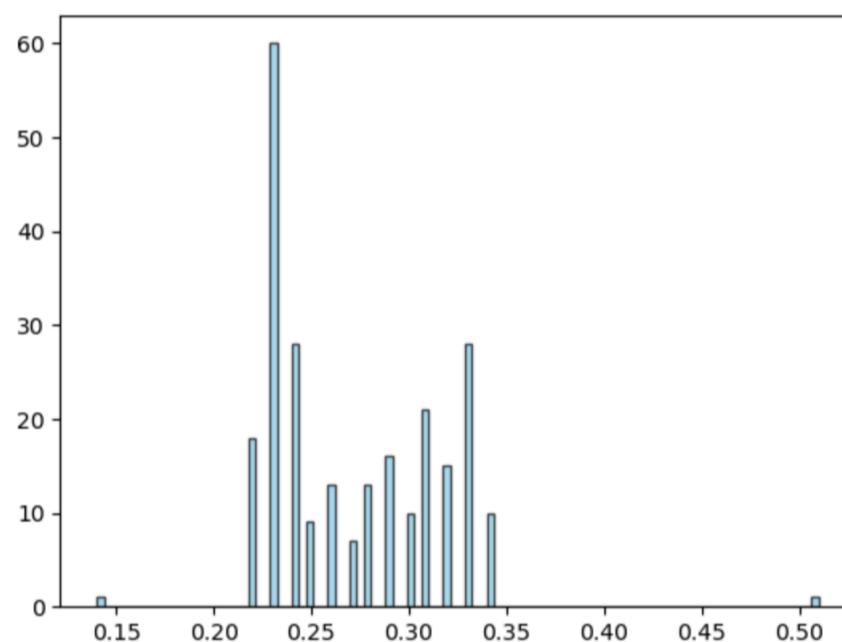


Project Report: Predictive Modeling with Market Dynamics

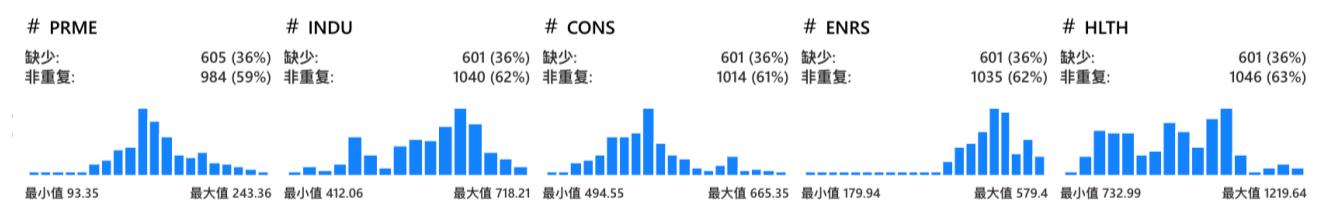
▼ 0. Data Preprocessing

- For all 12 features, there are no values in non-trading days, but LL100 has no missing values. So I assume LL100 is traded in all dates provided.

Although I notice that the weekend data follows approximately a uniform distribution:



The original feature missing values account for 35%-36%. Therefore, for the prediction model, I decided to use only trading days from **2016-01-01** to **2020-04-30** and filter out the missing values for feature input in weekends.



- After filtering out non-trading days, there are 0%-3% missing values for features, and then use forward fill for filling missing data.



For the strategy implementation part, I add a random value from the distribution to the Friday's data to simulate the LL100 weekend prediction value.

- All feature values are missing(except SP500) between **2020-05-01** and **2020-07-22**. Therefore this period is removed in future analysis to avoid introducing inaccuracies in feature inputs.

▼ 1. Non-Stationarity Analysis

▼ Stationarity Testing

- The Augmented Dickey-Fuller (ADF) test was applied to all features and the target variable LL100.
 - Null Hypothesis (H0): The series has a unit root and is non-stationary.
- Result:
 - LL100 is non-stationary, while all other 12 features are stationary.

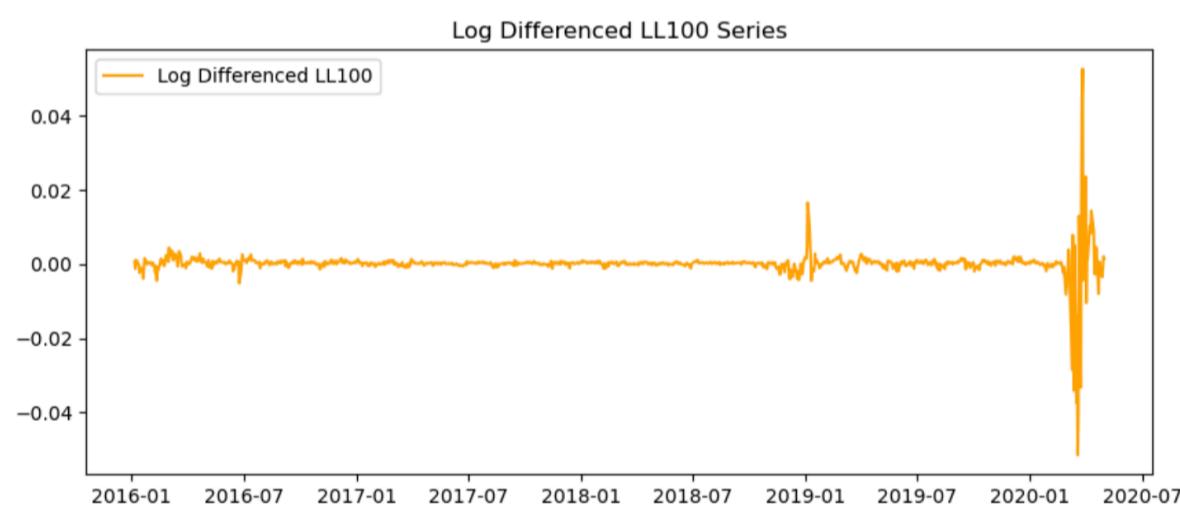
```
Testing stationarity of LL100:  
ADF Statistic: -2.406256531890945  
p-value: 0.13997832993914638  
Critical Value 1%: -3.4364992651202697  
Critical Value 5%: -2.8642551098431968  
Critical Value 10%: -2.5682156239065925  
The series is non-stationary.
```

This rules out the concerns for spurious regression when using OLS for prediction, but the stationary target variable can still cause inefficient coefficient result.

