

QF604 Group Project Report



Team member:

Yu Lingfeng
Wang Wenjie
Wang Haotong
Chen Liangrui
Liu Jing



Contents

1 Abstract

2 Replication

3 Improvement

4 Conclusion



Abstract



Original Paper –

Target: predict clean energy stock price under the impact of climate changes

Method: Random Forests, Extremely Randomized Trees, Gradient Boosting and Support Vector Machine outperform Lasso and Naïve Bayes

Result: prediction accuracy up to 85% under Random Forests, GBM methods

Improvement –

Method: Linear Regression (which gives best result and more stable)

Result: prediction accuracy around 60%

2

Replication



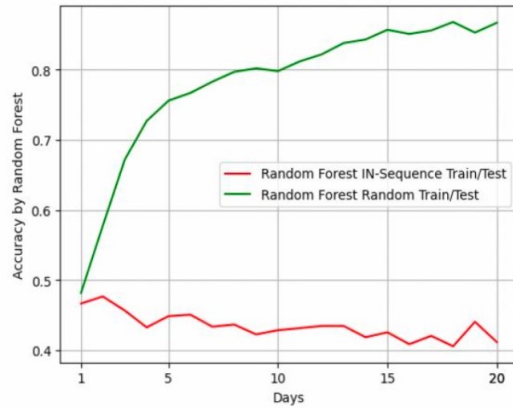
Study period: 1st, Jul 2022 -- 1st, Feb 2024

