

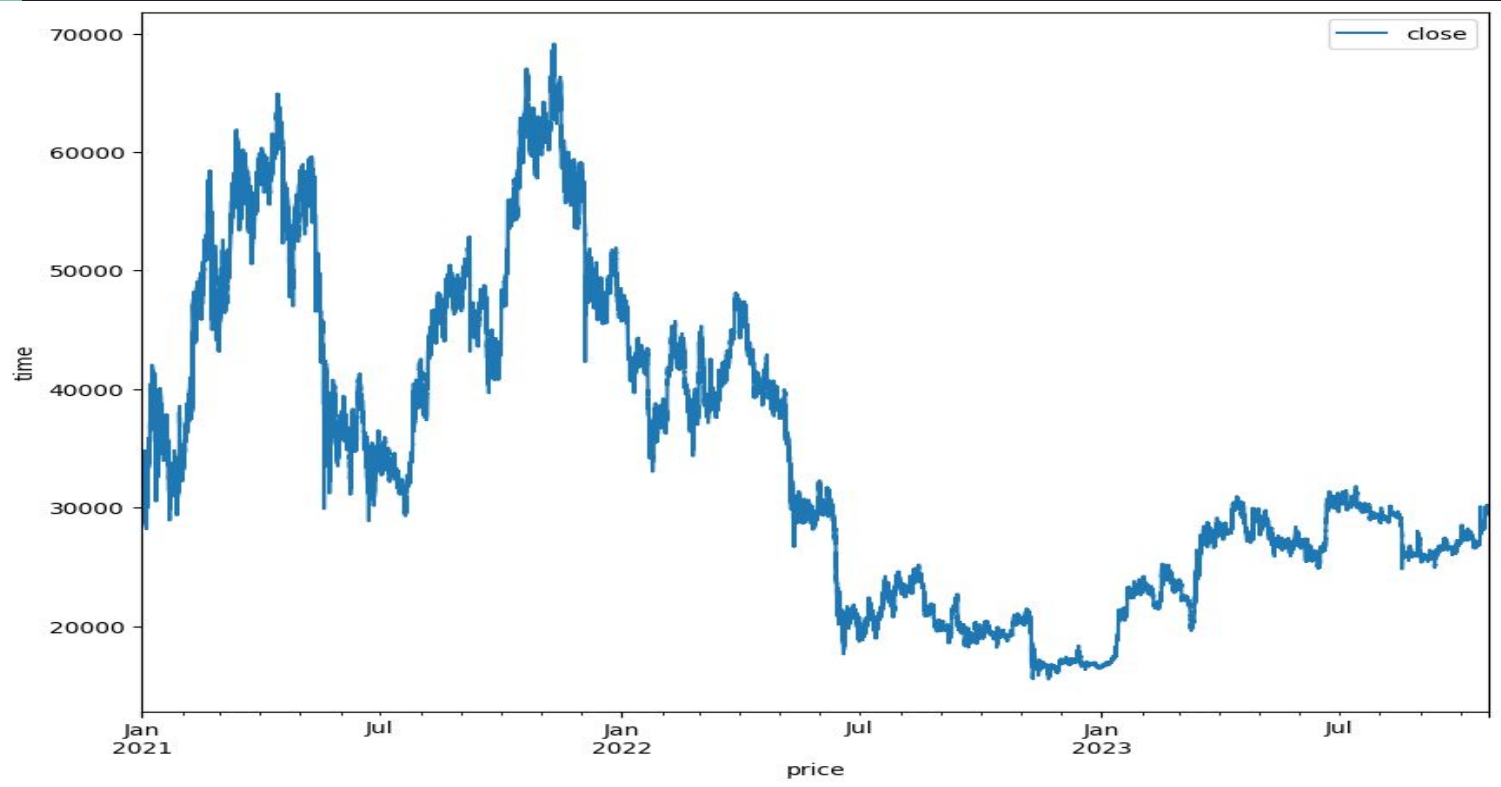


Using Price and Volume Data-based to Forecast Short-term Return in Bitcoin

by James Chiu, Johnson Hsiao
12, 22, 2023

Q: Why doing this topic?

A: Make money easily with stats.





Preview and abstract

Base(Benchmark) → Long a unit every minutes and close 1 hour later.

1. Win rate = win / total trades = 50.378%

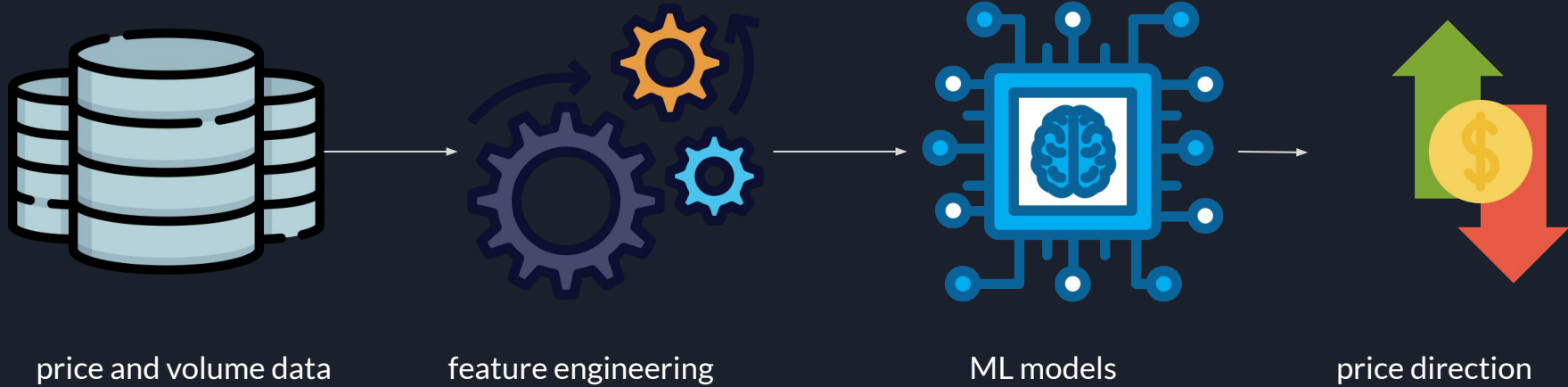
→ It makes sense because the market is the approximate 0-sum game.

2. Profit Factor = 1.014251,

→ It means you can only earn 1.01425 dollar by losing 1 dollar in 3 years

Hence, our target is predict the correct direction of hourly return to improve the win rate and PF.

Study Plan





DATA

Original

Open, High, Low, Close

Indicator and Algorithm

rsi = momentum

amount_spread = volume difference

bar_rtn_sum = momentum

volatility_0 = volatility

price_vol_corr = corr (price and volume)

zscore = rolling price change zscore

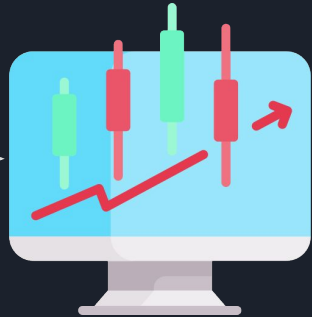
x_s_ratio = vwap / avg(rtn)

