



# OB Oracle

Group 6

# Our Project Goals

- Using machine learning models on a high frequency trading (HFT) crypto order book to predict price direction and price levels.
- Develop a strategy trading bot to place mock trades on predicted outcomes from ML models and compare actual versus realized strategy returns.

# Agenda

1. Order Book (OB) Overview
2. The Selected ML Models
3. Data Prep and Feature Selection
4. Our Approach
5. Results and Conclusions
6. Next Steps

# 1. Understanding Order Book Overview

# What is an Order Book?

- A **Limit order** is an order you place on the order book with a specific limit price
- **Top of Book** represents the highest bid and the lowest ask that time.
- A **bid-ask spread** is the amount by which the ask price exceeds the bid price for an asset in the market.
- **Market orders** let you purchase instantly at best price currently available.
- **Mid-price** is the price between the best price of the sellers offer price and best price of the buyers bid price.
- **Liquidity** refers to how rapidly shares of a stock can be bought or sold without substantially impacting the stock price.

Price	Quantity	Total Quantity
11254.9	0.66676358	18.97374802
11254.4	0.66666365	18.30698444
11253.5	0.25080000	17.64032079
11253.4	1.65441200	17.38952079
11253.2	0.56547391	15.73510879
11253.1	0.40000000	15.16963488
11252.9	0.26741127	14.76963488
11251.9	6.65222361	14.50222361
11251.8	7.85000000	7.85000000
11250.5 USD		
11250.3	0.26741784	0.26741784
11250.2	0.66666906	0.93408690
11249.3	0.88889208	1.82297898
11248.6	1.77400000	3.59697898
11247.5	0.66250000	4.25947898
11247.2	0.88909868	5.14857766
11247.1	0.88909536	6.03767302
11246.7	0.03750000	6.07517302
11246.0	0.10300000	6.17817302



# Order Book Data Source (Kaggle):

## Crypto Assets:

- Bitcoin (BTC)
- Ethereum (ETH)
- Cardano (ADA)

## Content:

- 12 days of limit order book data
- 15 best bid / ask price **levels**
- 5 min data / 1 min data / 1 sec data

<https://www.kaggle.com/datasets/martinsn/high-frequency-crypto-limit-order-book-data>

## High Frequency Crypto Limit Order Book Data

Bitcoin, Ethereum and Cardano



Crypto PROVIDED order book features/columns (\_x represent order book level)

- \* `midpoint` = the midpoint between the best bid **and** the best ask
- \* `spread` = the difference between the best bid **and** the best ask
- \* `bids/asks_distance_x` = the distance of bid/ask level x **from** the midprice **in %**
- \* `bids/asks_limit_x` = volume (= price \* quantity) of orders at bid/ask level\_x
- \* `bids/asks_notional_x` = (`asks_limit_notional_x` - `asks_market_notional_x` - `asks_cancel_notional_x`)

# Order Book Feature Engineering:

	midpoint	bids_distance_0	bids_price_0	bids_limit_quantity_0	bids_limit_cum_quantity_0	bids_limit_notional_0	bids_limit_cum_notional_0	asks_distance_0	asks_price_0	asks_limit_quantity_0	asks_limit_cum_quantity_0	asks_limit_notional_0
system_time												
2021-04-07 11:37:59.055697+00:00	1.17255	-0.000384	1.1721	2099.999925	2099.999925	2461.409912	2461.409912	0.000384	1.1730	8007.340453	8007.340453	9392.610
2021-04-07 11:42:59.055697+00:00	1.18390	-0.000591	1.1832	4633.062748	4633.062748	5481.839844	5481.839844	0.000591	1.1846	20004.229078	20004.229078	23697.009
2021-04-07 11:47:59.055697+00:00	1.17830	-0.000594	1.1776	2798.998045	2798.998045	3296.100098	3296.100098	0.000594	1.1790	500.000000	500.000000	589.500

1.1 New DERIVED order book features/columns developed in this notebook/python script

1.1 \* bids/ask\_price\_x = the price at bid/ask level\_x  
= midpoint \* (1 + distance) (where distance is represented as %)

1.1 \* bids/asks\_limit\_quantity\_x = quantity (= bids/asks\_limit/market/cancel\_notional\_x / bids/ask\_price\_x) of orders at bid/ask level\_x

1.1 \* bids/asks\_cum\_quantity\_x = Cumulative sum of quantities - i.e. bids/asks\_limit/market/cancel\_quantity\_x

Example. cum\_quantity\_0 = quantity\_0

cum\_quantity\_1 = quantity\_0 + quantity\_1

cum\_quantity\_2 = quantity\_0 + quantity\_1 + quantity\_2

...

1.1 \* bids/asks\_cum\_notional\_x = Cumulative sum of notionals - i.e. bids/asks\_limit/market/cancel\_notional\_x