Train/Test Split: 70%/30%

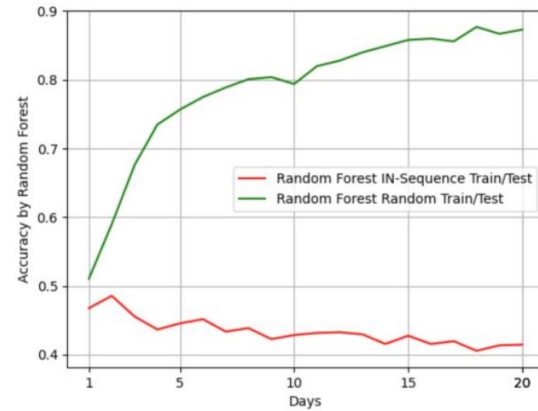


Random Forest

Number of trees: 100

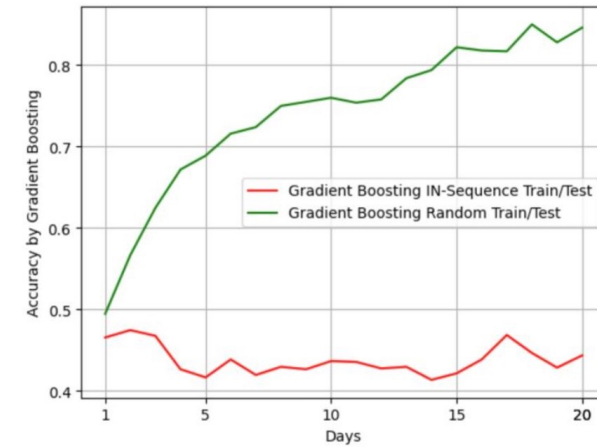


Number of trees: 500

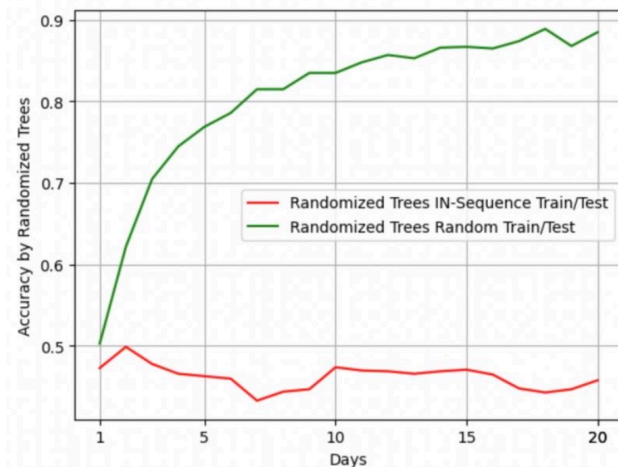


Gradient Boosting

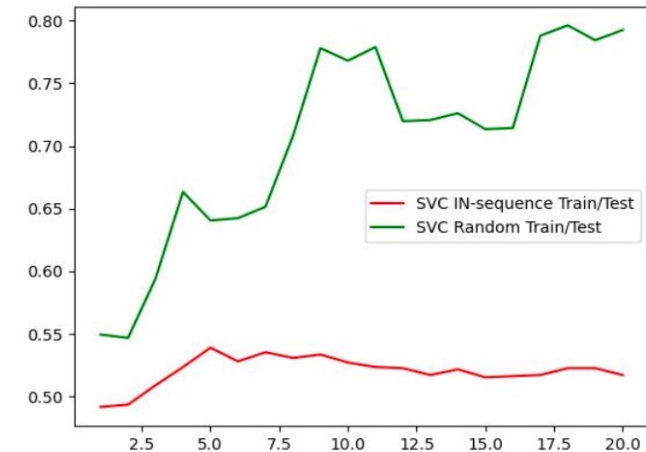
Number of boosting stages: 500 learning rate: 0.05