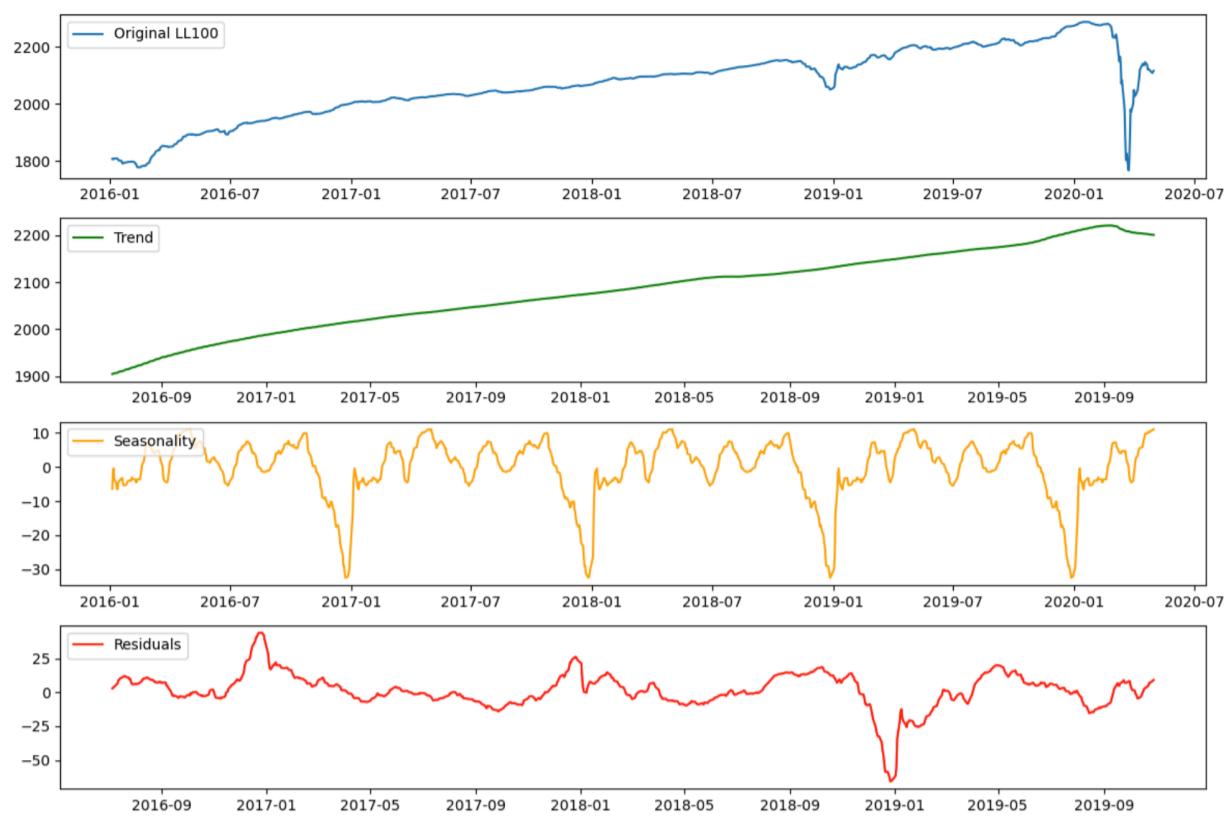
▼ Approaches to Handle Non-Stationarity

1. **Differencing:** Calculated the first difference of LL100, which successfully made the series stationary (ADF p-value < 0.05).
2. **Log Transformation + Differencing:** Log transformation followed by differencing also achieved stationarity.

```
Testing stationarity after log transformation:  
ADF Statistic: -8.88006675875351  
p-value: 1.325336476890943e-14  
Critical Value 1%: -3.436510851955201  
Critical Value 5%: -2.864260220574562  
Critical Value 10%: -2.5682183458999943  
The series is stationary.
```

