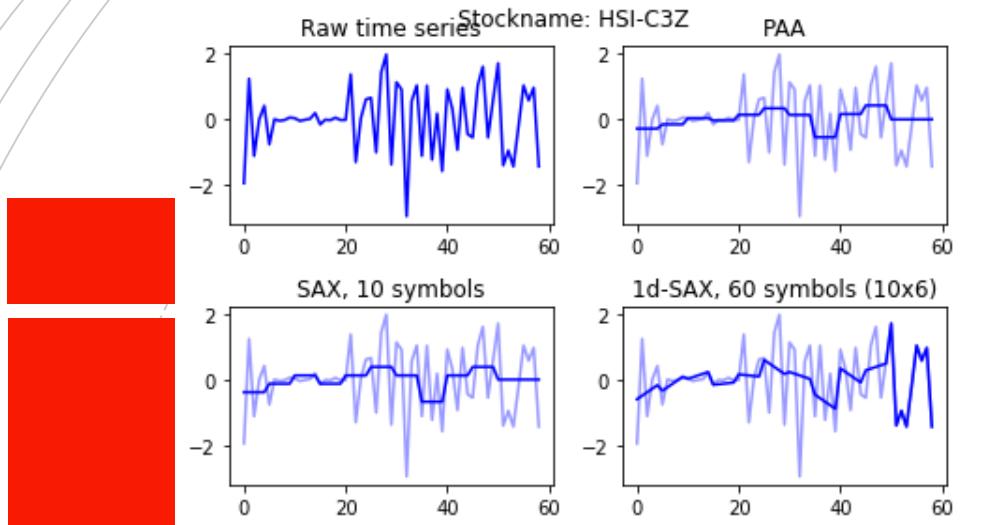
Clusters of stock: Volatility against Return



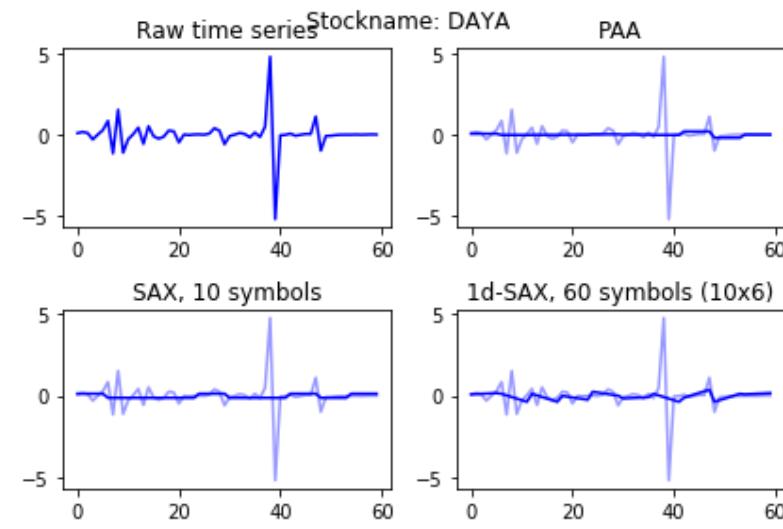
Most Positive correlation: **HSI-C3Z** **HSI-C5A** **0.998**