Extremely Randomized Tree

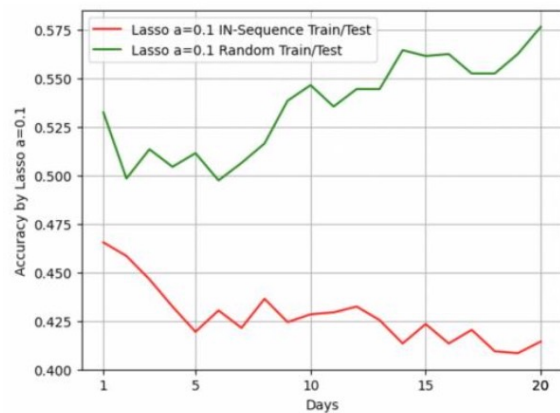


Support Vector Machine

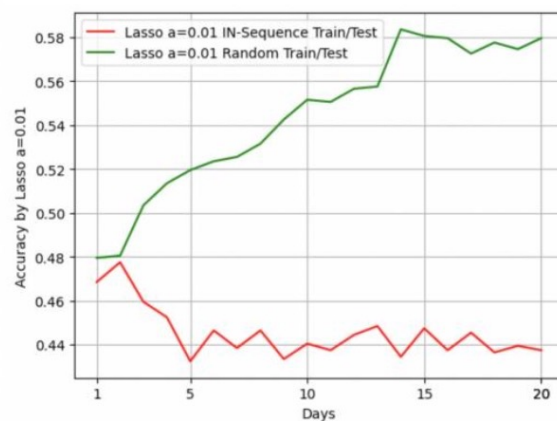


Lasso

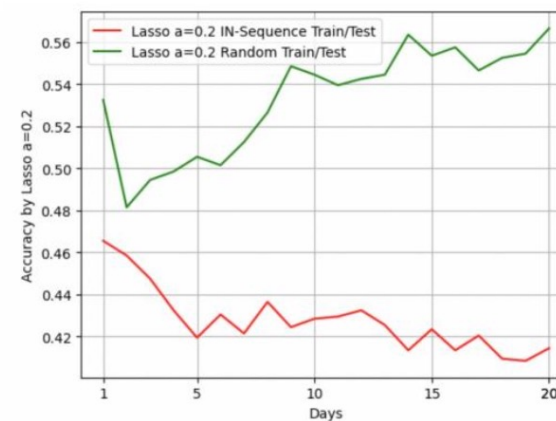
Alpha = 0.01



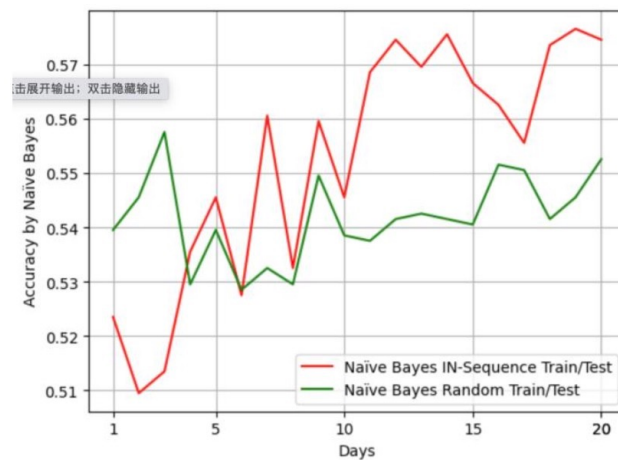
Alpha = 0.1

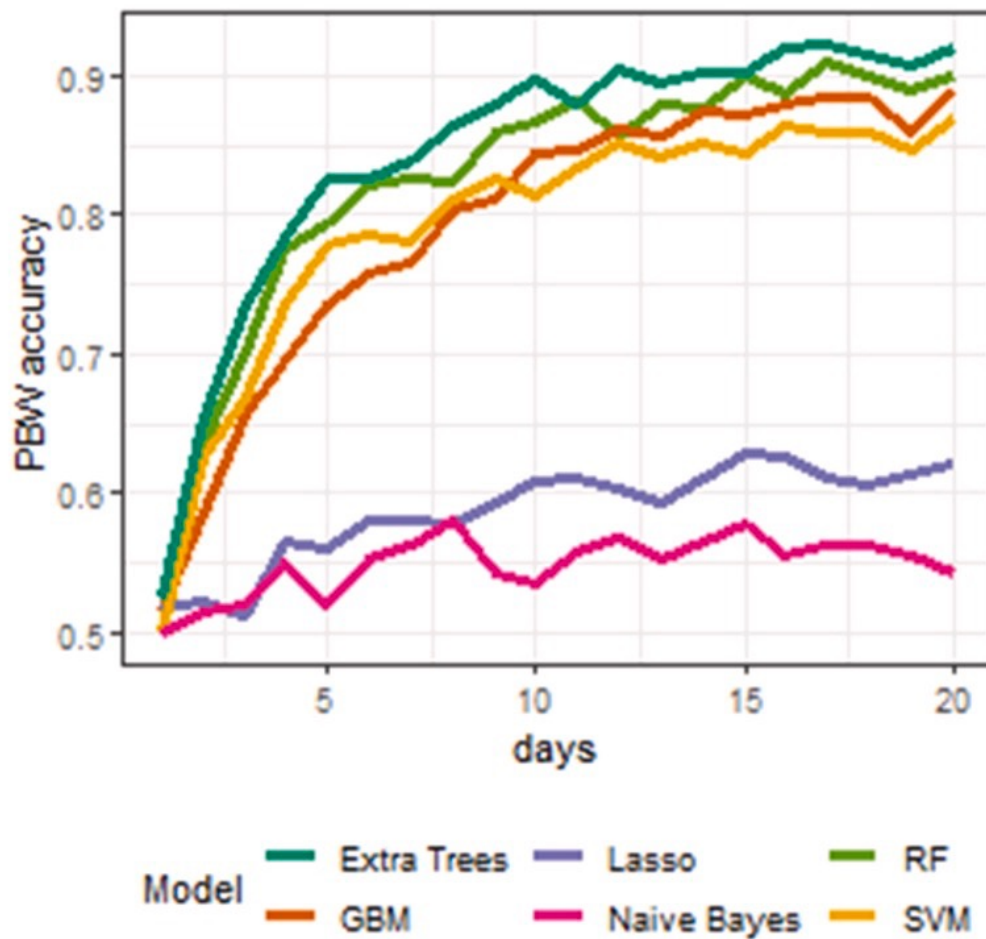


Alpha = 0.2



Naïve Bayes





Compare

The replicated results are **similar**, which is **as high as 90%**, just by doing a random partition on training and test set.