Example. cum\_notional\_0 = notional\_0

cum\_notional\_1 = notional\_0 + notional\_1

cum\_notional\_2 = notional\_0 + notional\_1 + notional\_2

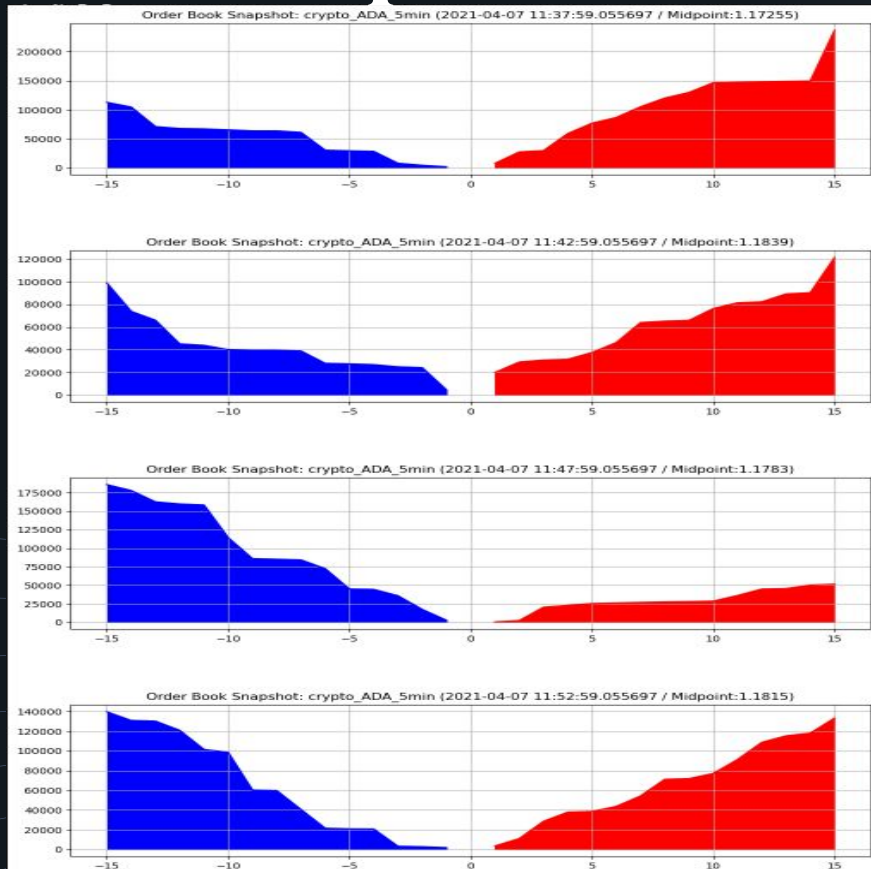
...

\* bid\_ask\_imbalance\_limit\_x = Bid Ask Imbalance = (bid notional / (bid notional + ask notional))