Most Negative correlation: **DAYA** **ORION-WA** **-0.106**

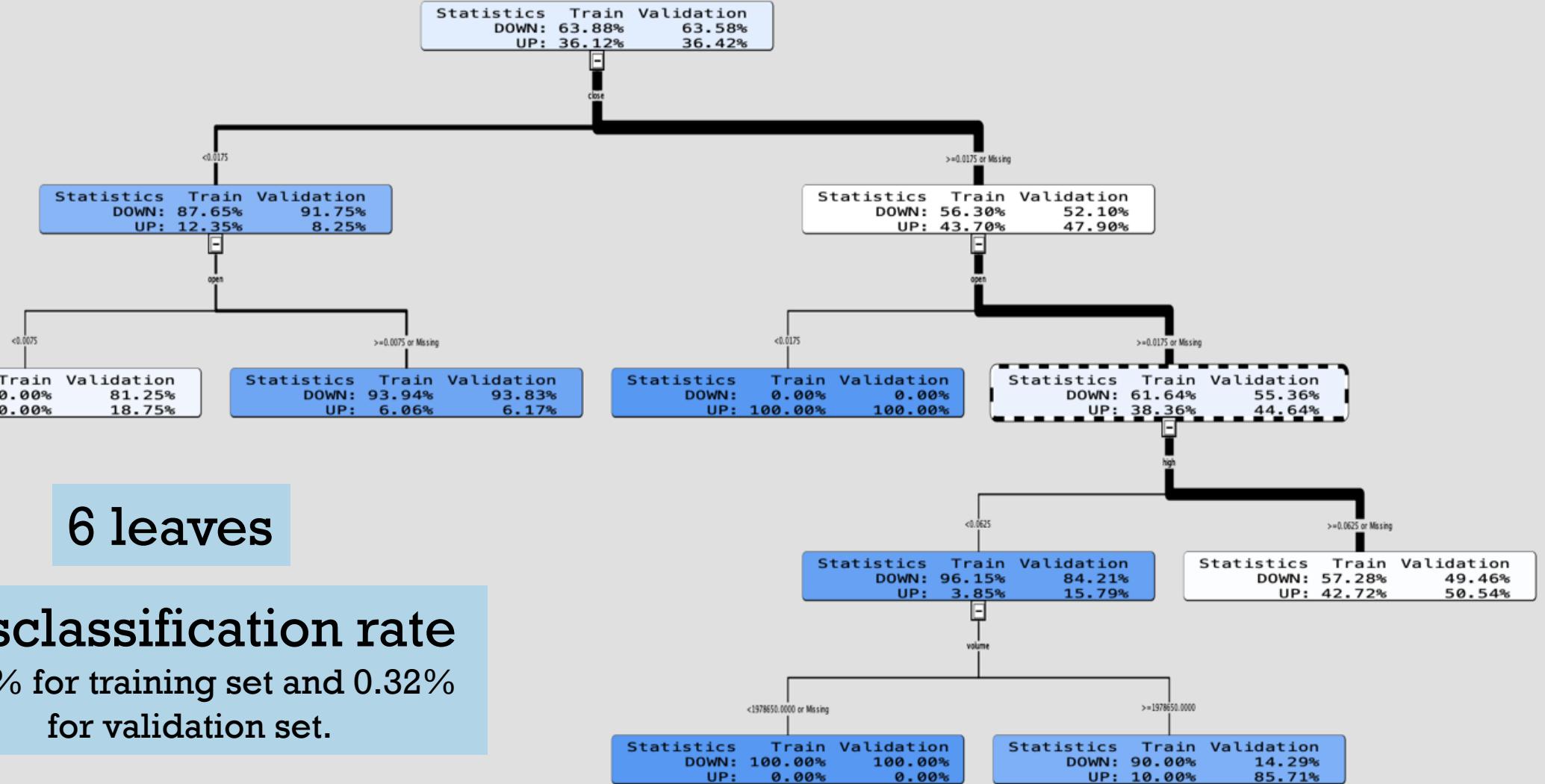