3

Improvement



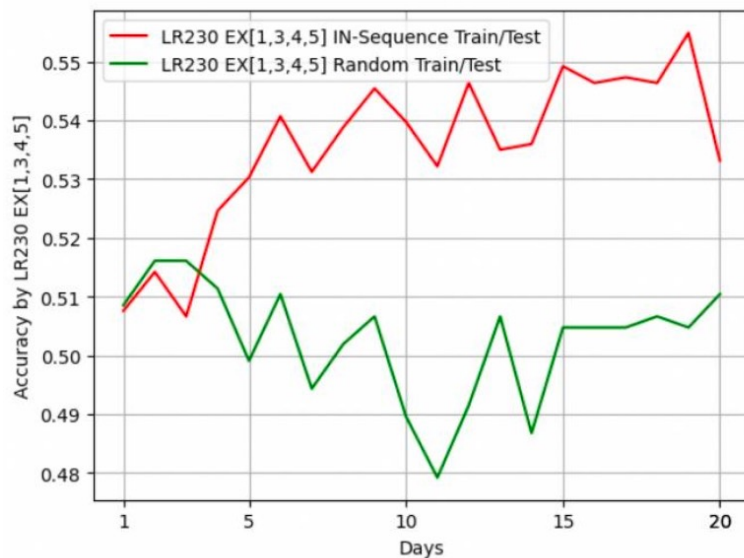
Improvement – Linear Regression

Input features: **9 features** (exclude EPU, EMU, IDT)

Do **correlation analysis** to identify effective predictors.

Do **mean correlation cut off** to select effective predictors.

Generate **511 combinations** and select **best two** performances.