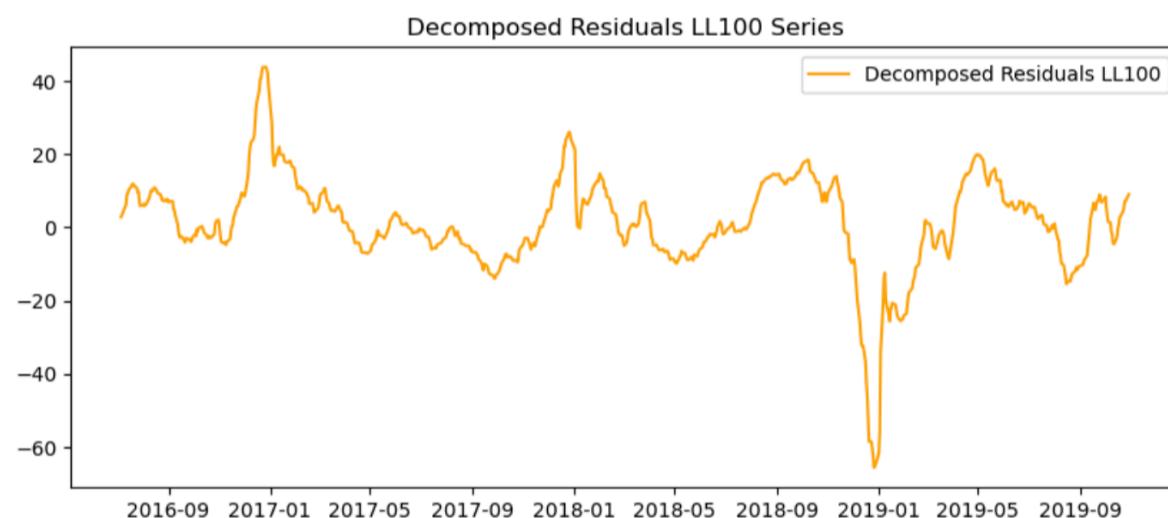


3. **Decomposition:**



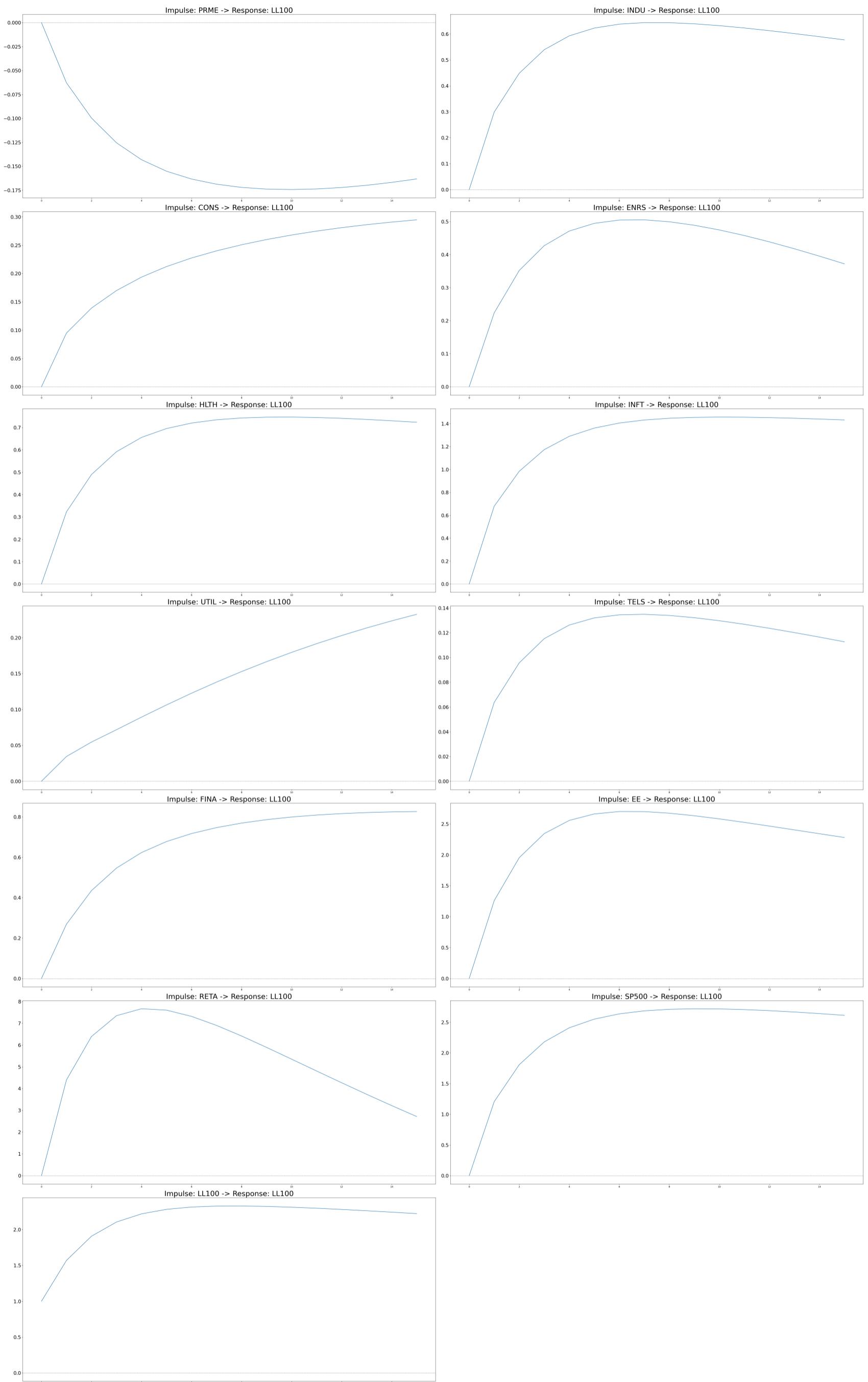
Seasonal decomposition of revealed:

- A clear upward trend: LL100 shows an overall upward trend, but there was a sharp decline from the end of 2019 to the beginning of 2020, possibly triggered by covid market crash.
- Seasonal pattern: LL100 exhibits clear periodic fluctuations, indicating that certain fixed time factors have a significant impact on the data.
- The residuals fluctuate around 0 for most of the time period, indicating that the trend and seasonality have explained the data well.



▼ 2. Causal Relationship Mapping

▼ Correlation Analysis



▼ Results

VAR modeling revealed the following impacts on LL100:

- **Short-term impact:** RETA and EE.
 - Retail index has a significant short term impact (8) on LL100. Emerging Economies index has a short term impact (2.5) on LL100. Retail and emerging markets are more reactive to current economic conditions, making them significant short-term drivers of LL100. These sectors capture immediate changes in consumer behavior and global market sentiment.
- **Long-term impact:** SP500 and INFT.
 - SP500, Information Technology index has a significant long term impact(2.5, 1.4) on LL100. The S&P 500 and the Information Technology index reflect structural and macroeconomic trends that shape the economy over time. Their influence on LL100 underscores the importance of equity markets and technological innovation in driving sustained economic performance.
- **Negative impact:** PRME.
 - Prime Index is the only one with a small negative impact on LL100.

▼ 3. Model Architecture

In this section, I explore econometric, machine learning and state-of-the-art models to predict LL100 index. I created a evaluation metrics with several indicators, aiming to identify which approach better captures market regime changes. I use the first 80% as train dataset, and the last 20% as test dataset.

▼ Evaluation Metrics

see eval.py

- **Mean Squared Error (MSE)**
- **Mean Absolute Error (MAE)**
- **Rolling Window MSE**
 - Use a rolling window of 30 days
- **R-squared (R²)**
- **Explained Variance**
 - Measures the proportion of variance in the target variable that is explained by the model's predictions. Similar to R², but focuses on the variance of the predicted values relative to the actual values.
- **Directional Accuracy**
 - Measures the percentage of times the model correctly predicts the direction of change (up or down) in LL100.

▼ OLS Regression

```
Ljung-Box Test Results:  
    lb_stat  lb_pvalue  
10  6269.362169      0.0  
ADF Test Results:  
ADF Statistic: -3.317074972178901  
p-value: 0.014132263180800748  
Critical Values:  
1%: -3.437923659686726  
5%: -2.8648832361839442  
10%: -2.5685501889710864
```

	OLS_train	OLS_test
MSE	537.92	7199.04
MAE	18.9	51.68
Rolling Window MSE	542.23	3369.62
R-squared (R ²)	0.95	0.09
Explained Variance	0.95	0.39
Directional Accuracy (%)	55.52	60.83

- Residuals were stationary (ADF p-value < 0.05) but showed autocorrelation (Ljung-Box p-value = 0.0), indicating OLS coefficient is consistent but inefficient.
- Performance on the test set is poor, prompting exploration of regularization techniques.