HSI-C3Z vs HIS-C5A

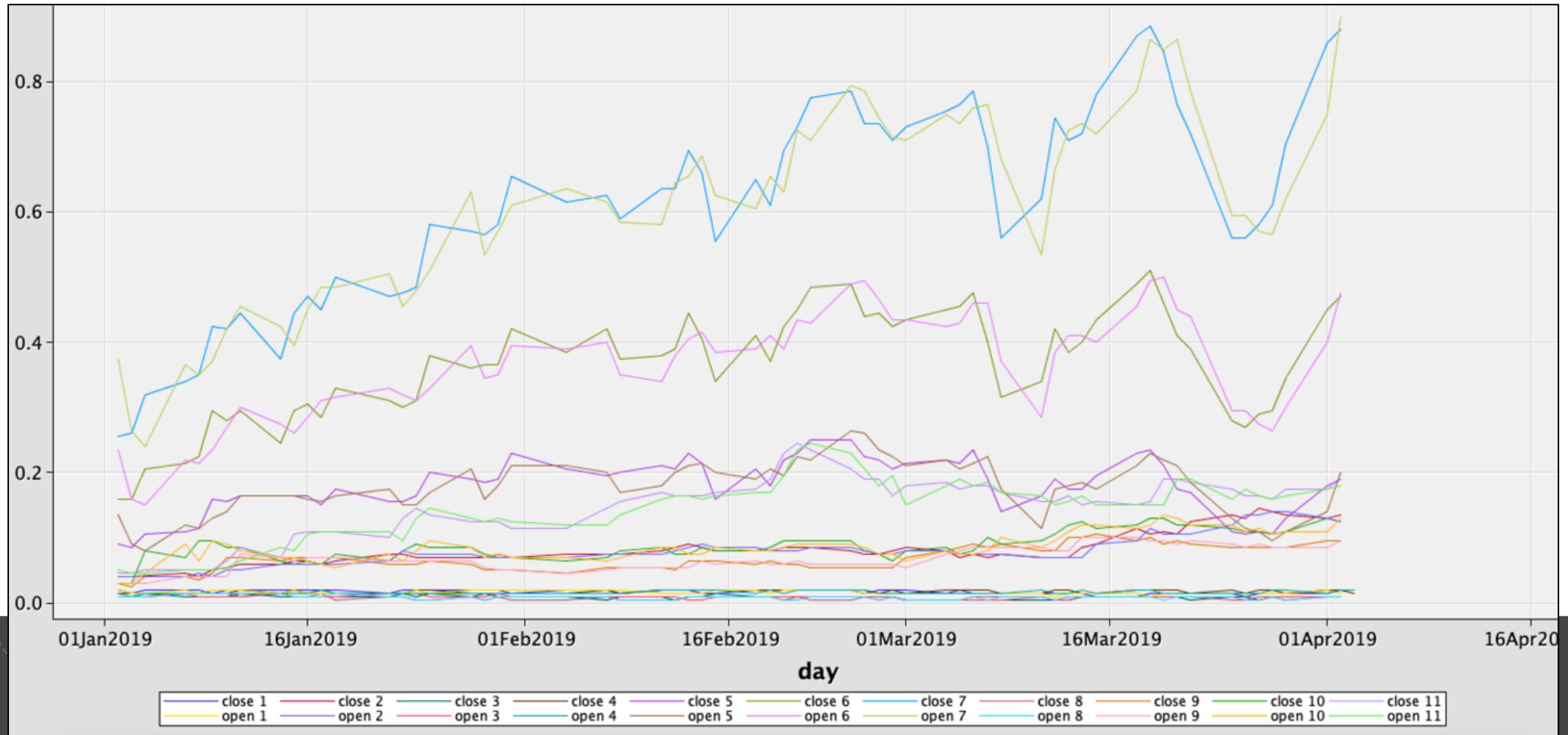
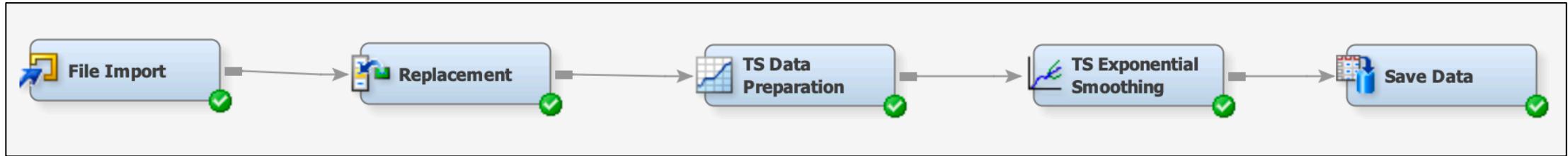