high_low_dis = high and low distance rolling sum

Feature engineering Method 1



raw OHLCV data

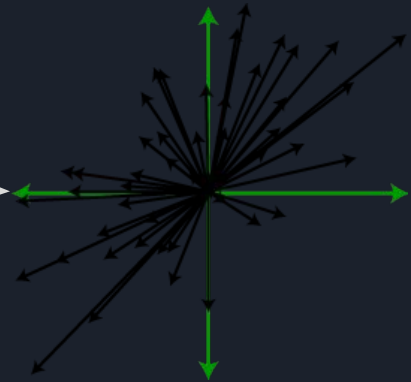


technical indicators

long

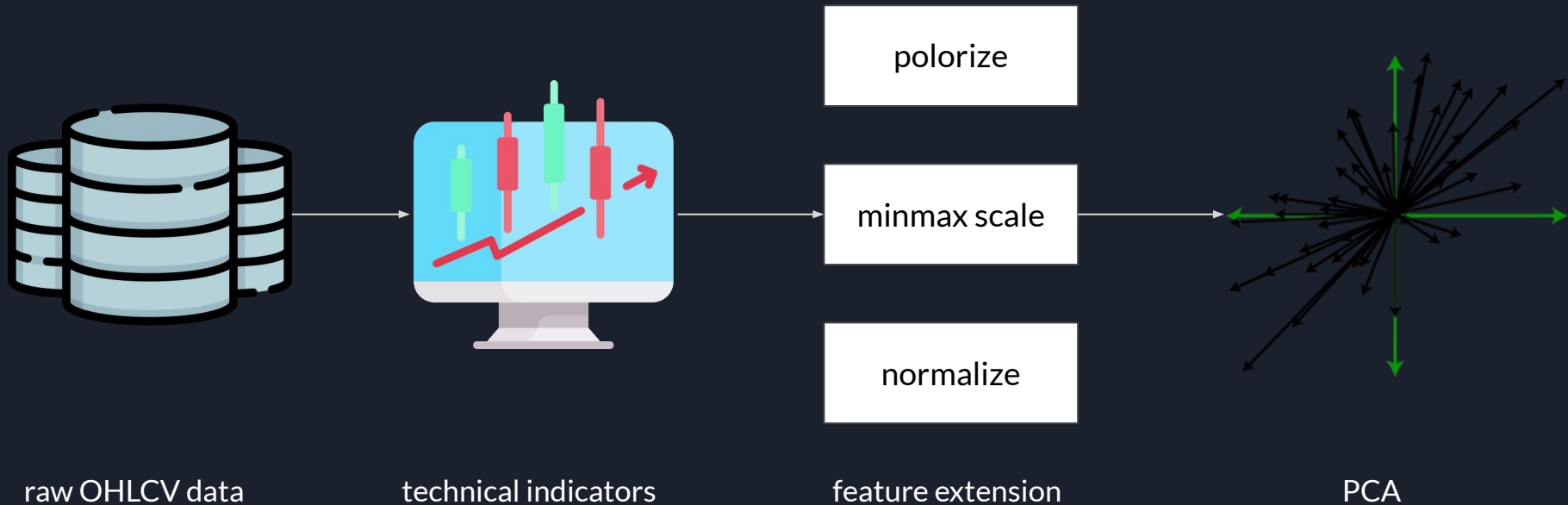
short

Strategies

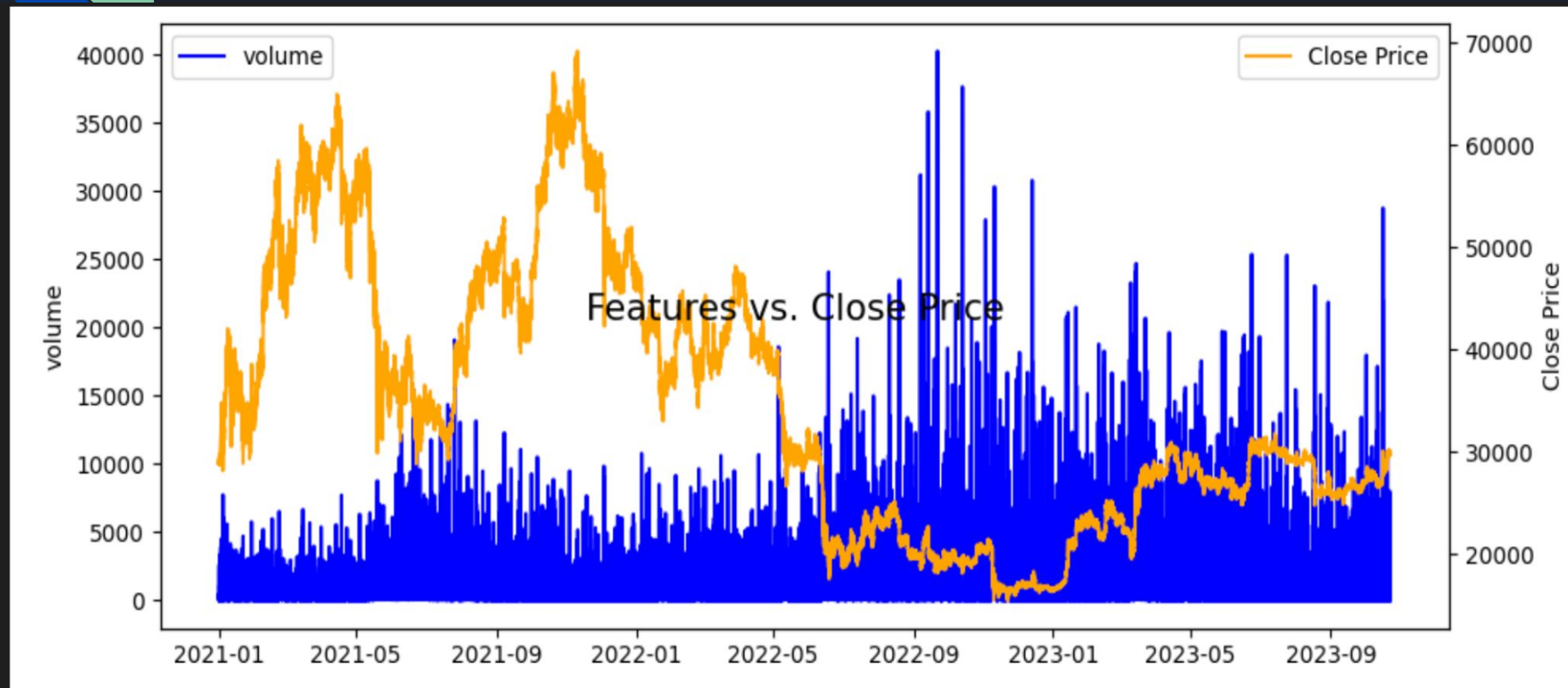


PCA

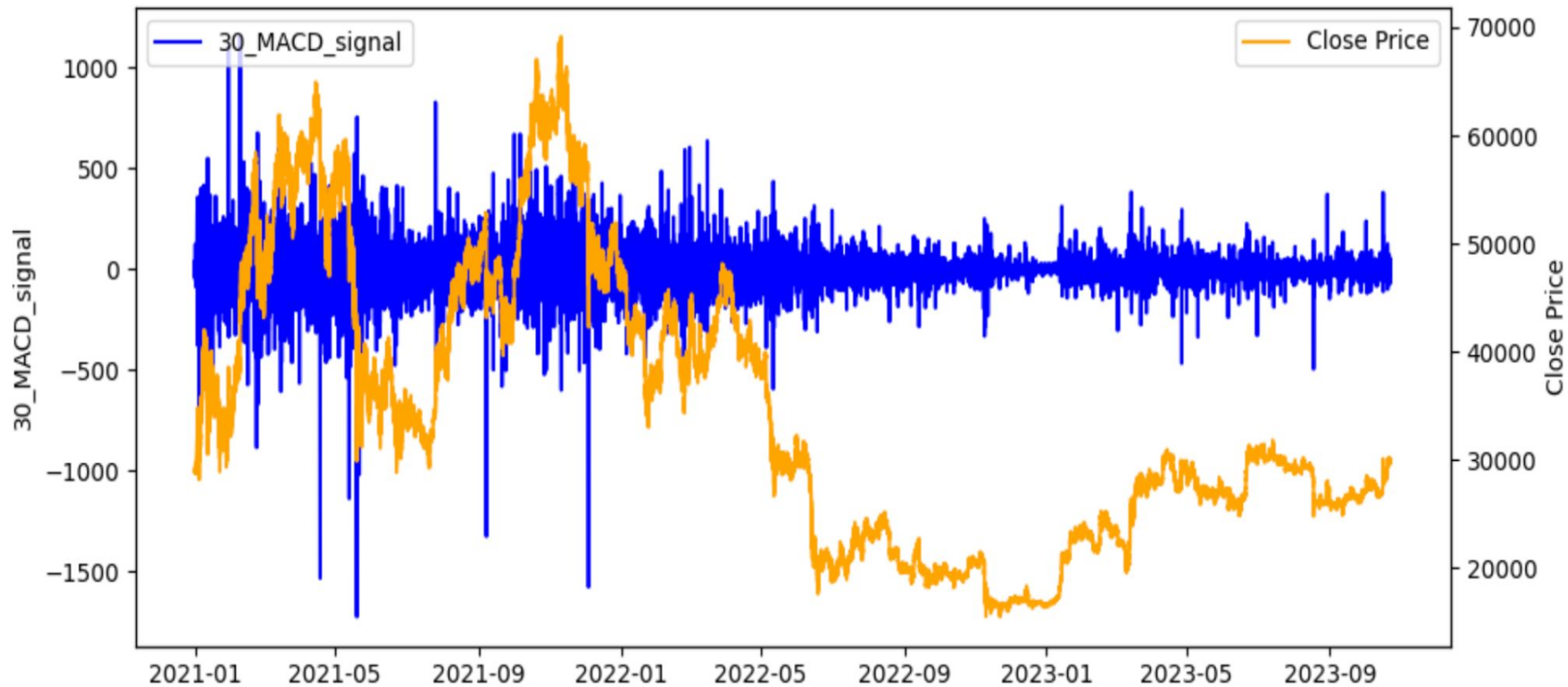
Feature engineering Method 2



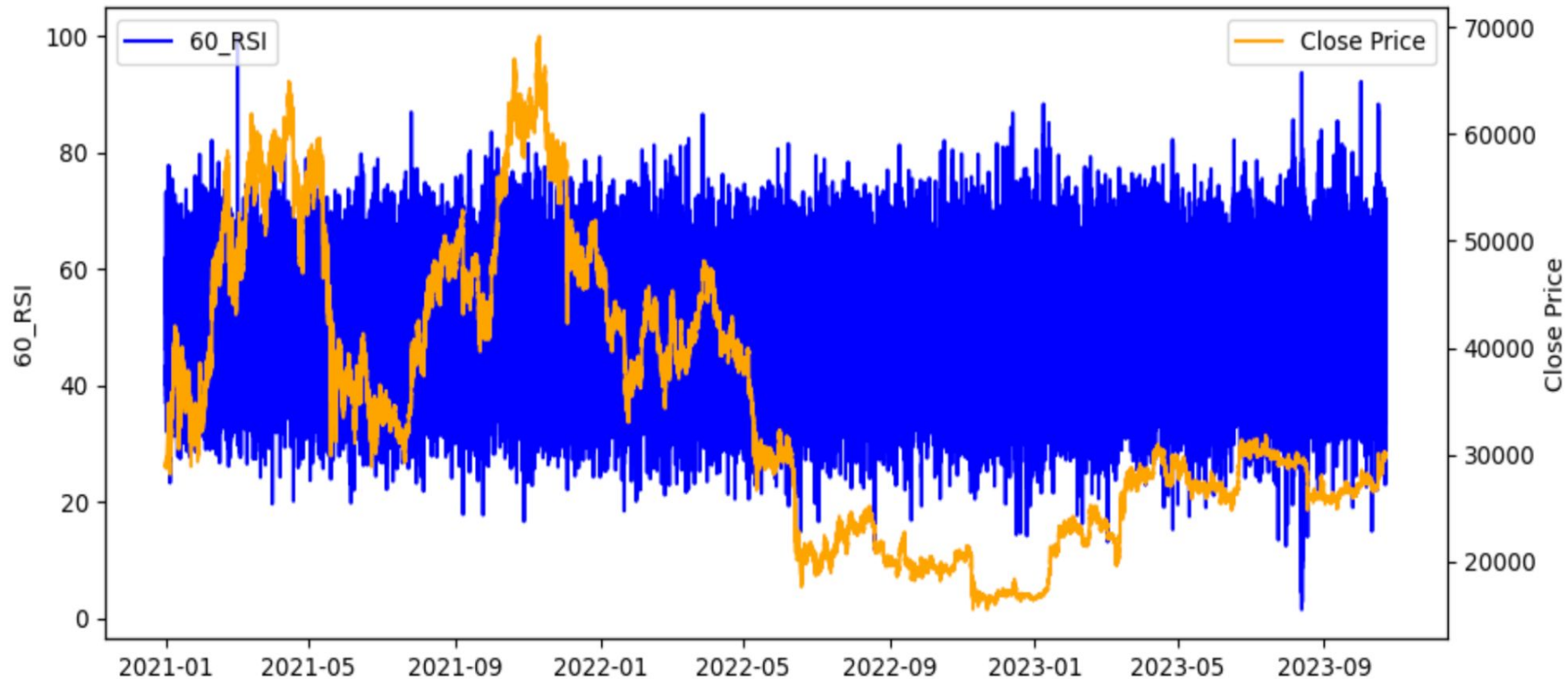