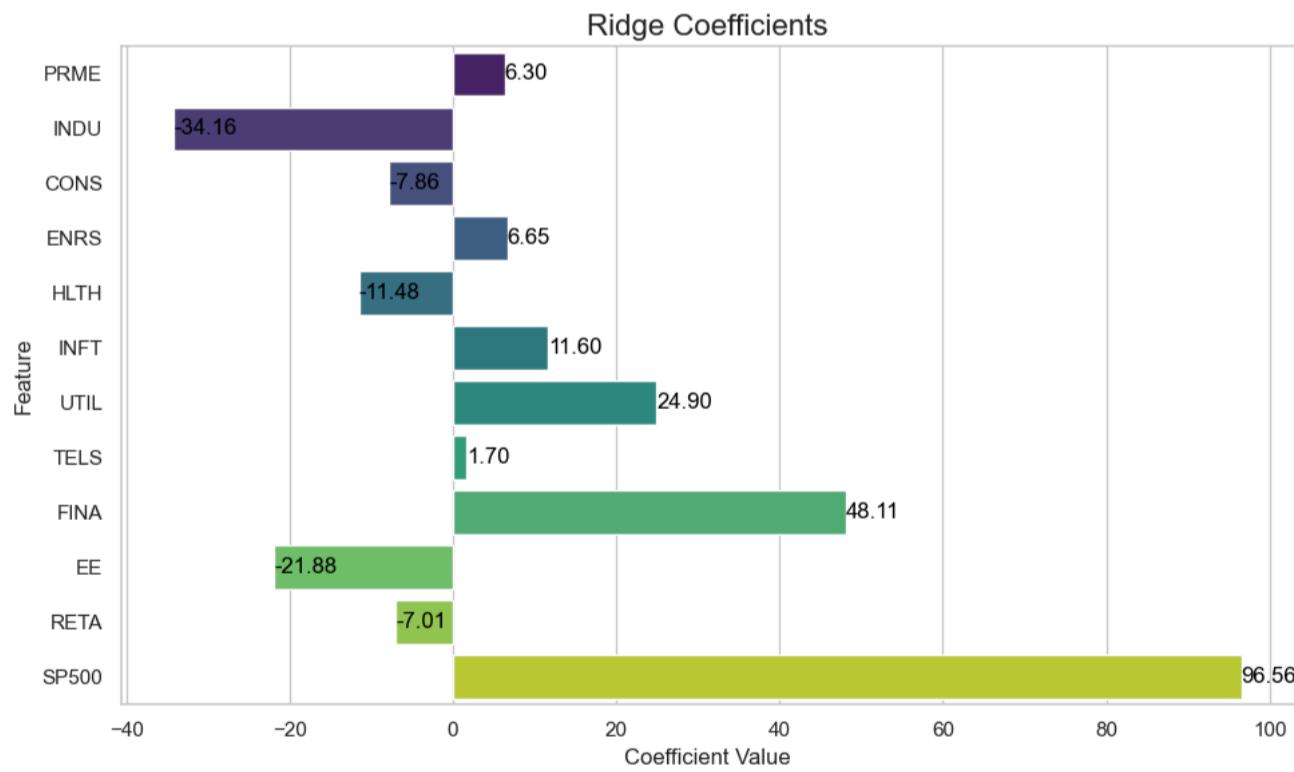
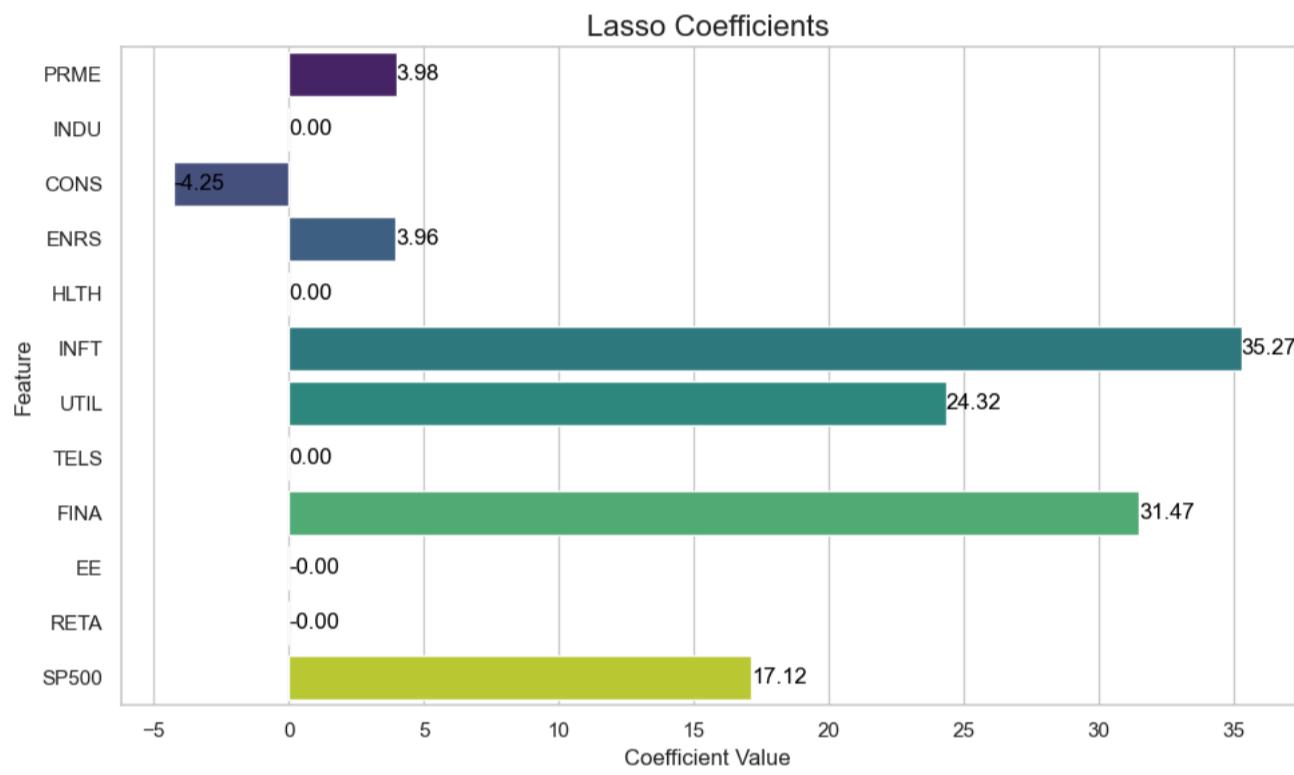
▼ Lasso, Ridge and Elastic Net

- Standardized 12 features as input

- The MSE and Rolling window MSE is large due to the data pattern in 2020 (test set) with a significant market crash. In the case, MAE is probably a better indicator than MSE for model evaluation.

	OLS_train	OLS_test	Lasso_train	Lasso_test	Ridge_train	Ridge_test	Elastic Net_train	Elastic Net_test
MSE	537.92	7199.04	653.59	2975.72	563.61	4892.34	869.4	1793.84
MAE	18.9	51.68	20.78	40.79	19.67	58.88	23.81	28.85
Rolling Window MSE	542.23	3369.62	654.67	3211.12	569.76	5193.78	822.55	1489.32
R-squared (R^2)	0.95	0.09	0.94	0.63	0.95	0.38	0.92	0.77
Explained Variance	0.95	0.39	0.94	0.76	0.95	0.82	0.92	0.8
Directional Accuracy (%)	55.52	60.83	56.21	63.13	57.01	62.67	56.9	61.75

- Elastic Net Achieved the best performance among linear models with $R^2 = 0.77$ and Explained Variance Score = 0.80 on the test set. Since it combines the strengths of both Lasso and Ridge regression. The regularization parameter alpha controls the trade-off between the L1 and L2 penalty terms.