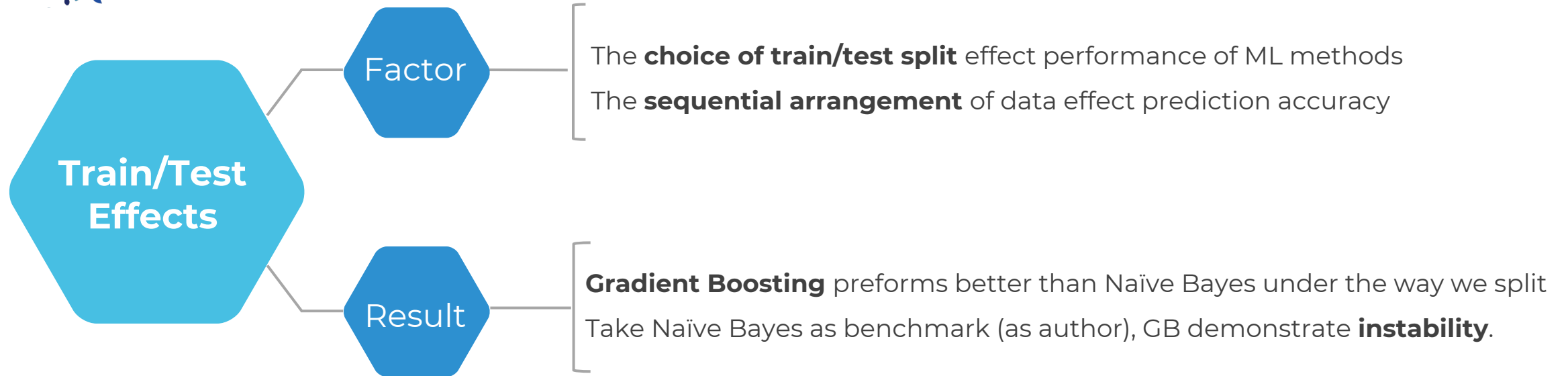
4

Conclusion

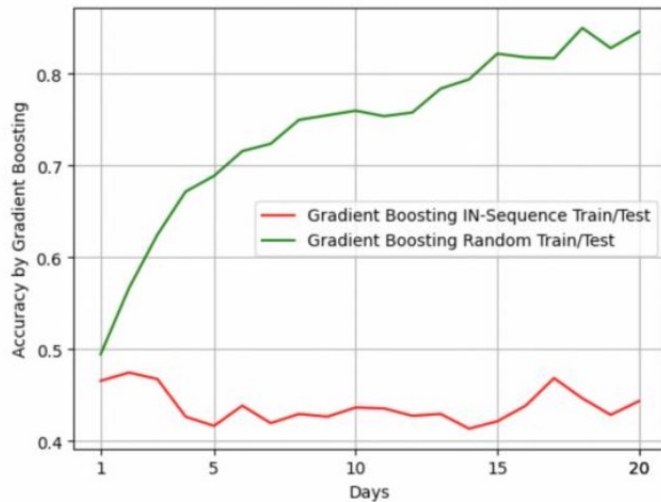




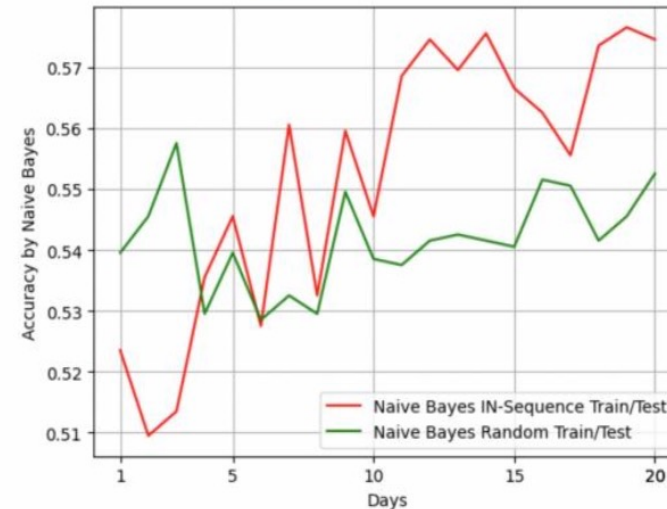
Conclusion 1 – Train/Test Effects



```
plot_result(PBW_GB_acc, 'Gradient Boosting')
```



```
plot_result(PBW_NB_acc, 'Naive Bayes')
```



04 Conclusion 2 – Features Set Effects

Features set Effects

Replication Part

High performance due to **improper partitioning train/test set**.
The assessment of feature significance remains **inconclusive**.

Some features exhibit **improper significance**. IDT is susceptible to **look-forward bias**.
Most technical Indicators have **high degree of freedom** and is subject to **overfitting**.

Improvement Part

Feature selection and tuning can improve prediction accuracy.

The testing window we introduced, with a downward trend, limit the testing results.

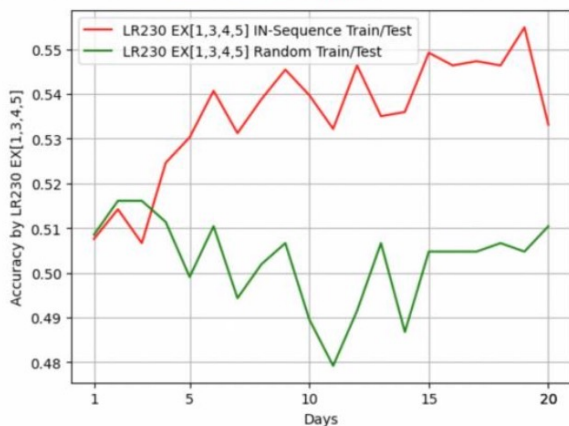
The best linear model is feature combinations EX1, EX3, EX4, EX5.

The features selected by the author have **lower prediction accuracy** than we added.