Overview



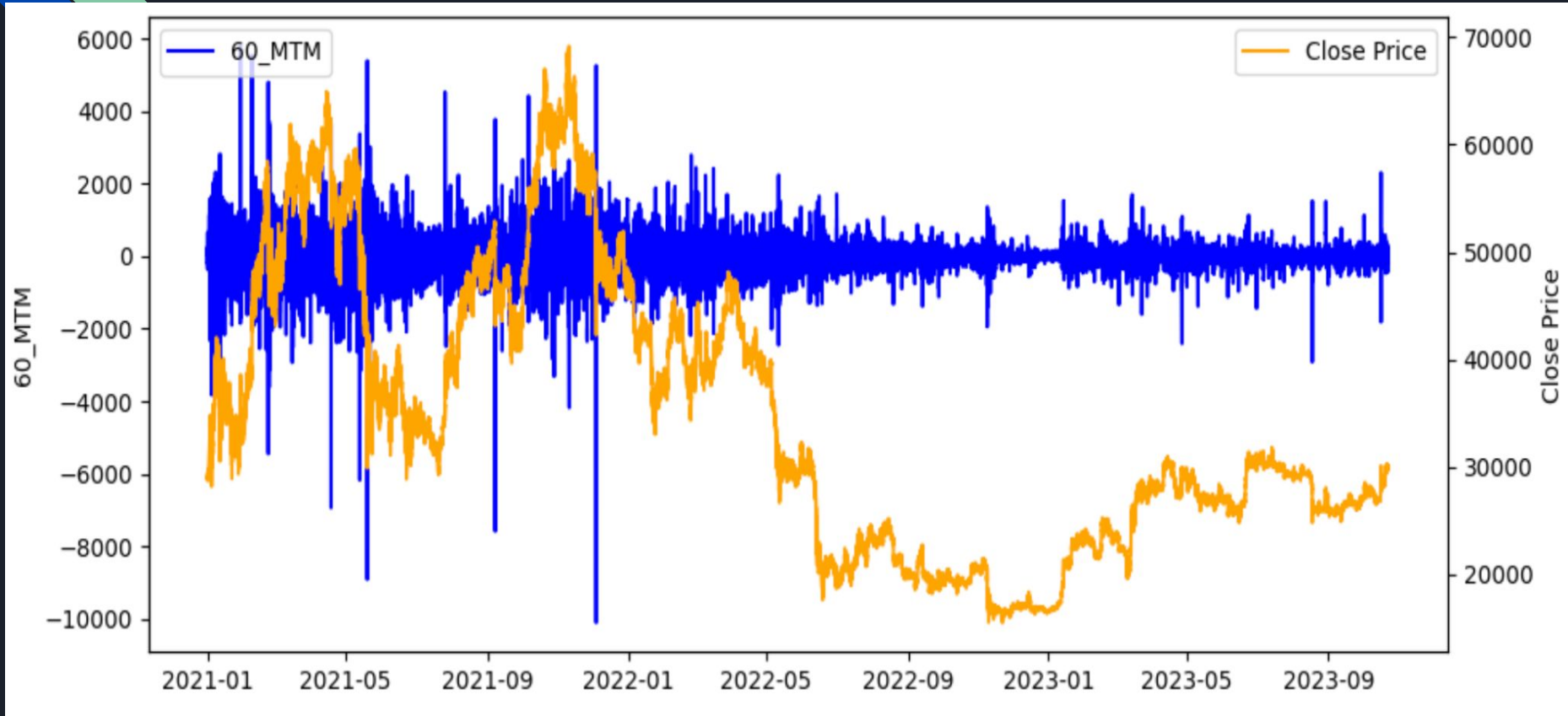
Overview



Overview



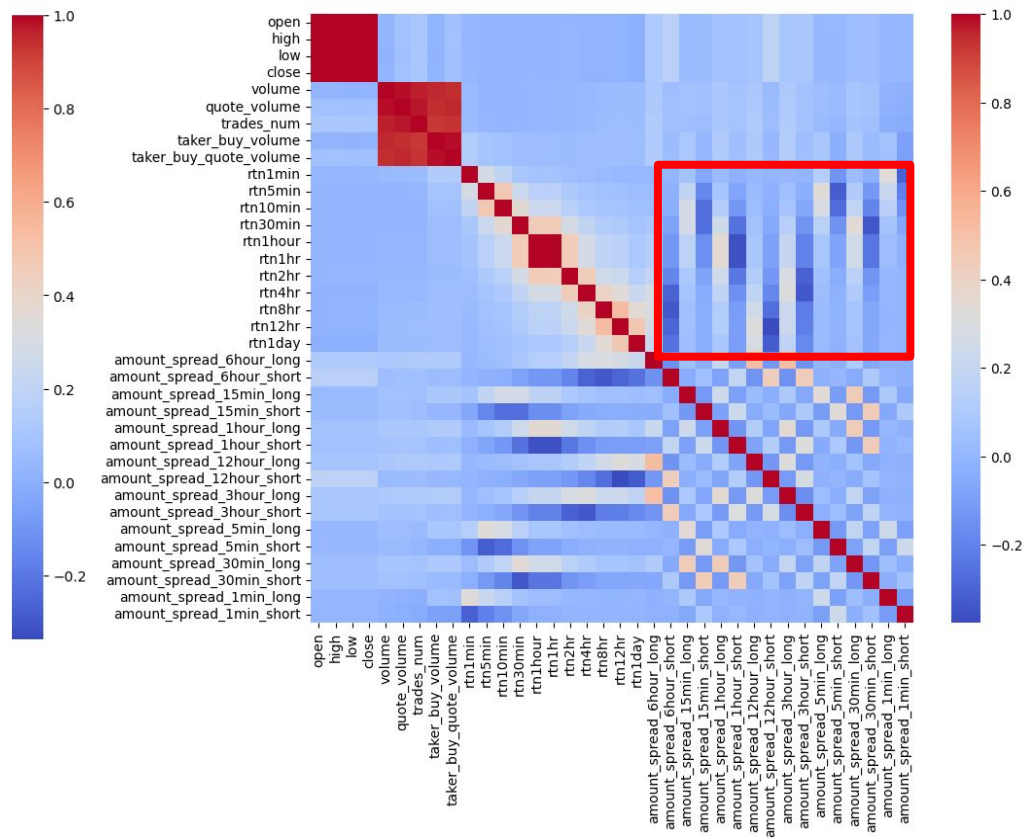
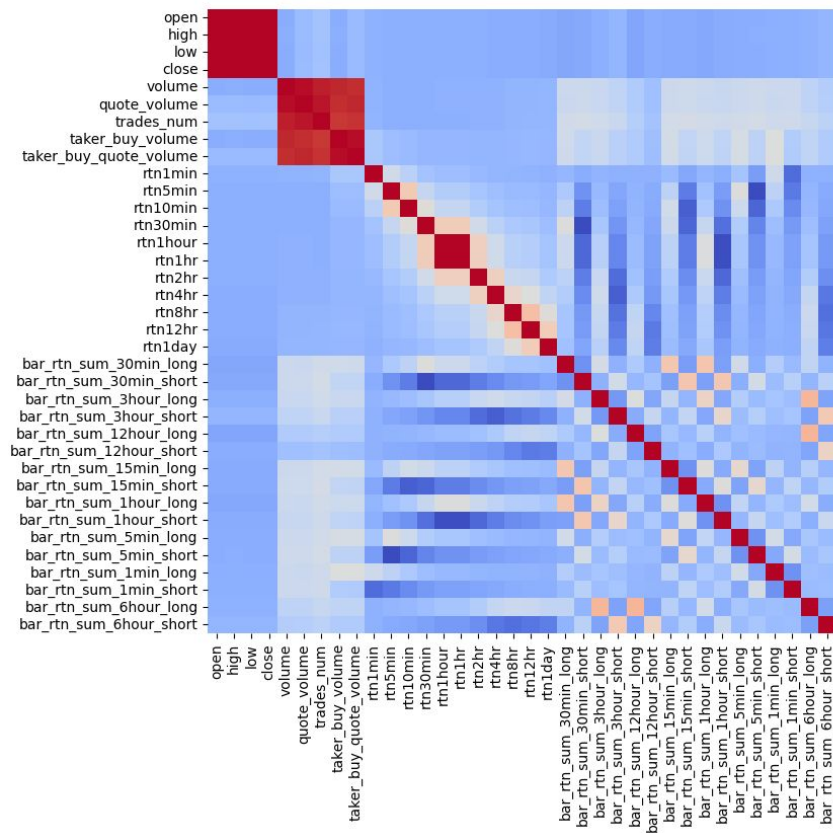
Overview