- Expanding Window

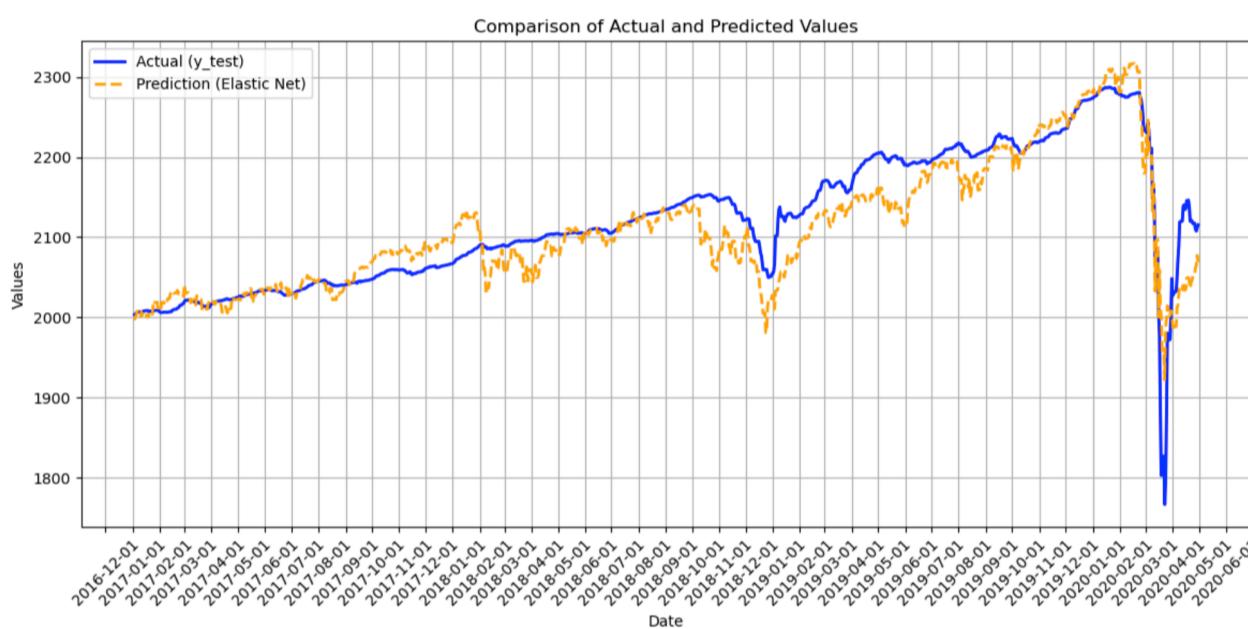
Start with data in `2016` as train set to predict the LL100 of `2017-01-01`, expanding the train set for daily LL100 prediction. Then shift the prediction for one day to avoid using future data.

- Use the distribution of `ll100_diff` to add values to fill NaN values for weekends, accounting for weekend transaction for LL100.

Performance Metrics for Elastic Net (test)

Elastic Net_test

MSE	1089.40
MAE	23.87
Rolling Window MSE	1048.68
R-squared (R^2)	0.84
Explained Variance	0.85
Directional Accuracy (%)	69.91



▼ Momentum-based Strategy

The trading strategy uses a rolling prediction signal to identify potential trends in the market and generates buy or sell signals based on the momentum of the predictions. The strategy incorporates risk management techniques, including `trailing stop-loss`, `take-profit levels`, and a `maximum holding period`. All trades are executed using the entire available capital, and the strategy is benchmarked against a simple buy-and-hold approach.

▼ Trading Logic

Momentum Signal:

The signal is calculated as the difference between the current and previous rolling predictions:

$$\text{Signal} = \text{Rolling Prediction}_t - \text{Rolling Prediction}_{t-1}$$

- If the signal > 0 , it indicates upward momentum, generating a buy signal.
- If the signal < 0 , it indicates downward momentum, generating a sell signal
- If the signal $= 0$, no trade executed.

Risk Management:

The strategy exits a position based on the following conditions:

- Trailing Stop-Loss:
 - For a long position, if the price falls below the peak price by more than the specified `stop_loss_trailing` percentage, the position is exited.
 - For a short position, if the price rises above the peak price by more than the specified `stop_loss_trailing` percentage, the position is exited.

- Take-Profit:
 - For a long position, if the price rises above the entry price by the specified `take_profit` percentage, the position is exited.
 - For a short position, if the price falls below the entry price by the specified `take_profit` percentage, the position is exited.

Capital Allocation

At the time of entry, the entire available capital is allocated to the position. The position size is calculated as:

$$\text{Position Size} = \frac{\text{current capital}}{\text{entry price}}$$

▼ Parameters

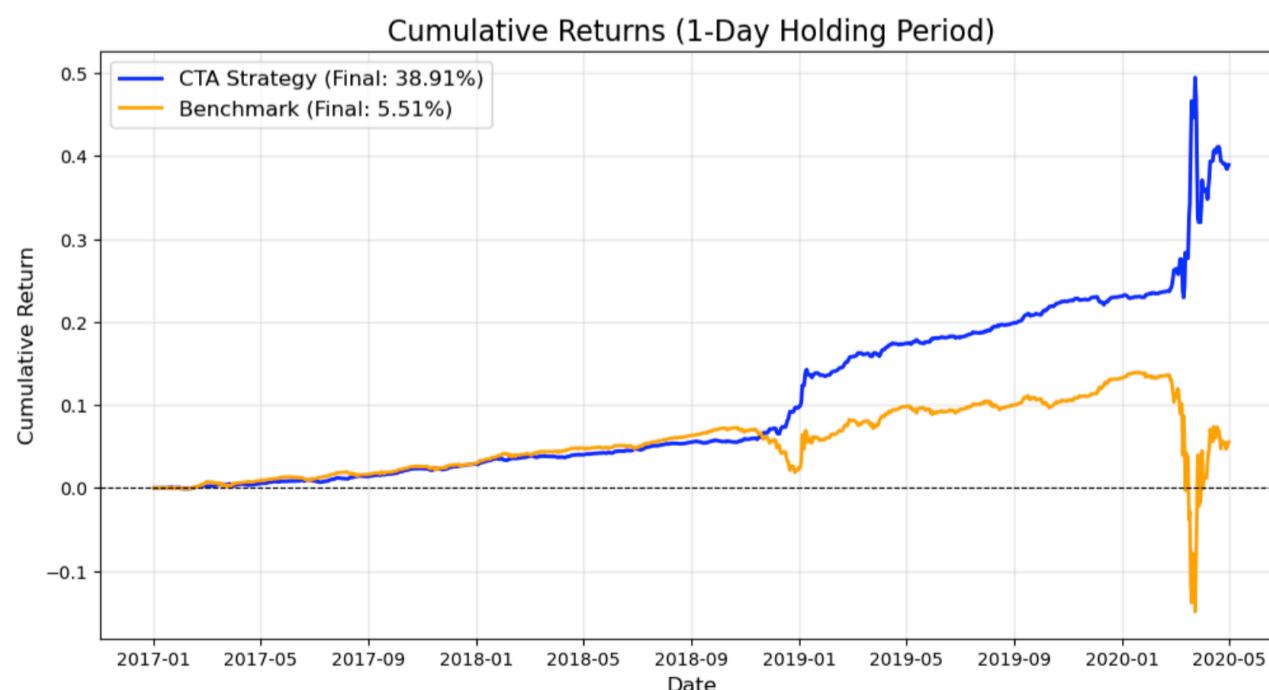
1. Initial Capital (`initial_capital=100000`):
2. Rolling Window (`rolling_window=10`):
 - Used to calculate the average of the LL100 predictions.
3. Trailing Stop-Loss (`stop_loss_trailing=0.05`):
 - Exits a trade if the price moves against the position by 5% from the peak price.
4. Take-Profit (`take_profit=0.1`):
 - Exits a trade when the price moves in favor of the position by 10%
5. Maximum Holding Period (`holding_period`):
 - The maximum number of days a position can remain open.
6. Transaction Cost (`fee_rate=0.1%`)
 - Assume transaction cost is 0.1% of the capital at each trade