= bids\_notional\_x / bids\_notional\_x + asks\_notional\_x

# Order Book Perspectives

- BID LEVELS  
- ASK LEVELS



5min  
Snapshots

CumQty



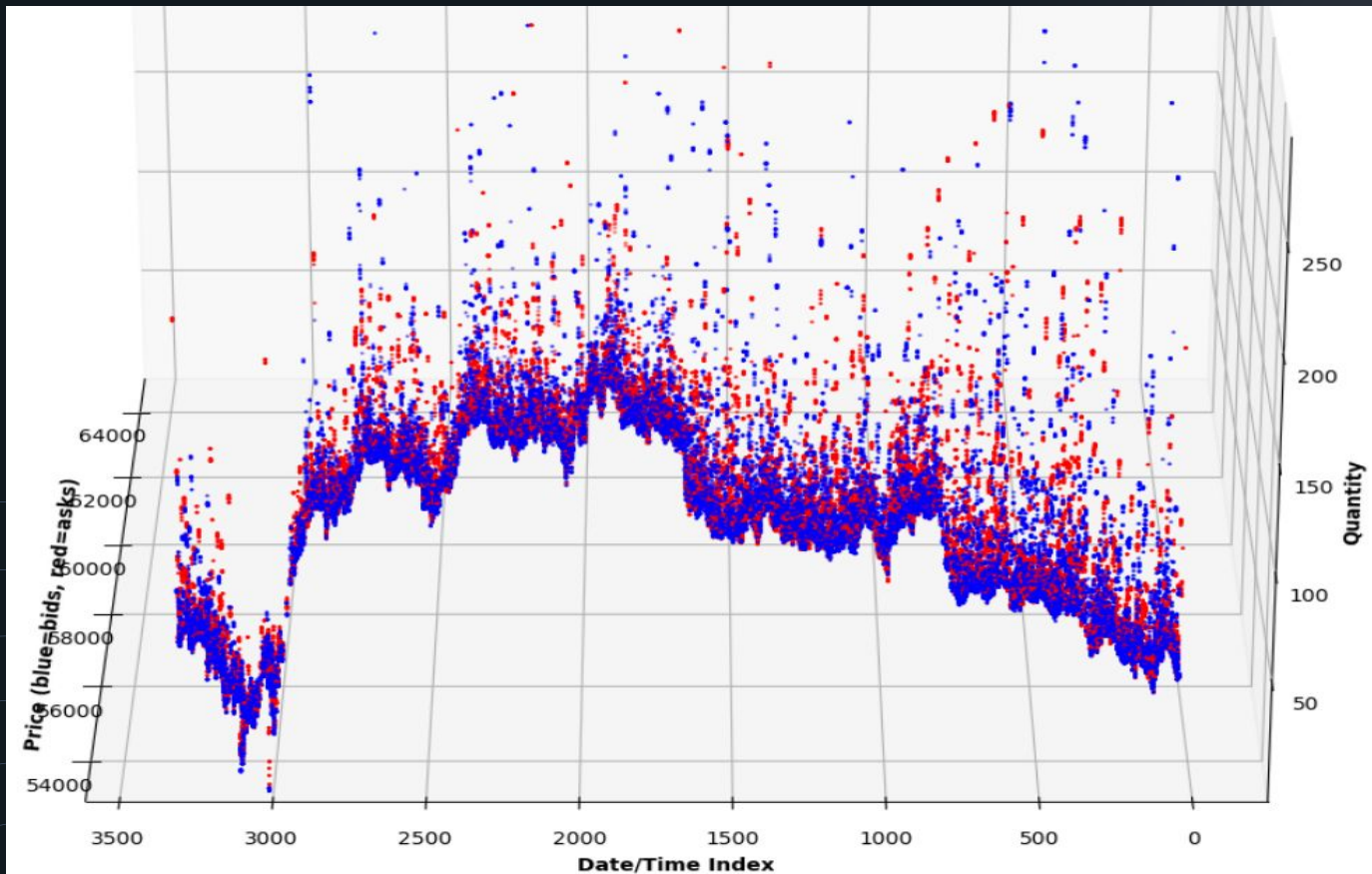
Price





# Order Book Perspectives

- BID LEVELS  
- ASK LEVELS

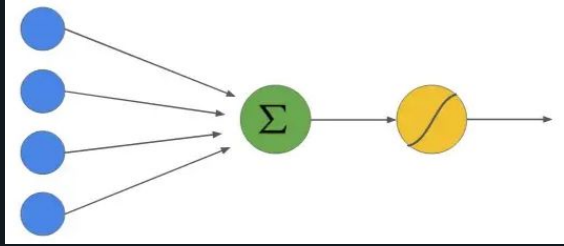


CumQty  
Price  
Time

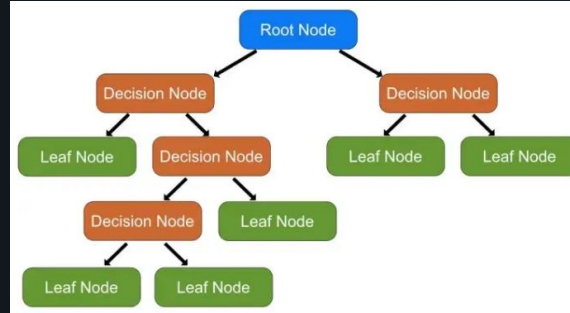
## 2. The Selected ML Models

# ML Supervised Learning: Models

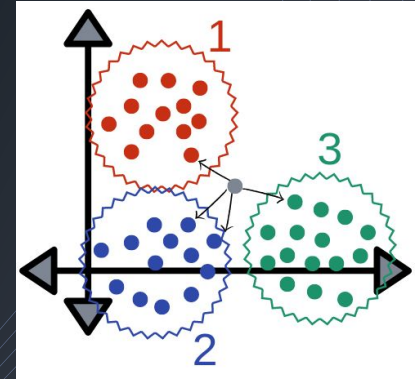
Logistic Regression



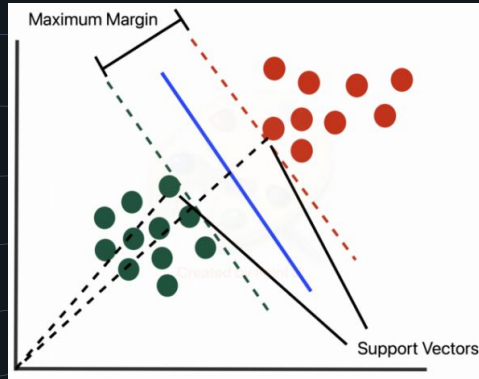
Decision Trees



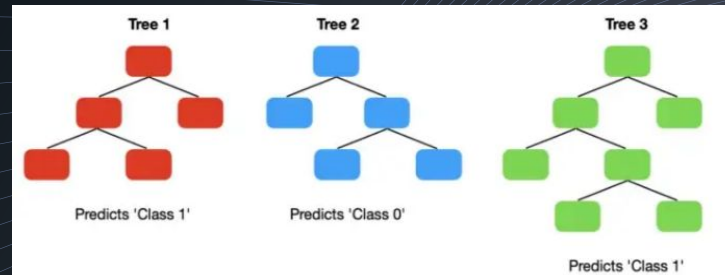
KNN



SVM

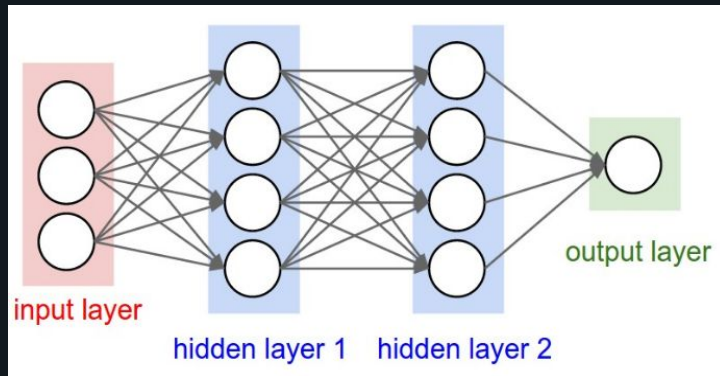


Random Forest

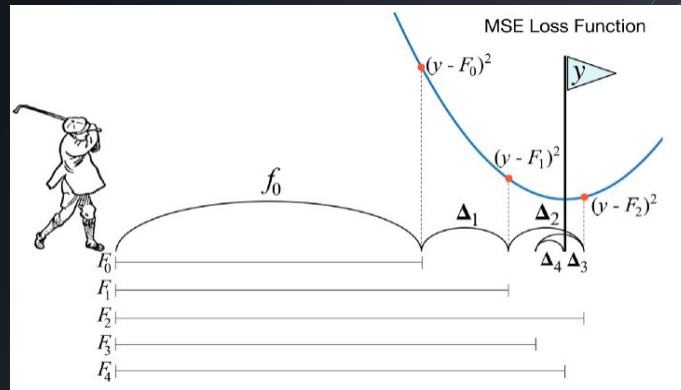


# ML Neural Networks: Models

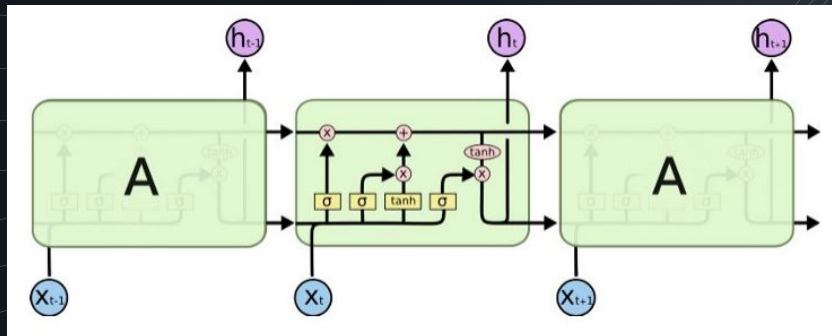
Sequential/Dense



Gradient Boosting\*



LSTM



# ML Supervised Learning: Classification

Predict  
Up Or Down



midpoint  
price

Significant **support** wall at  
top of order book (**bid** side)