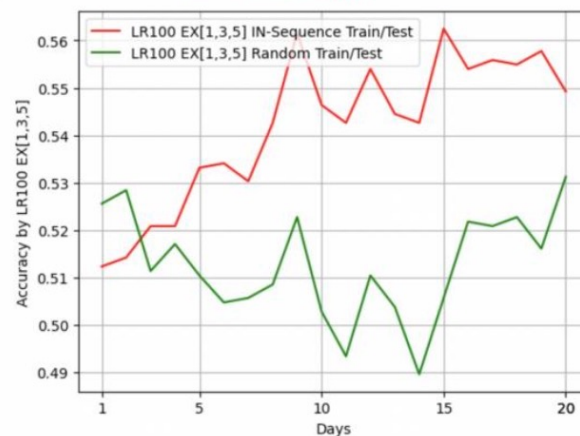
```
plot_result(PBW_LR_all, 'Linear Regression original')
```



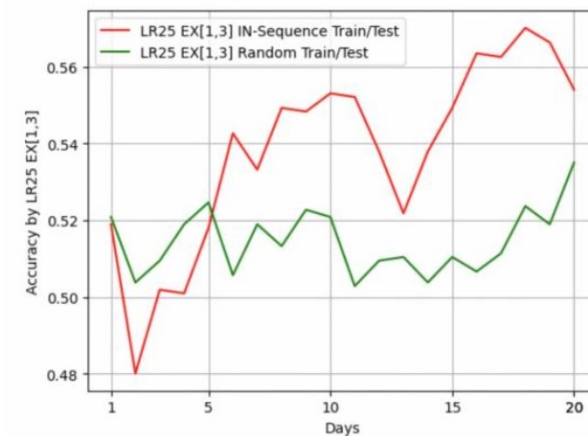
```
# LR230 ('EX1', 'EX3', 'EX4', 'EX5')  
plot_result(all_acc[230], 'LR230 EX[1,3,4,5]')
```



```
plot_result(all_acc[100], 'LR100 EX[1,3,5]')
```

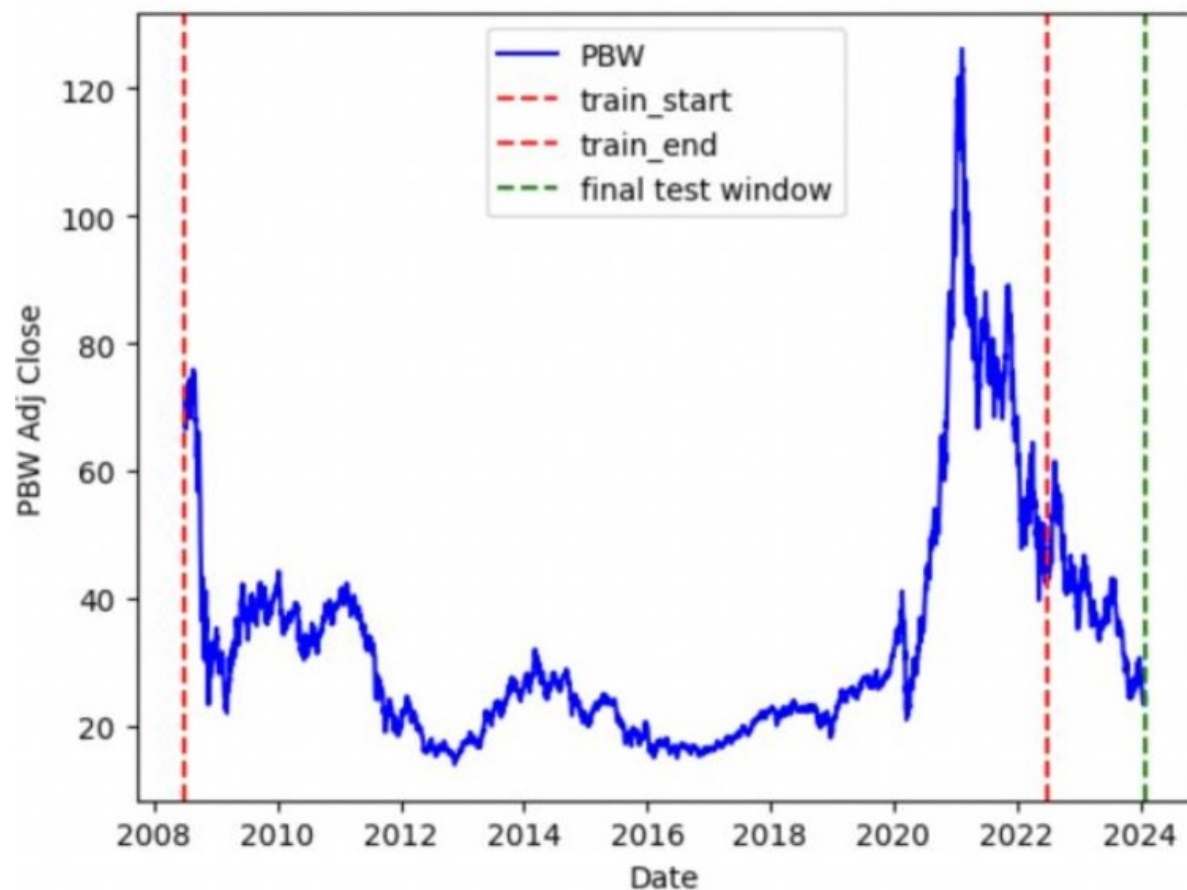


```
plot_result(all_acc[25], 'LR25 EX[1,3]')
```