Daya vs Orion WA



SAX & PAA



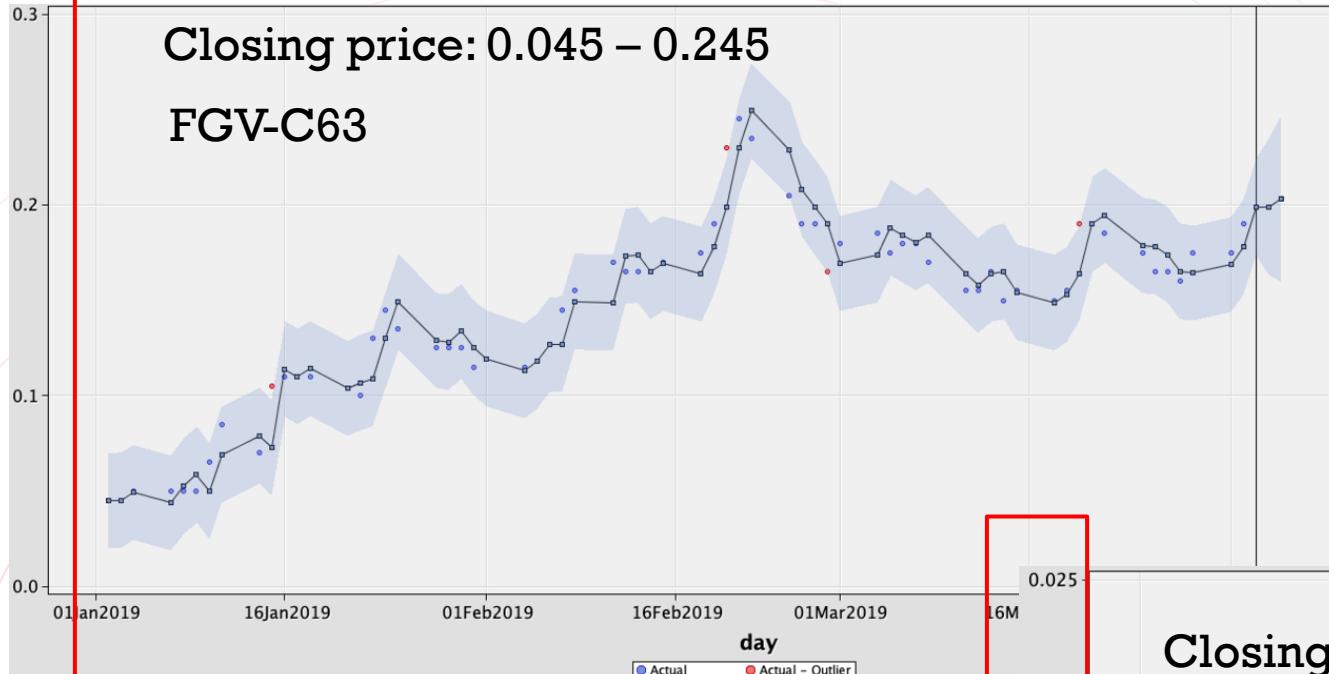


| Company Name | RMSE | R-Squared |
|---------------------|-------------|------------------|
| NETX | 0.0024 | 0.3385 |
| ORION-WA | 0.0063 | 0.9366 |
| DAYA | 0.0024 | 0.2178 |
| TRIVE | 0.0023 | 0.1650 |
| HSI-C3Z | 0.0236 | 0.6805 |
| HSI-C5A | 0.0423 | 0.7422 |
| HSI-C5B | 0.0590 | 0.8549 |
| SUMATEC | 0.0026 | 0.2563 |
| ZELAN | 0.0068 | 0.8714 |
| DBHD-WA | 0.0109 | 0.7599 |
| FGV-C63 | 0.0124 | 0.9330 |