# ML Supervised Learning: Regression

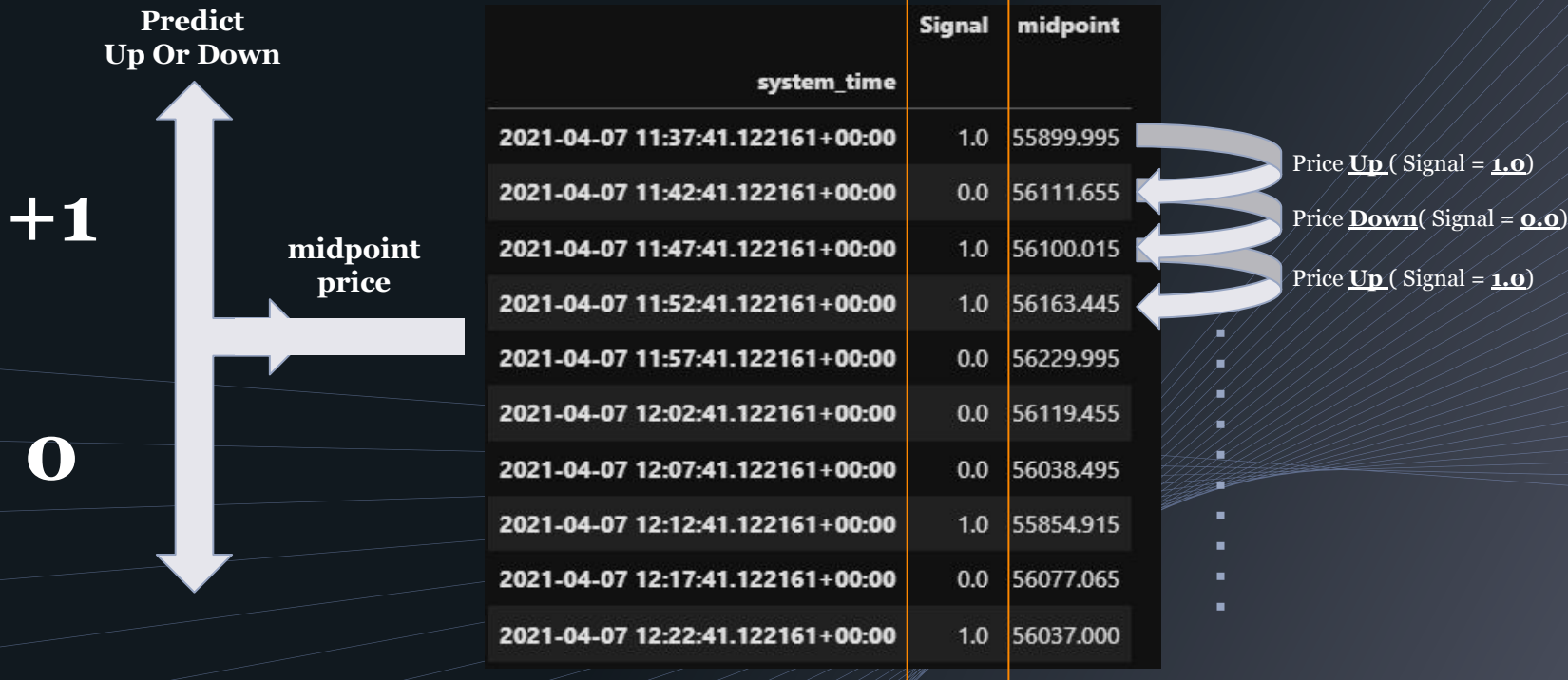
Predict  
Price = X

Predict  
Price = Y



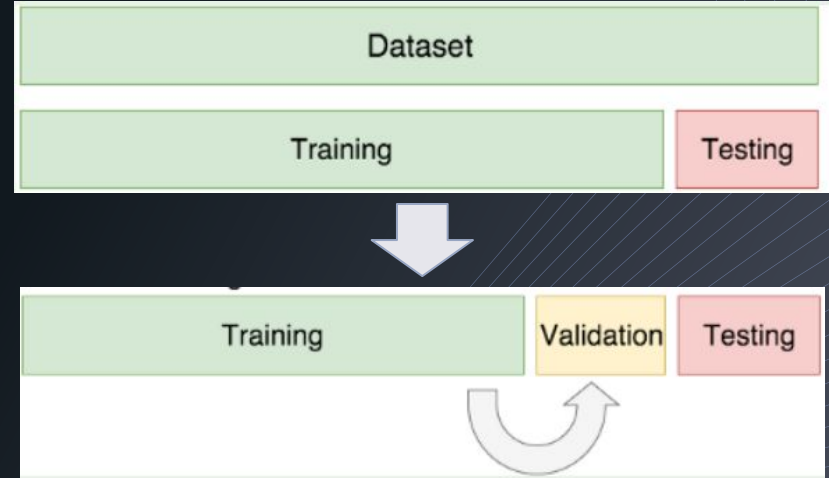
# 3. Data Prep and Feature Selection

# Data Preparation: Generate Target Signal



# Data Preparation: Training/Test Data Set

- **Training Set:** Complete training set that will be applied to the ML models to train and fit the data.
- **Validation Set:** Divide training set into a train set and validation set. Based on validation test results, the model can be trained (e.g. changing parameter classifiers). This will help us get the most optimized model.
- **Testing Set:** Data set used to perform blind predictions and evaluations on the ML model.