Conclusion 3 – Final Model Testing with Out-of-Sample testing

Add a final test on out-of-sample data set from **07/01/2022 – 01/02/2024**



04 Conclusion 3 – Final Model Testing

Serve as a **benchmark** for comparison with other models

Naïve Bayes

Gradient Boosting

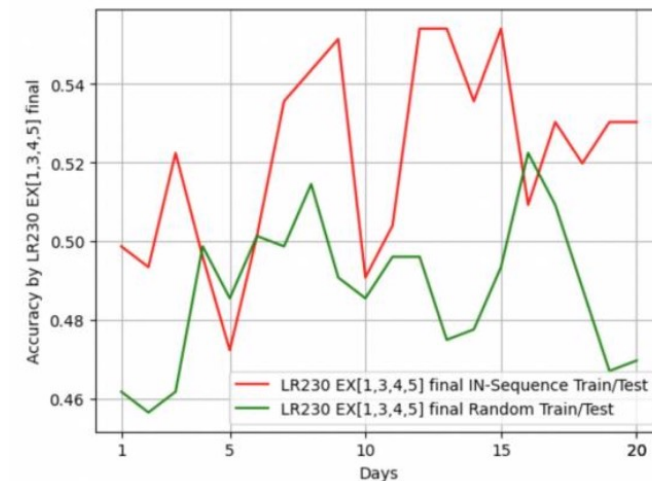
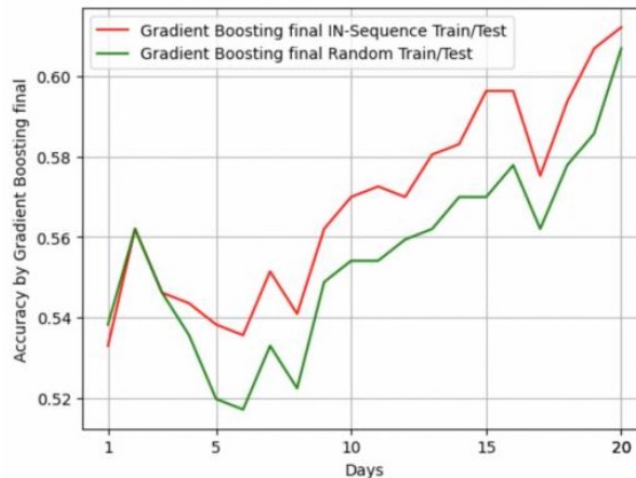
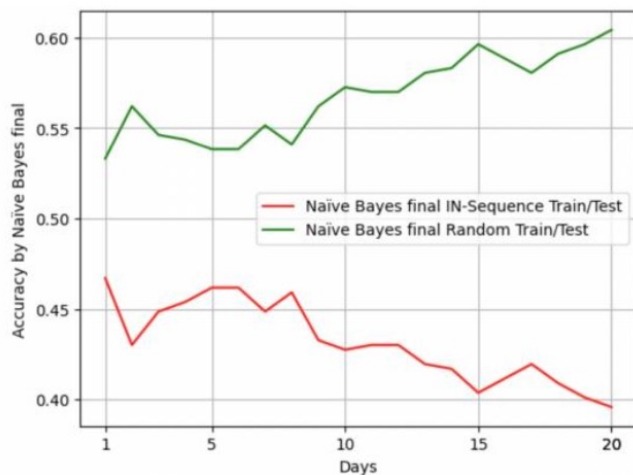
Stand out as top performing ML methods