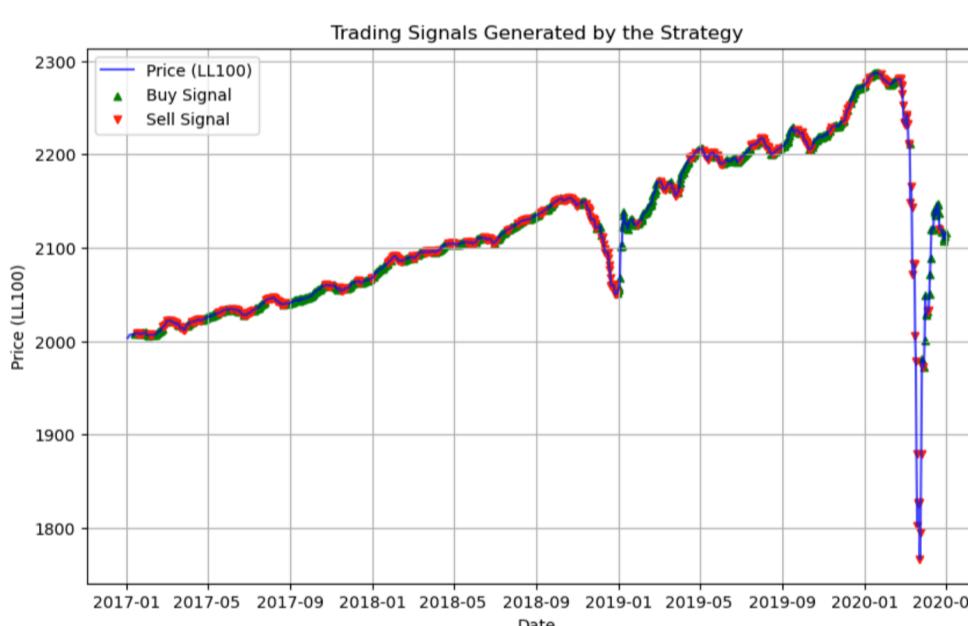
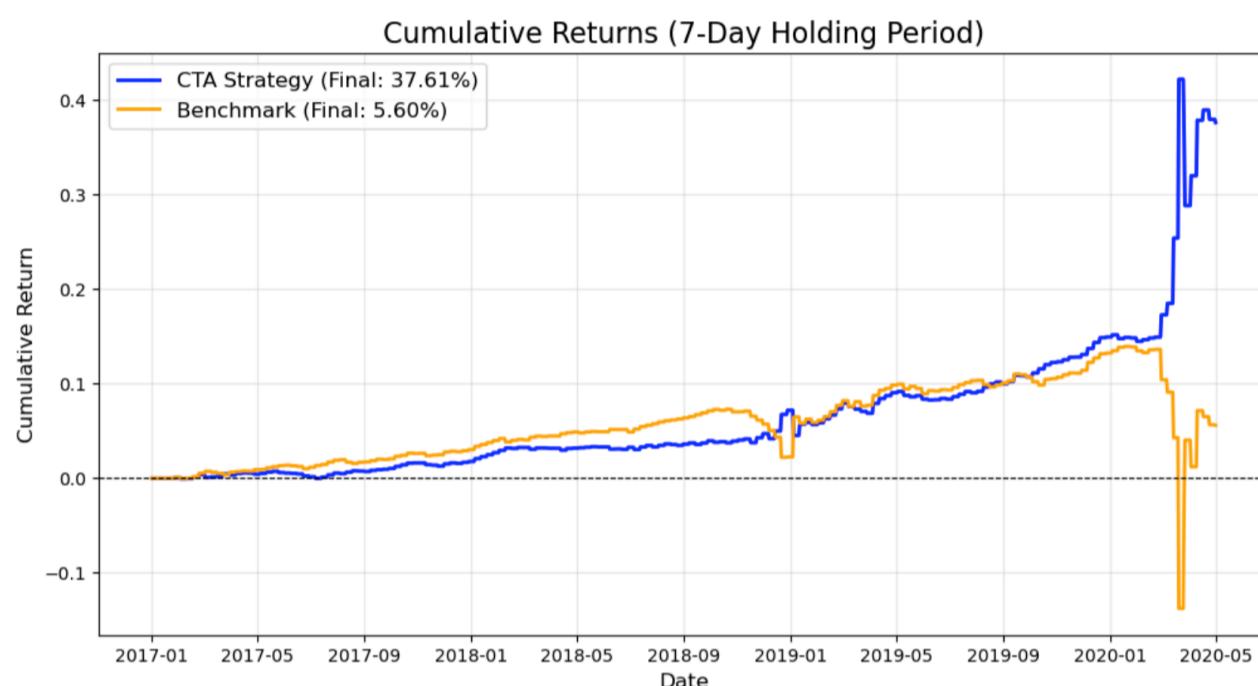
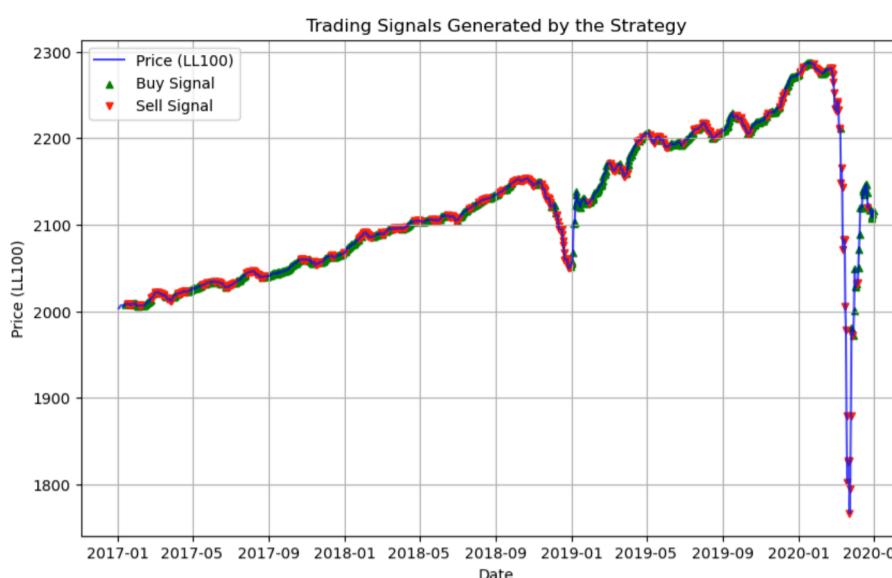
▼ Backtest Result

see `backtest.py`, Benchmark is a simple buy-and-hold strategy.