- **Data Partition Permutations:**
  - Training Set: 50% → 90%, 10% increments
  - Validation Set: 0% → 30%, 5% increments
  - Testing Set: 10% → 50%, 10% increments

# Data Preparation: Feature Selection / Cleanup

```
# Level selection
# Select number of orderbook levels to include in feature set (max = 15)
MAX_ORDER_BOOK_LEVELS = 15    # Do not modify
NUM_ORDER_BOOK_LEVELS = 15
```

```
# Optional removing of feature/columns from provided crypto order book
# Remove cancel features/columns
REMOVE_CANCEL_FEATURES = True
# Remove market features/columns
REMOVE_MARKET_FEATURES = True
# Remove price distance feature/column
REMOVE_PRICE_DISTANCES = False
# Remove notional feature/column
REMOVE_FINAL_NOTIONAL = True
```

```
# Select crypto assets interested in evaluating
INCLUDE_BTEC = True
INCLUDE_ETH = True
INCLUDE_ADA = True
```

```
# Select the orderbook timeframe's interested in
INCLUDE_5MIN_DATA = True
INCLUDE_1MIN_DATA = True
INCLUDE_1SEC_DATA = False
```

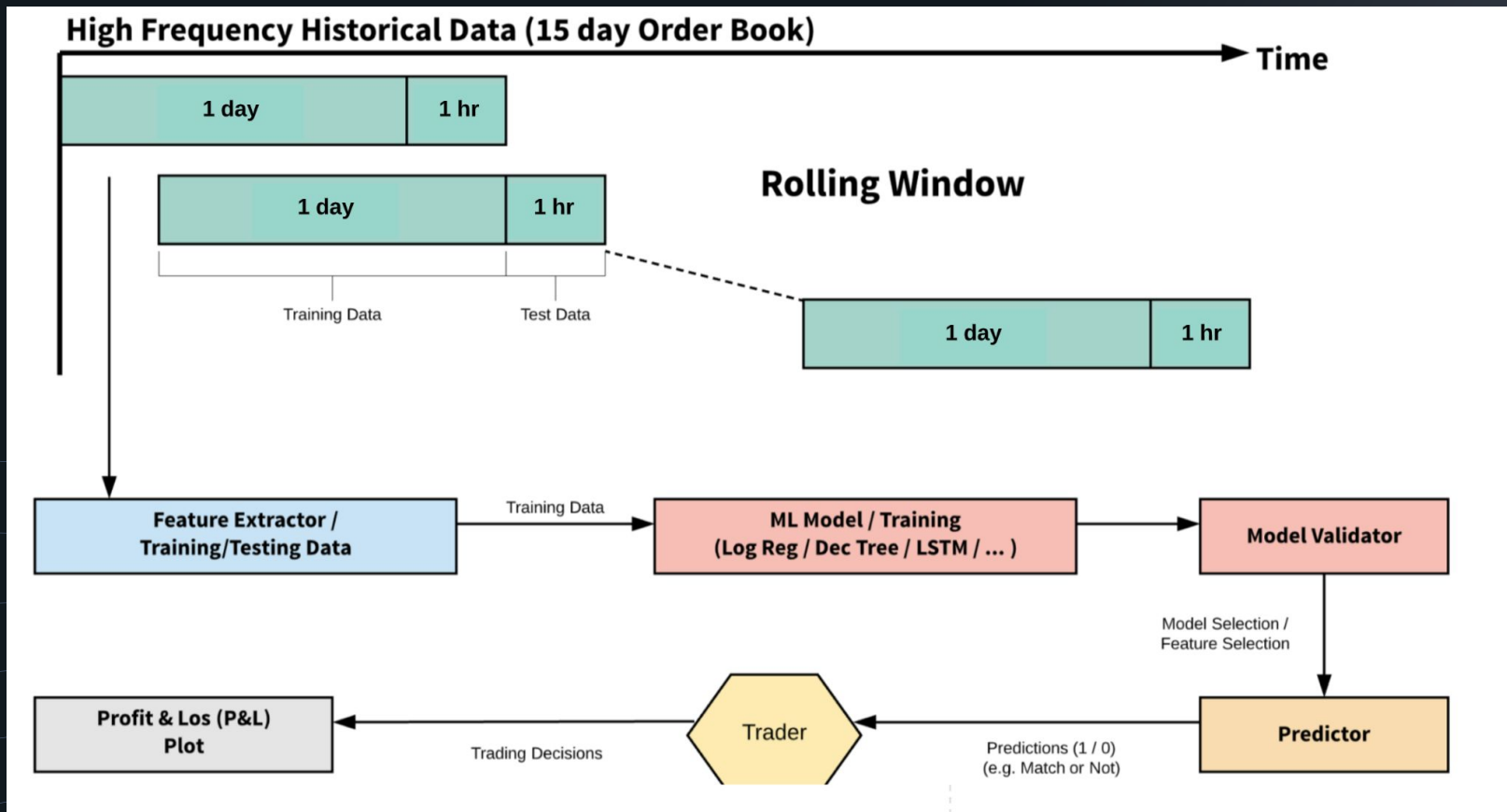
```
# Select which "derived" features to include in ML Model feature set

# Optional NEW derived feature/columns that can be added to dataframe order book
# Option to add absolute price levels to feature set
# (default dataframe shows delta price with mid price = distance)
ADD_ABSOLUTE_PRICE_LEVELS = True
# Option to add quantity at each level
ADD_QUANTITY_FEATURE = True
# Option to add CUMULATIVE quantity "across" levels
ADD_CUM_QUANTITY_FEATURE = True
# Option to add CUMULATIVE quantity "across" levels
ADD_CUM_NOTIONAL_FEATURE = True
# Option to add bid and ask IMBALANCE
ADD_BID_ASK_IMBALANCE = True
# ToDo: ADD MORE FEATURES (e.g. Logarithmic returns)
```