Feature

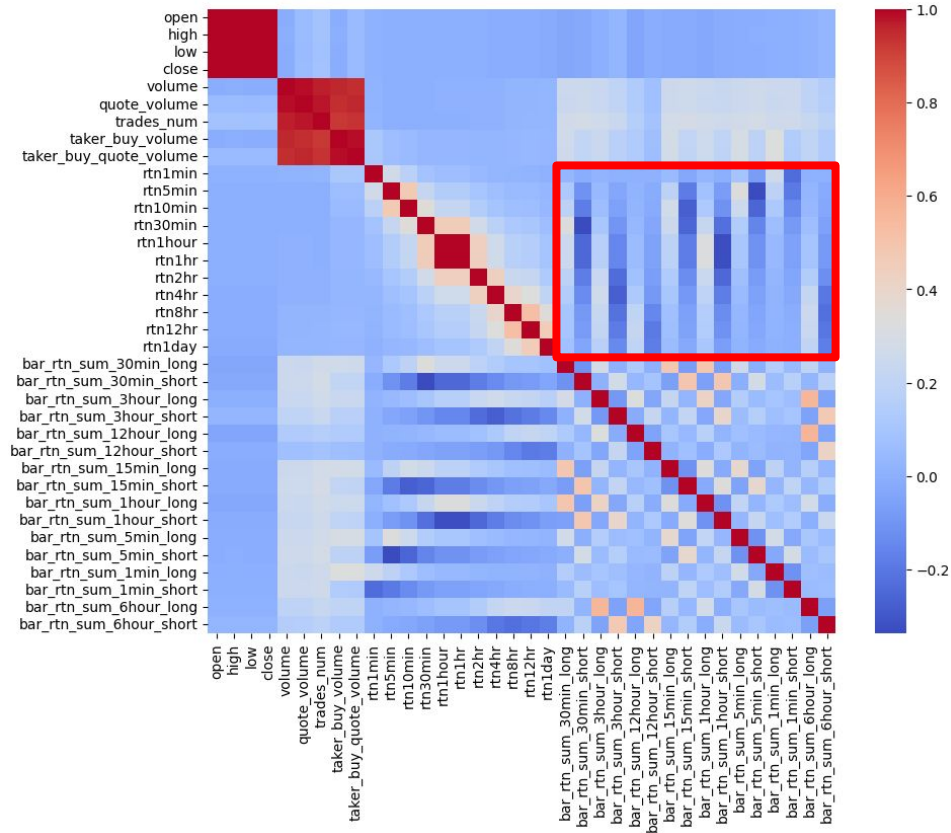
RSI

amount_spread

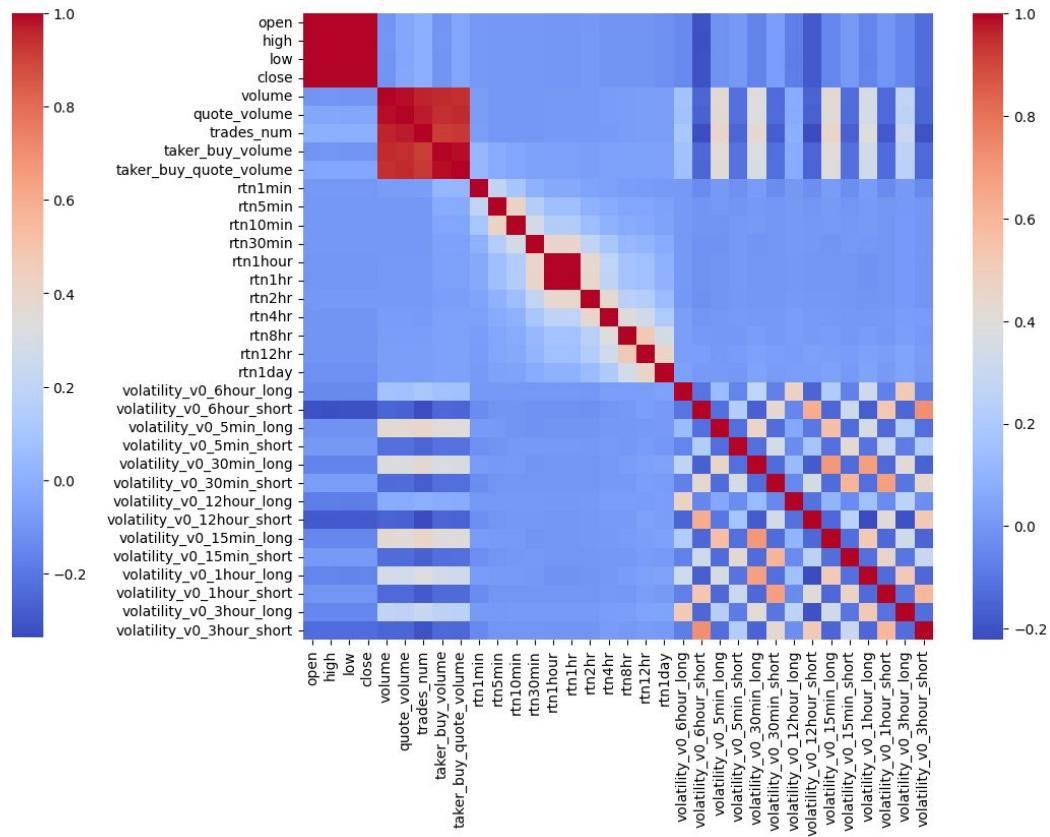


Feature

Kbar_rtn_sum

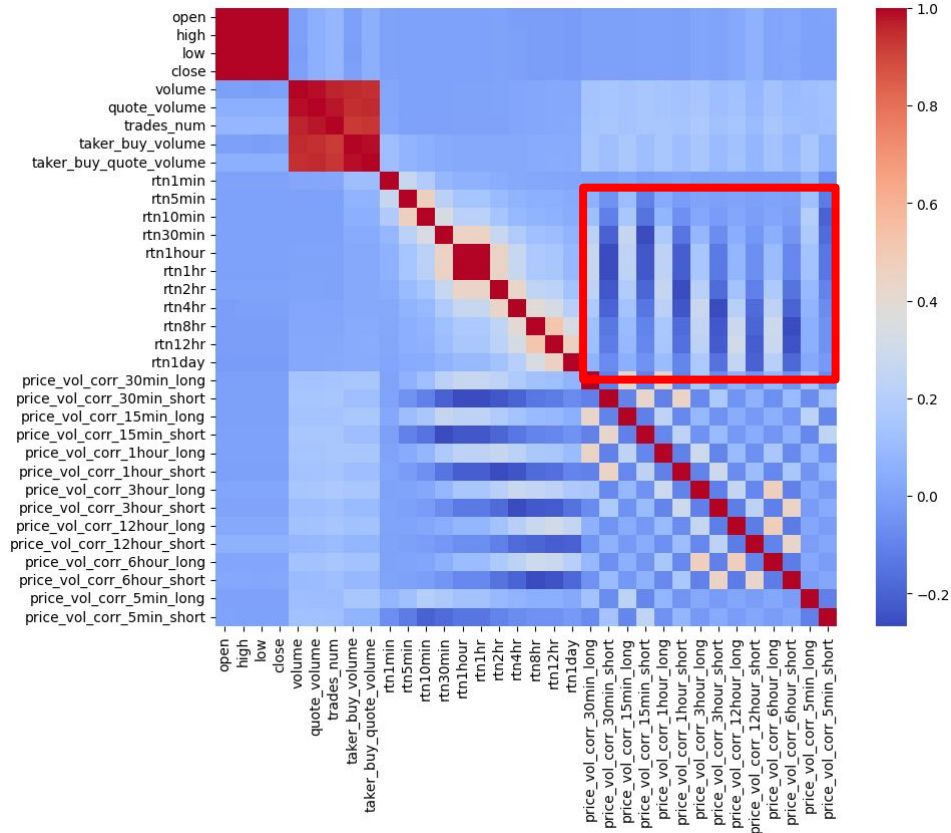


Volatility

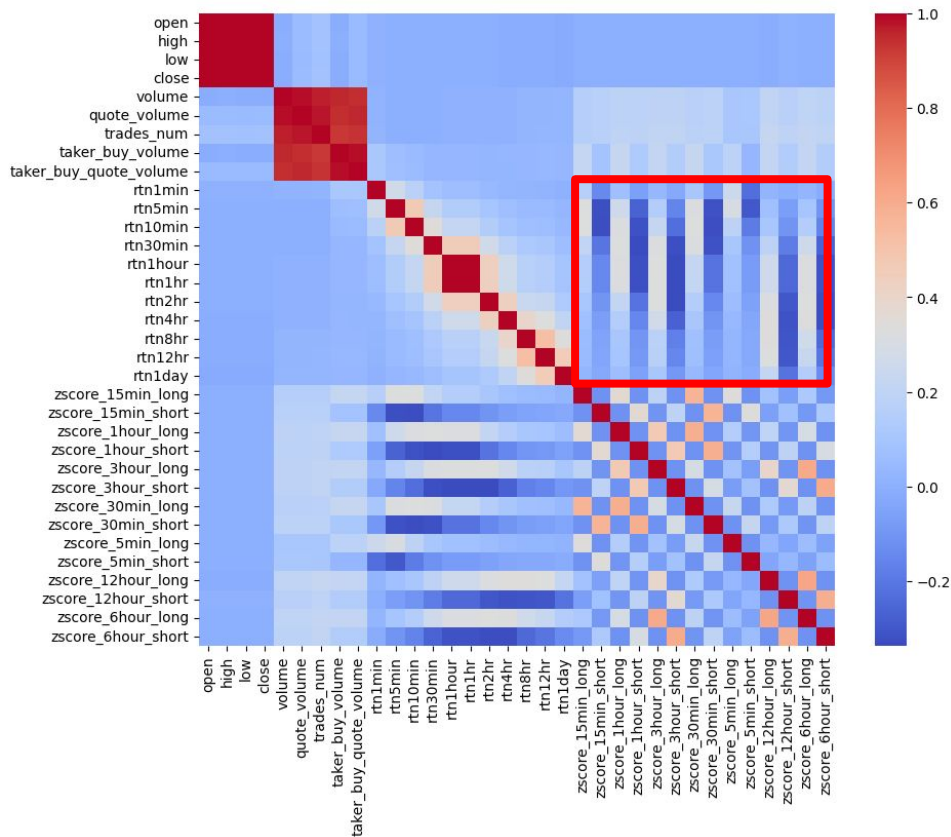


Feature

Price_vol_corr

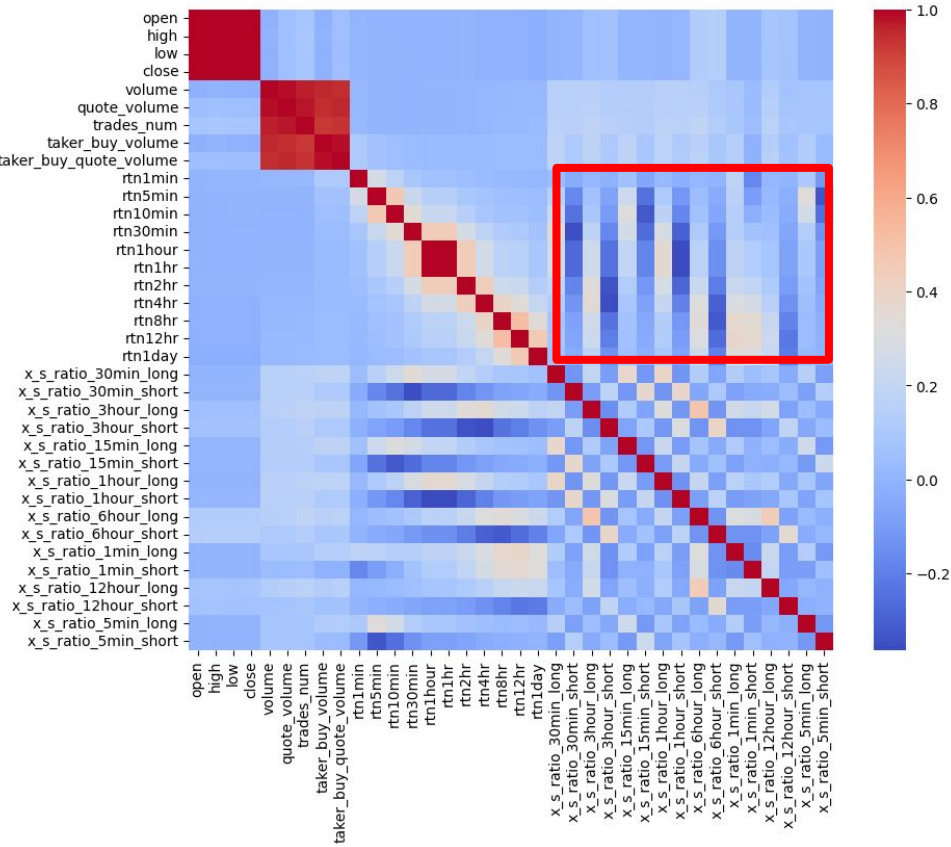


Zscore

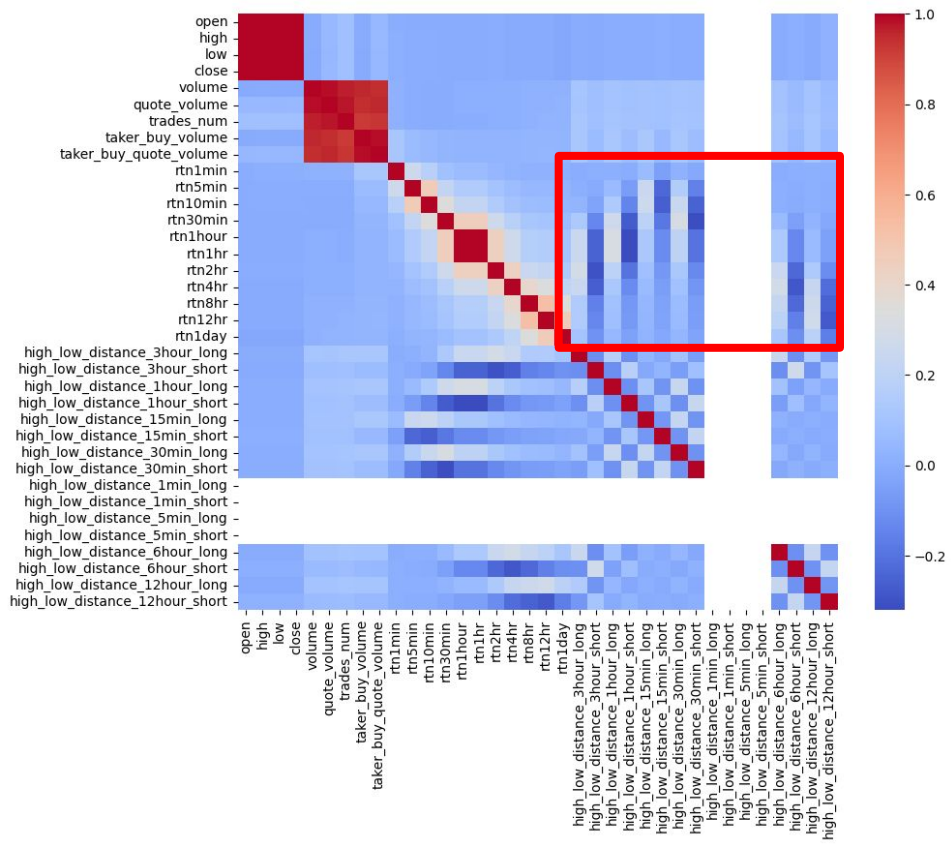


Feature