▼ 2017-01 - 2020-04 :

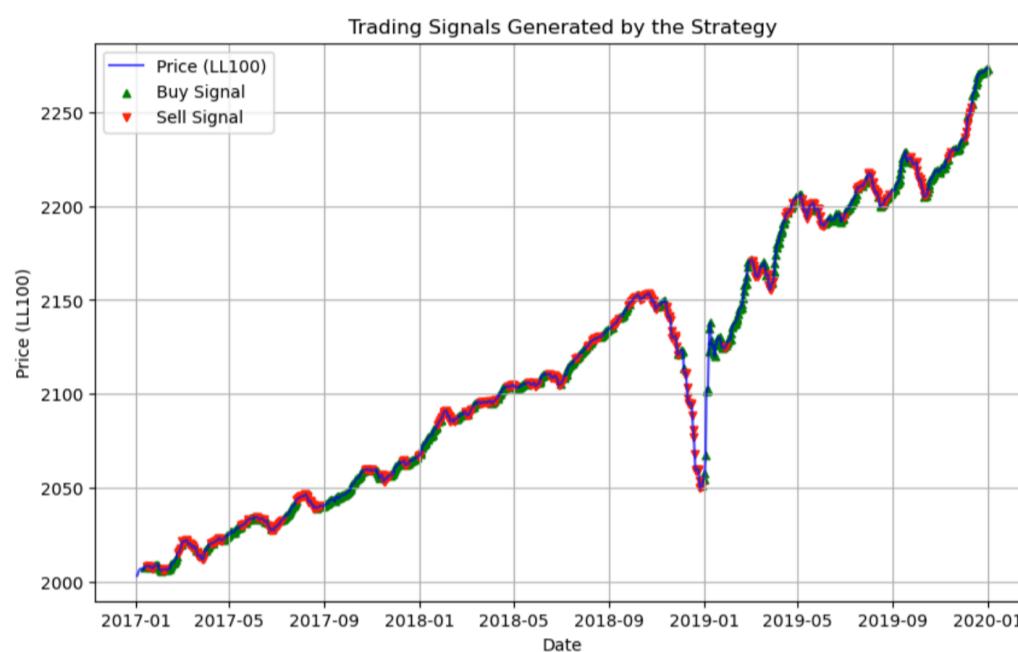
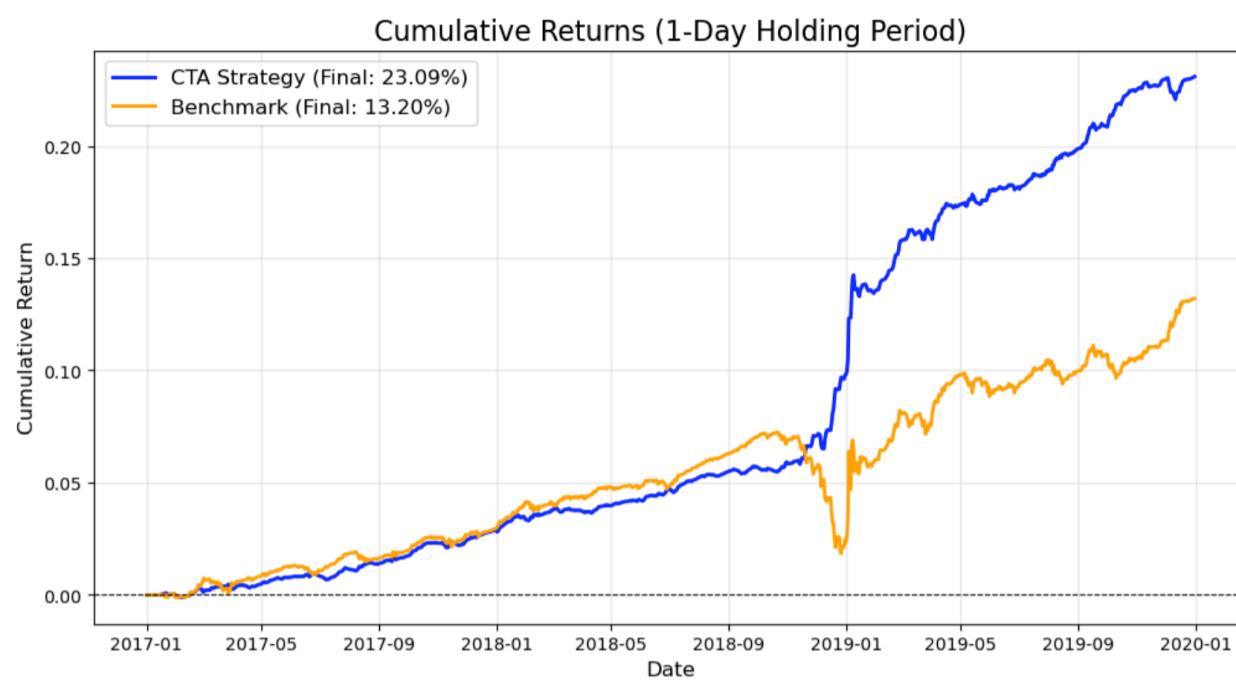
201701-202004	MOM_1d	MOM_7d	Benchmark
Annualized Return	0.1065	0.1055	0.0265
Sharpe Ratio	1.1351	0.8302	0.0445
Max Drawdown	-0.1168	-0.0940	-0.2432
Volatility	0.0715	0.0967	0.1374





▼ 2017-01 - 2019-12 :

201701-201912	MOM_1d	Benchmark
Annualized Return	0.0718	0.0427
Sharpe Ratio	2.9833	0.7039
Max Drawdown	-0.0084	-0.0503
Volatility	0.0166	0.0310



▼ Limitations

1. Sideways Markets: The strategy may underperform in range-bound or choppy market conditions, where trends are weak or nonexistent.
2. Parameter Sensitivity: The performance of the strategy is sensitive to the choice of parameters (e.g., rolling window size, stop-loss percentage).
3. Transaction Cost and Slippage: While I set a fixed transaction cost rate, real-world markets often experience slippage, especially during periods of high volatility or low liquidity. This can further impact net returns and should be accounted for in the backtest.