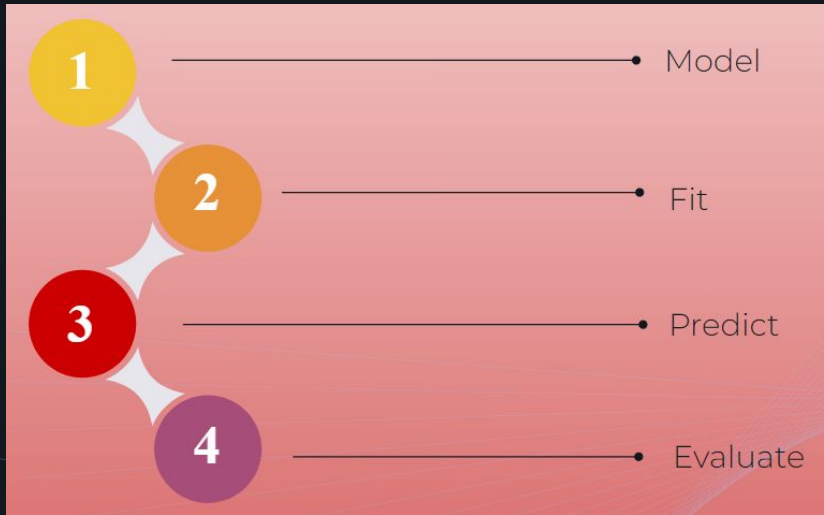


# 4. Our Approach

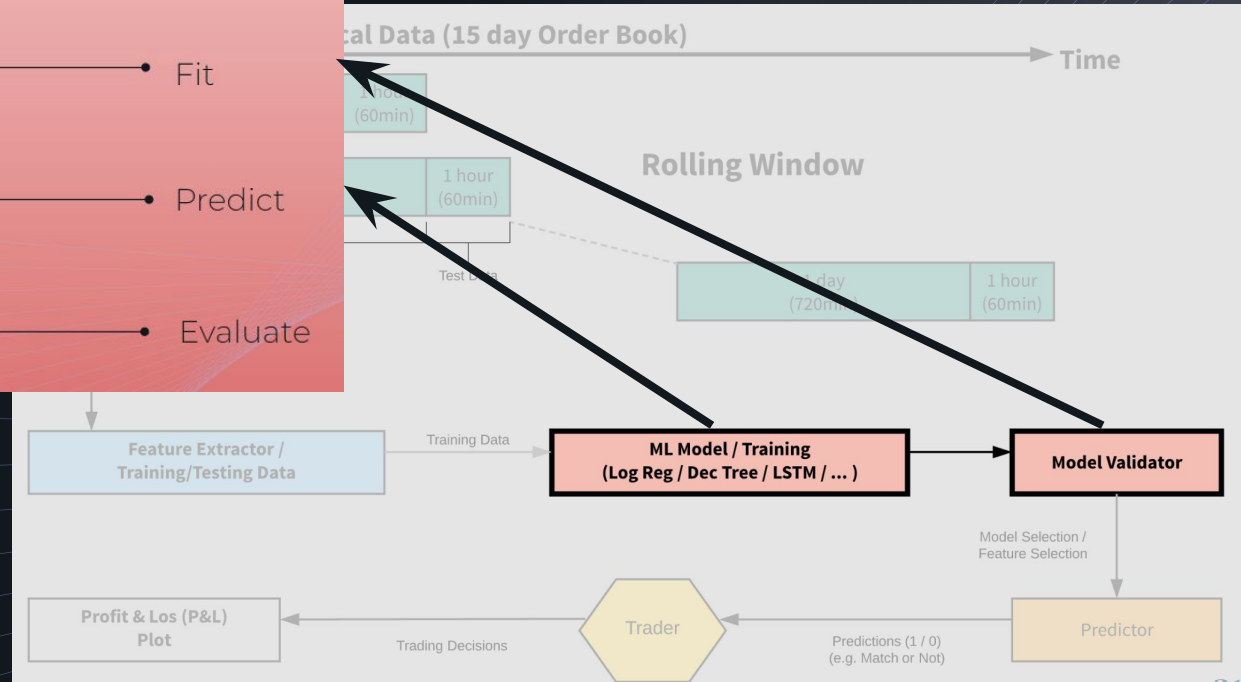
# Framework for ML Limit Order Strategy Submission



# ML Model Training / Validation Steps



- Logistic Regression
- Decision Trees
- KNN
- Random Forest
- SVM
- Sequential/Dense
- LSTM
- Gradient Boosting



# ML Neural Network Model Validation

## ML Parameter Permutations:

→ activation function = {relu, sigmoid, exp, swish}  
→→ # hidden layers = {1, 2, 3}  
→→→ # nodes, hidden layer 1 = {32 : 32 : 128}  
→→→→ # nodes, hidden layer 2 = {16 : 16 : 64}  
→→→→→ # nodes, hidden layer 3 = {4 : 4 : 16}  
→→→→→→→ # of epochs = {50 : 50 : 200}

Total Permutations = 4096 runs

Running Permutation #27

# Input Features = 199  
Use Model Activation Function = relu  
# Hidden Layers = 3  
# Layer 1 Nodes = 96  
# Layer 2 Nodes = 48  
# Layer 3 Nodes = 8  
# of Output Nuerons = 1  
# Epochs = 150

\*\*\* Permutation:27 ==> Accuracy: 0.5231316685676575 \*\*\*  
\*\*\* Best Running Results: Permutation:7 ==> Accuracy: 0.5338078141212463 \*\*\*

Running Permutation #28

# Input Features = 199  
Use Model Activation Function = relu  
# Hidden Layers = 3  
# Layer 1 Nodes = 96  
# Layer 2 Nodes = 48  
# Layer 3 Nodes = 8  
# of Output Nuerons = 1  
# Epochs = 200

\*\*\* Permutation:28 ==> Accuracy: 0.5397390127182007 \*\*\*  
\*\*\* Best Running Results: Permutation:28 ==> Accuracy: 0.5397390127182007 \*\*\*