Model Assessment

Closing price: 0.045 – 0.245

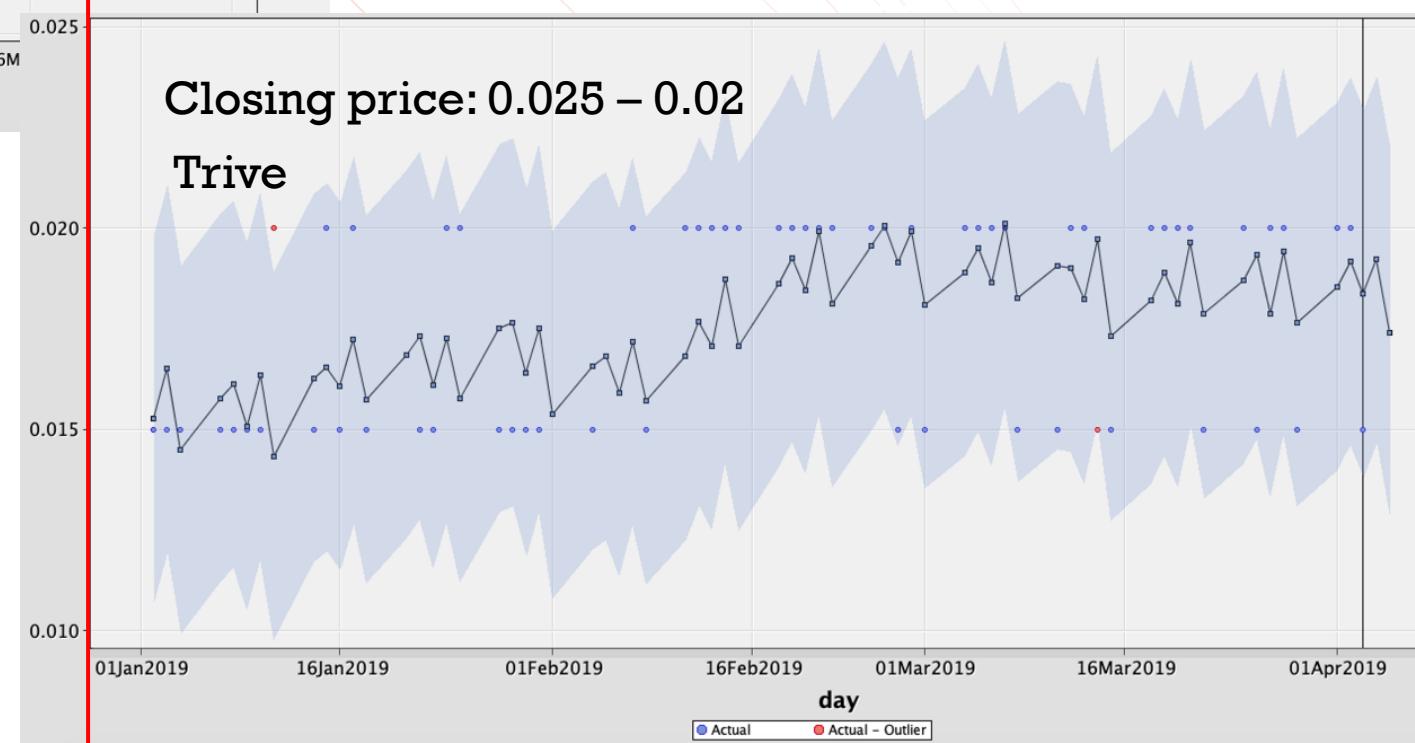
FGV-C63



| Date | Company | Predicted | Actual | % Difference |
|--------|---------|-----------|--------|--------------|
| 03-Apr | FGV-C63 | 0.199 | 0.190 | 5% |
| 04-Apr | FGV-C63 | 0.199 | 0.180 | 10% |
| 05-Apr | FGV-C63 | 0.203 | 0.180 | 11% |
| 03-Apr | Trive | 0.015 | 0.015 | 0% |
| 04-Apr | Trive | 0.019 | 0.015 | 22% |
| 05-Apr | Trive | 0.017 | 0.020 | -15% |

Closing price: 0.025 – 0.02

Trive



Conclusions

- It can be observed that the model performs differently depending on the value of the dependent variable.
- The performance of the machine learning models can be further improved with more data.