x_s_ratio



high_low_dis





ML Model

Logistic Regression

Functionality:

Models the probability that an instance belongs to a particular category.

Random Forest

Functionality:

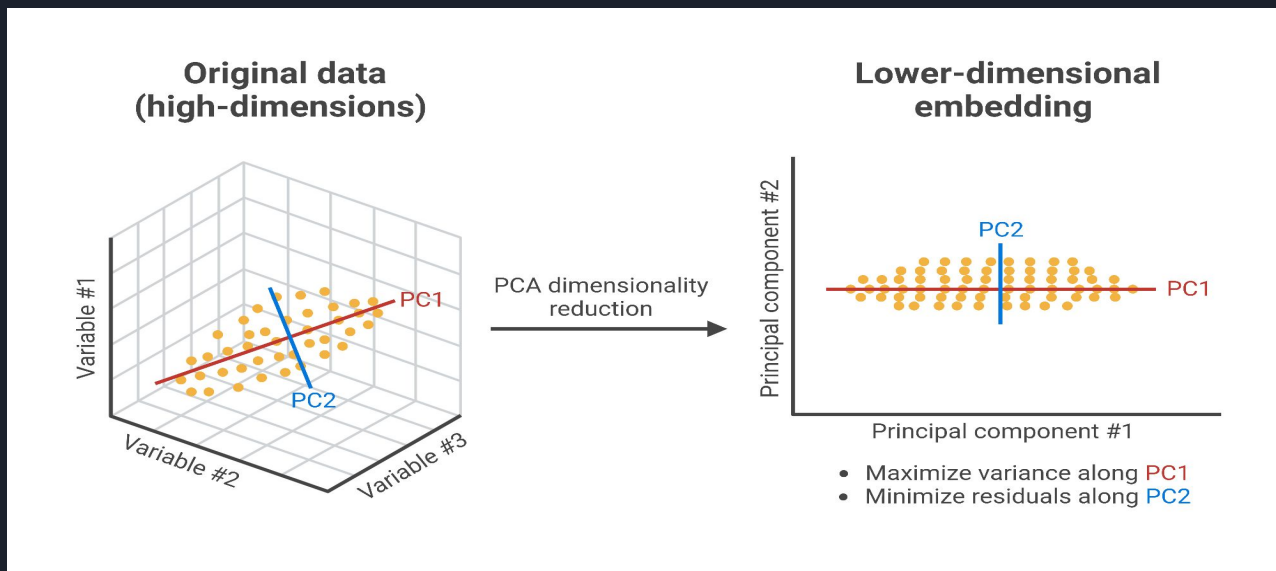
Constructs a multitude of decision trees during training and outputs the mode of the classes or mean prediction of the individual trees.