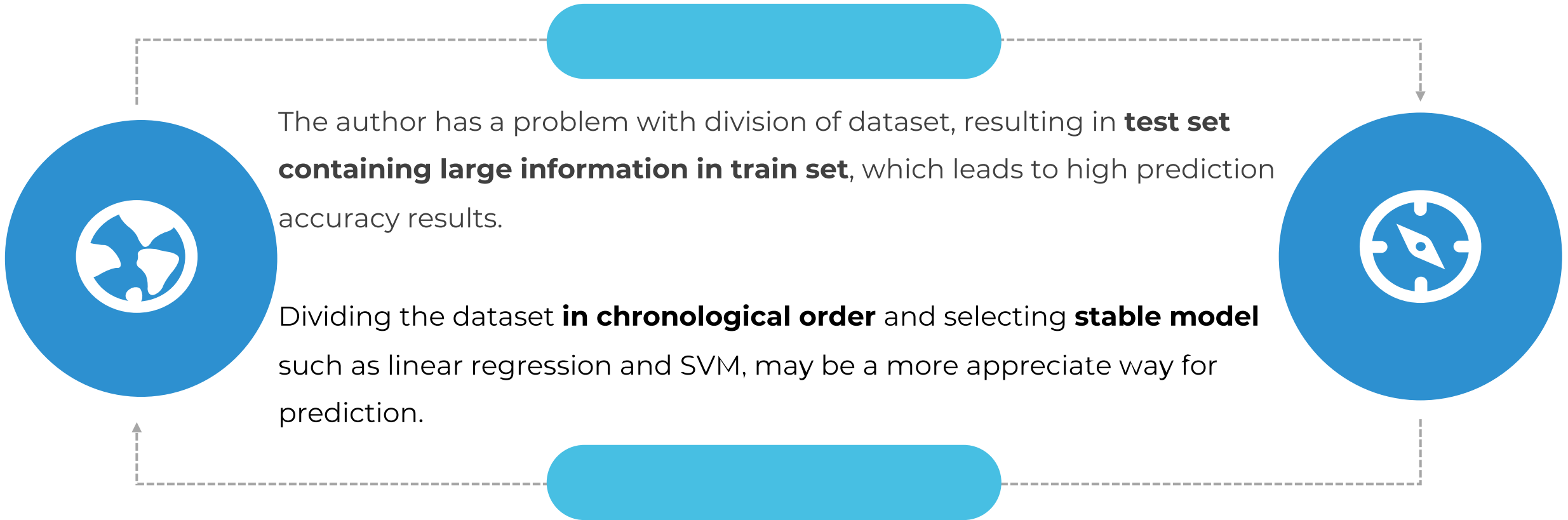
The In-Sample test with random partition gives an 80% accuracy; in-sequence partition gives around 45%.

The Out-of-Sample test with random partition gives an 60% accuracy; in-sequence partition gives around 60%.

The performance **dropped significantly**, and implied **non-consistency**.

The result is better than final model testing, yet this could be due to **market regime change**.





THANK YOU!