Conclusion

- The MOM_1d strategy consistently outperformed the benchmark, achieving higher annualized returns, Sharpe ratios, and lower drawdowns and volatility.
- During 2017-2019, MOM_1d demonstrated exceptional risk-adjusted performance with a Sharpe ratio of 2.9833, far surpassing the benchmark's Sharpe ratio (0.7) with low drawdown.

▼ 5. Model Improvements

- Additional Data Sources
 - Macroeconomic indicators (e.g., GDP, interest rates, inflation).
 - Stock prices related to LL100 index.
 - Commodity prices related to LL100 index.
 - Alternative data (e.g., news sentiment, social media trends).
- Use larger datasets with longer time periods to train state-of-the-art models.
- Incorporate cointegration analysis to inform pairs-trading strategy.

▼ Cointegration Analysis

```
''' Cointegration Test '''
from statsmodels.tsa.vector_ar.vecm import coint_johansen

variables = data[feature_names]
result = coint_johansen(variables, det_order=0, k_ar_diff=1)

...
Johansen Cointegration Test
H0: r=0 ; r<=1; r<=2; ...; r<=11
...
print(result.lr1) # Trace statistic - r (number of cointegration)
print(result.cvt) # Critical values - [90%, 95%, 99%]
```

```
[3.95189408e+02 2.78529126e+02 2.09988919e+02 1.63554260e+02
 1.22907502e+02 8.89812094e+01 6.18562227e+01 4.04860355e+01
 2.44366853e+01 1.12517888e+01 4.74811961e+00 3.54398118e-01]
[[326.5354 334.9795 351.215]
 [277.374 285.1402 300.2821]
 [232.103 239.2468 253.2526]
 [190.8714 197.3772 210.0366]
 [153.6341 159.529 171.0905]
 [120.3673 125.6185 135.9825]
 [ 91.109 95.7542 104.9637]
 [ 65.8202 69.8189 77.8202]
 [ 44.4929 47.8545 54.6815]
 [ 27.0669 29.7961 35.4628]
 [ 13.4294 15.4943 19.9349]
 [ 2.7055 3.8415 6.6349]]
```

From the Johansen Cointegration test above, 8898<190.8714, indicating there're at least 5 cointegration relationship in 12 features, suggesting possible pairs trading opportunities.

▼ Z-Score Strategy

Below use the entire dataset 2016-2020 to determine the cointegration coefficient for pair-trading strategy, which is no realistic in real life. Suppose there are previous data from before 2016, we can use the same strategy to avoid any future data leakage.

The z-score-based strategy described in the code employs a statistical arbitrage approach, focusing on trading three assets: **LL100**, **EE**, and **INFT**. The strategy is designed to exploit mean-reversion behavior in the residuals of these assets. Below is a summary of the trading logic.

▼ Trading Logic

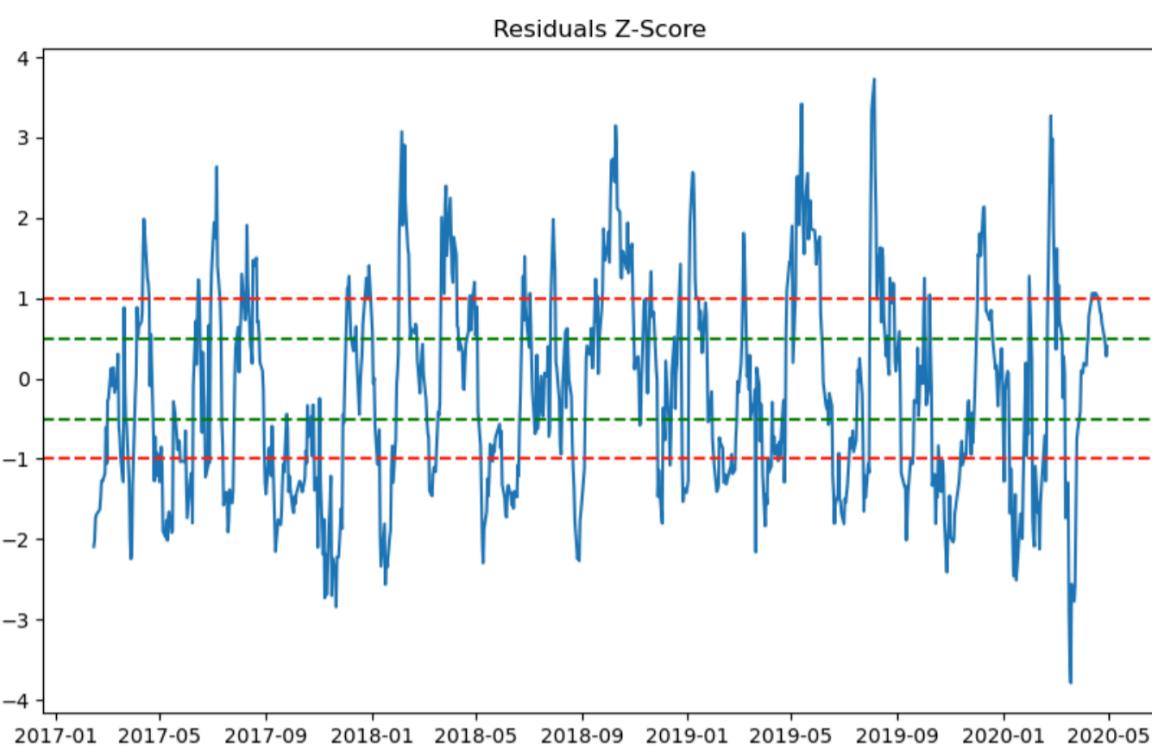
- Found significant cointegration relationships of **LL100** with **INFT** and **EE**.

ADF Statistic for residuals: -2.9543750819242813
p-value for residuals: 0.039380490868017495

- **Z-Score Calculation:**

$$z_score = \frac{\text{residuals} - \text{rolling mean}}{\text{rolling standard deviation}}$$

where window is 30



- **Entry Conditions:**

- **`z_score > 1:`**

- Short 1 unit of LL100. Long on 0.2 units of EE and 0.2 units of INFT (proportional to the available capital).

- **`z_score < -1:`**

- Long 1 unit of LL100. Short 0.2 units of EE and 0.2 units of INFT.

- **Exit Condition:**

- **`0.5 < z_score < 0.5:`**

- Close the position when the z-score reverts to the range. At this point, calculate the pnl for all three assets. Add the total pnl to the capital.

▼ **Backtest Result**

201701-202004	Z-score_Strat	Benchmark
Annualized Return	0.0547	0.0265
Sharpe Ratio	0.7585	0.0445
Max Drawdown	-0.0296	-0.2432
Volatility	0.0438	0.1374