XGBoost

Functionality:

Builds a series of weak learners (usually decision trees) sequentially, where each new learner corrects the errors made by the previous ones.

Principal Component Analysis



Logistic Regression

Logistic regression

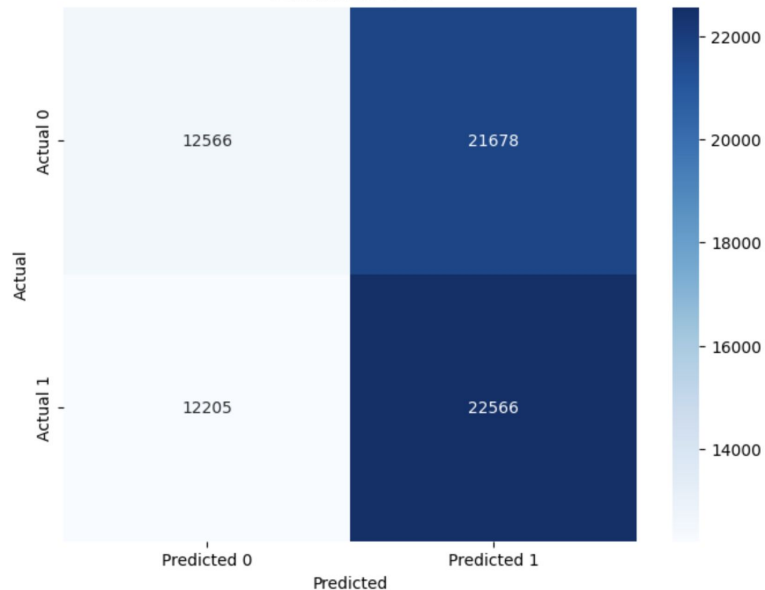
Model Accuracy: 0.5090

Precision: 0.5100

Recall: 0.6490

F1-Score: 0.5712

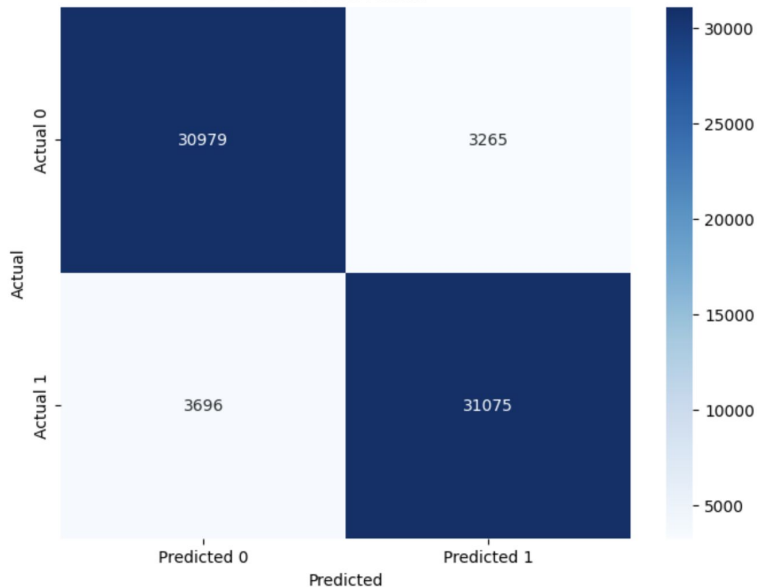
Confusion Matrix



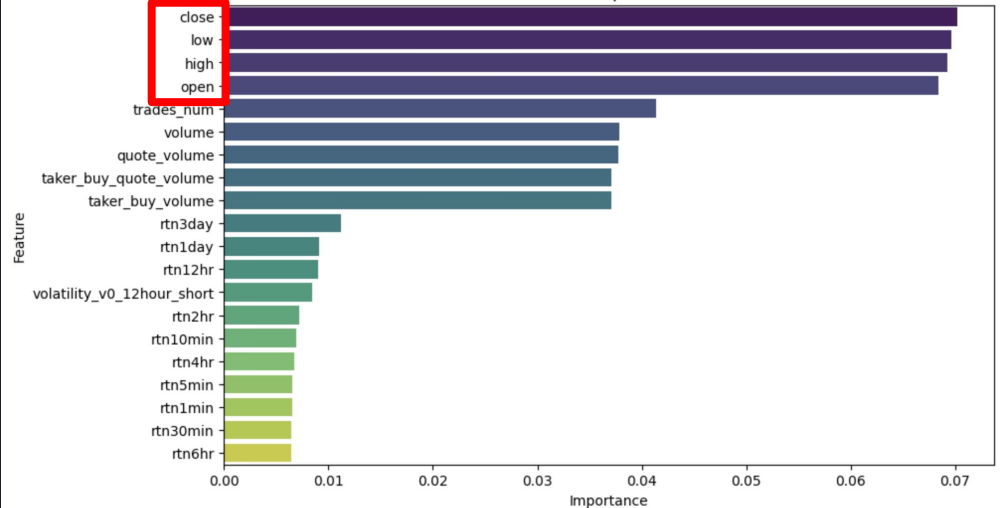
Random Forest

Model Accuracy: 0.8991
Precision: 0.9049
Recall: 0.8937
F1-Score: 0.8993

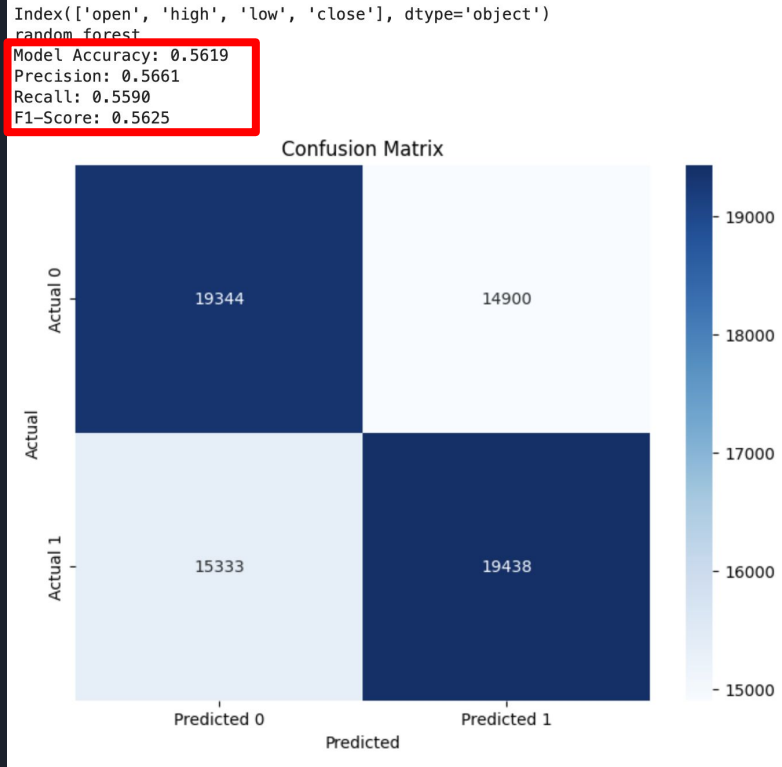
Confusion Matrix



Feature Importances



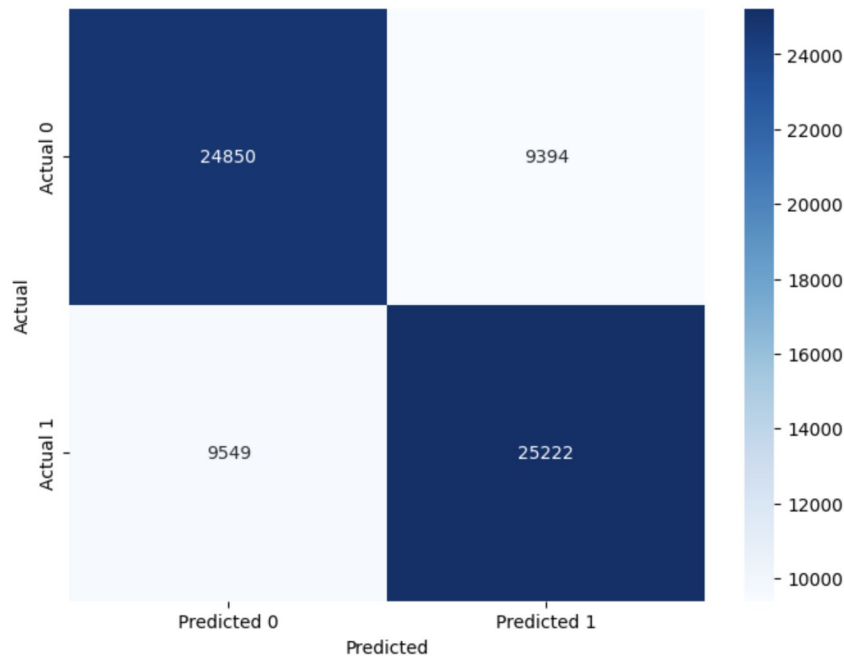
Random Forest (only use OHLC)



XGBoost

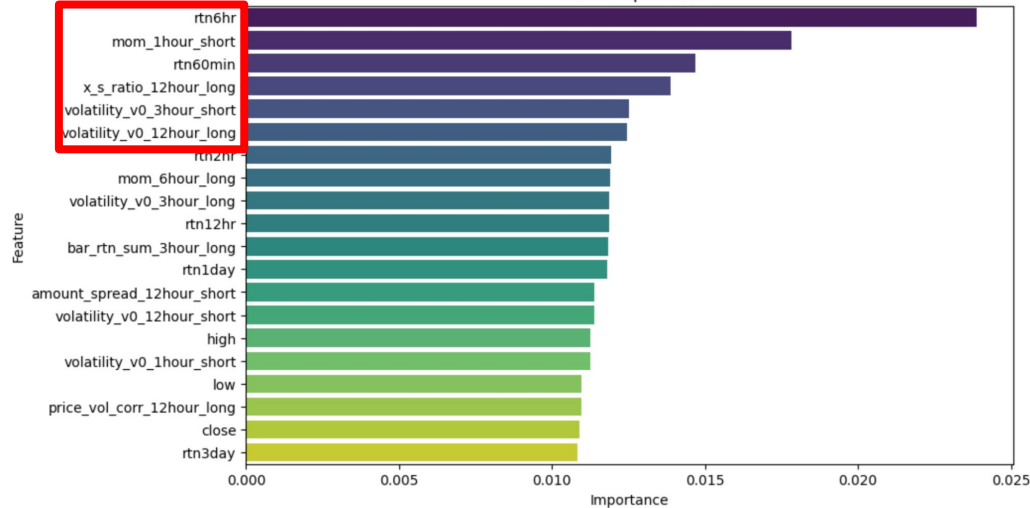
Model Accuracy: 0.7255
Precision: 0.7286
Recall: 0.7254
F1-Score: 0.7270

Confusion Matrix



XGBoost

Feature Importances

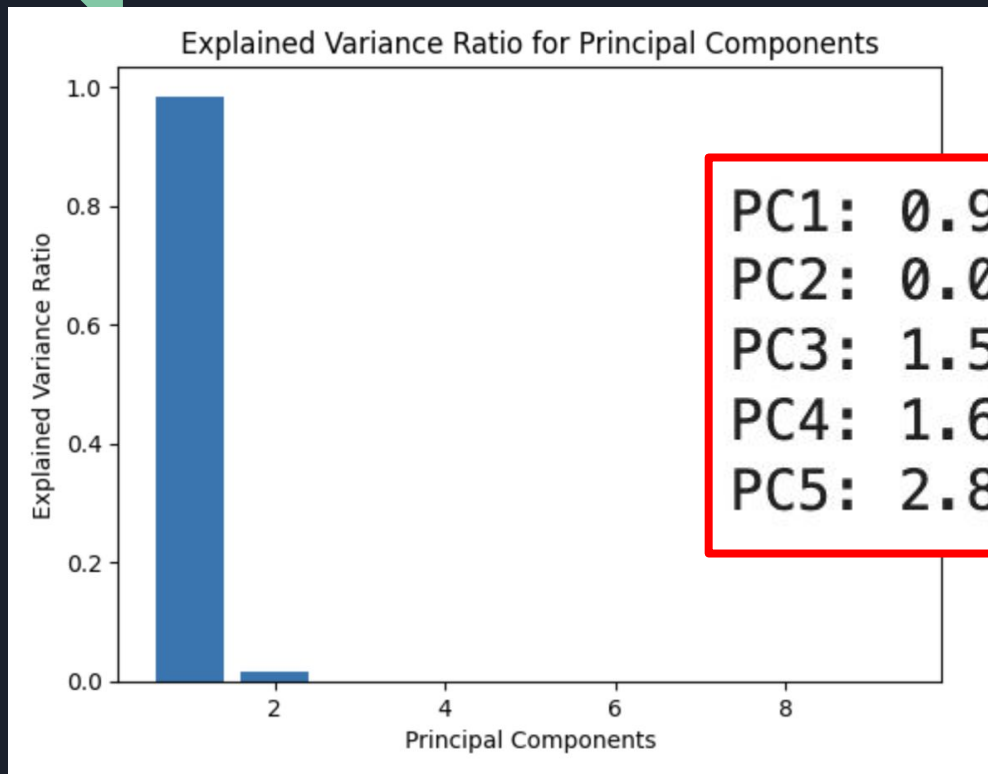




Result (Without PCA)

	Accuracy	Precision	Recall	F1 score
Logistic Regression	0.509	0.5100	0.6490	0.5712