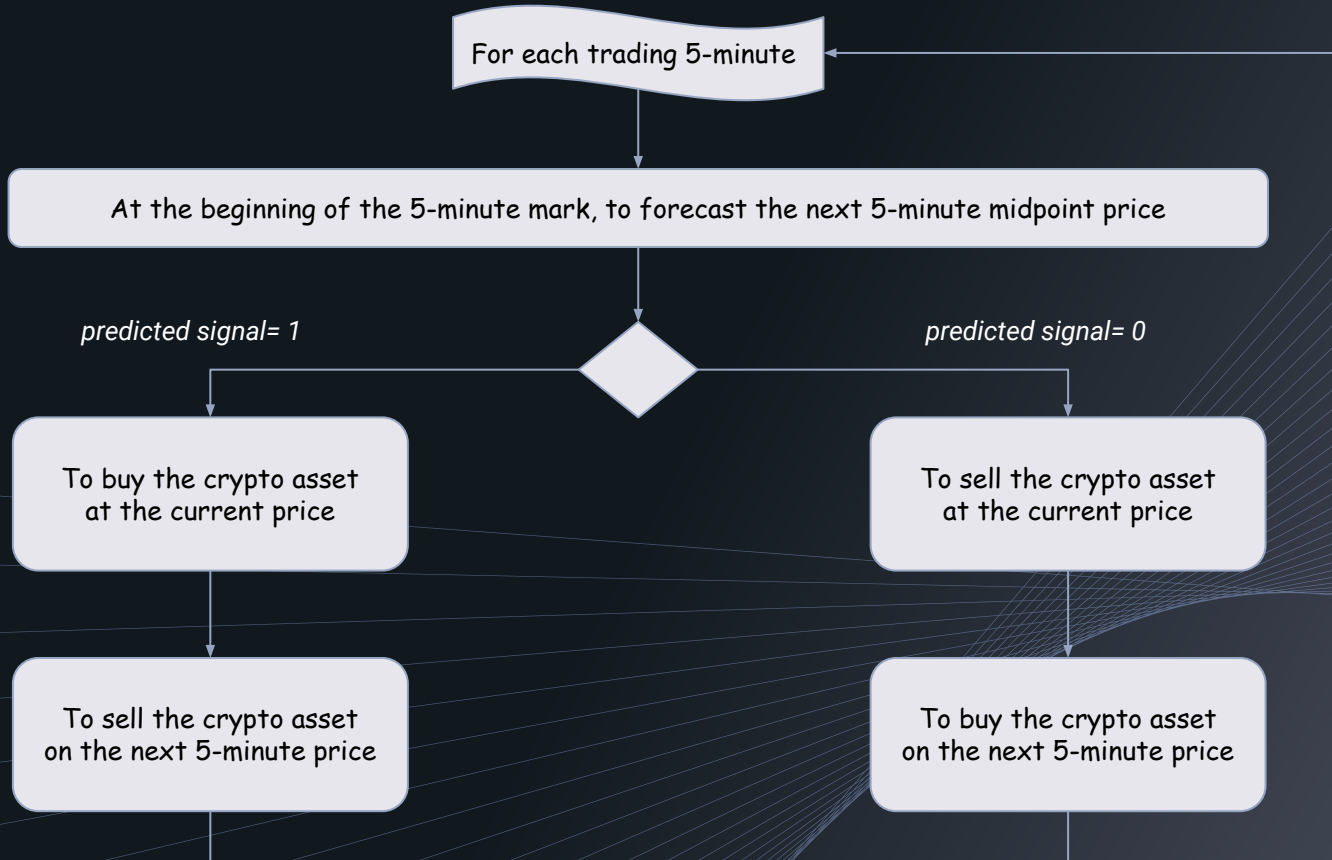
Running Permutation #29

# Input Features = 199  
Use Model Activation Function = relu  
# Hidden Layers = 3  
# Layer 1 Nodes = 96  
# Layer 2 Nodes = 48  
# Layer 3 Nodes = 13  
# of Output Nuerons = 1  
# Epochs = 50

\*\*\* Permutation:29 ==> Accuracy: 0.4922894537448883 \*\*\*  
\*\*\* Best Running Results: Permutation:28 ==> Accuracy: 0.5397390127182007 \*\*\*