Random Forest	0.8991	0.9049	0.8937	0.8993
XGBoost	0.7255	0.7286	0.7254	0.7270

PCA



PC1: 0.9847721605154578

PC2: 0.015227678278119187

PC3: 1.595094682428055e-07

PC4: 1.6674910168497199e-09

PC5: 2.8728703150190032e-11



Choosing N

	Logistic Regression	Random Forest	XGBoost
10 principal component	0.5067	0.5532	0.5563
50 principal component	0.5067	0.6952	0.6263
100 principal component	0.5067	0.7286	0.6365
158(all) principal component	0.5067	0.7584	0.6608



Result (With PCA)

	Accuracy	Precision	Recall	F1 score
Logistic Regression	0.5067	0.5165	0.3282	0.4014
Random Forest	0.8991	0.9049	0.78937	0.8993
XGBoost	0.6608	0.6643	0.6608	0.6625

Result

Strategy: If predict is 1, long the position and close it 1 hour later.

Win rate is improved from 50% to 69.35%

Profit Facotr is improved from 1.01 to 2.996

Which is the more profitable strategy by Machine Learning.

$$PF = \frac{\text{avg}(\text{profit per trades})}{\text{avg}(\text{loss per trades})} = 2.99615391136621$$

$$\text{Win rate} = \text{win} / \text{total tades} = 0.693475$$