*Keeping track of model parameters that provides the best accuracy performance*

# Strategy Trading Bot Flowchart (Classification Model)



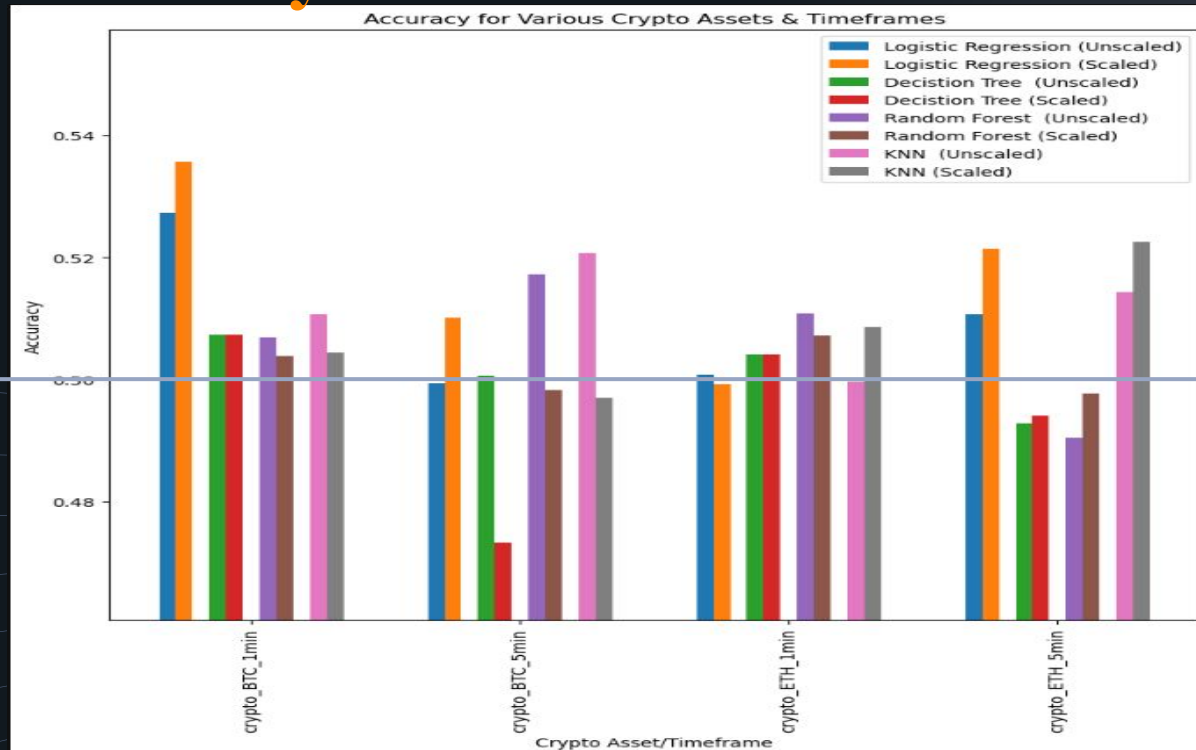


# 5. Results and Conclusions

# ML Supervised Learning Model Results

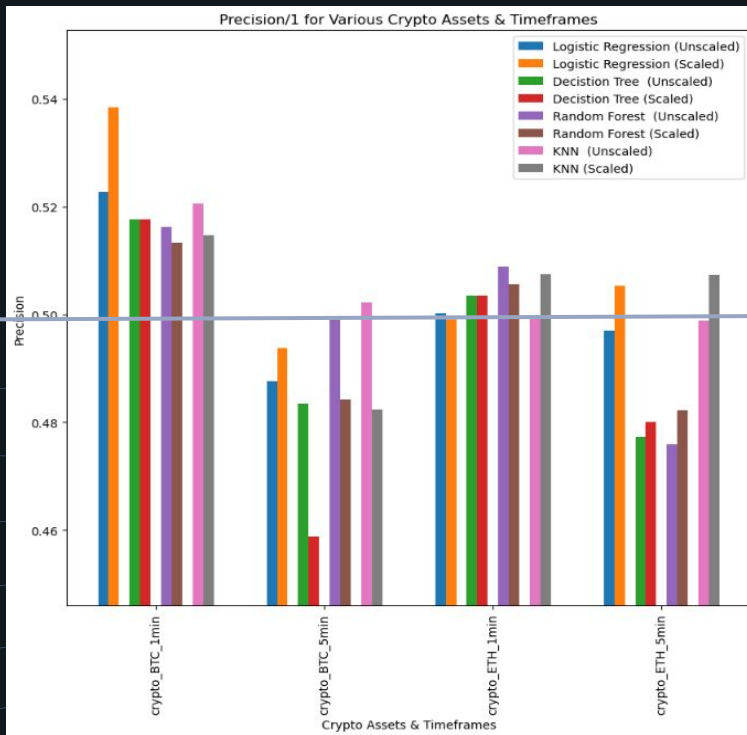
## Accuracy

50%



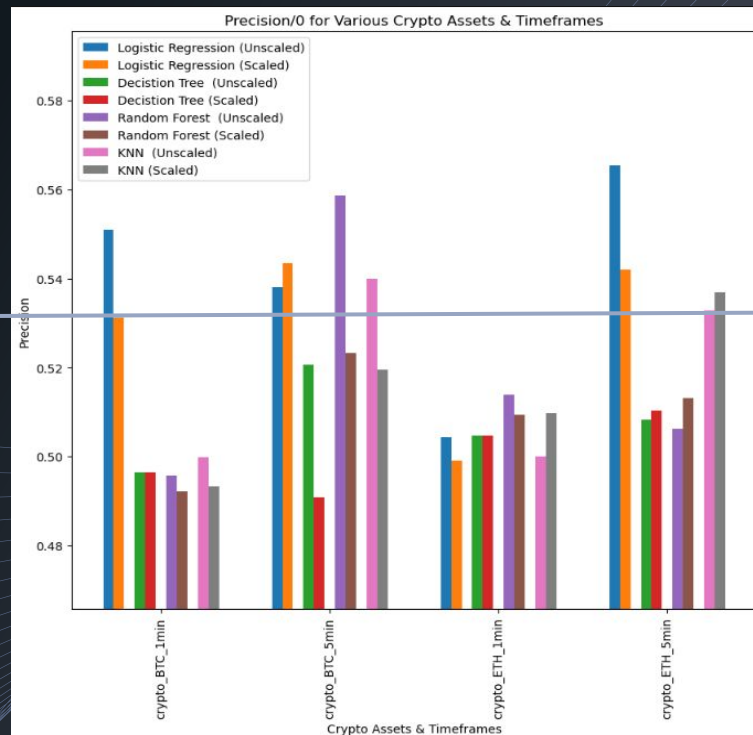
# ML Supervised Learning Model Results

## Precision-1



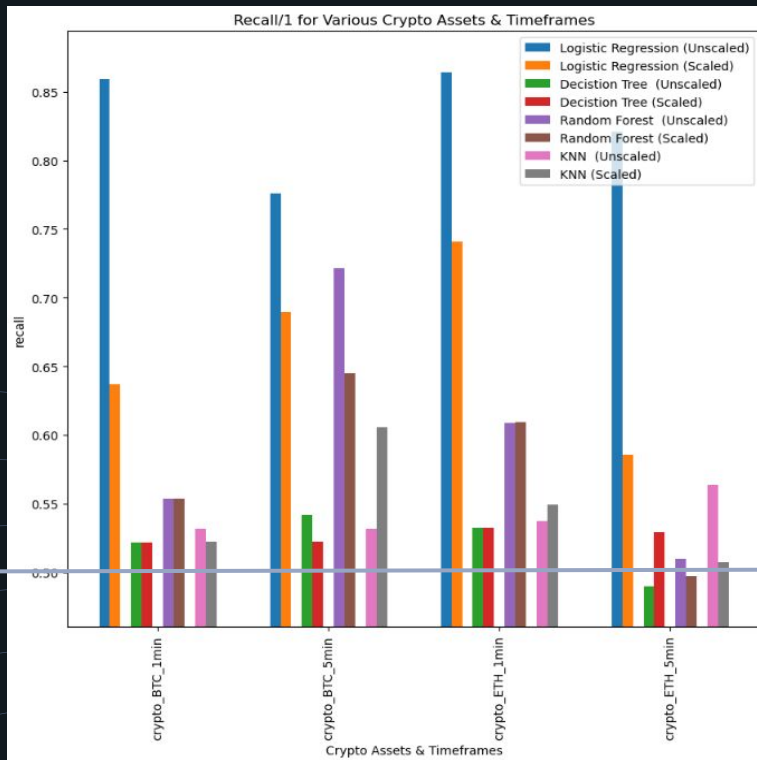
50%

## Precision-0



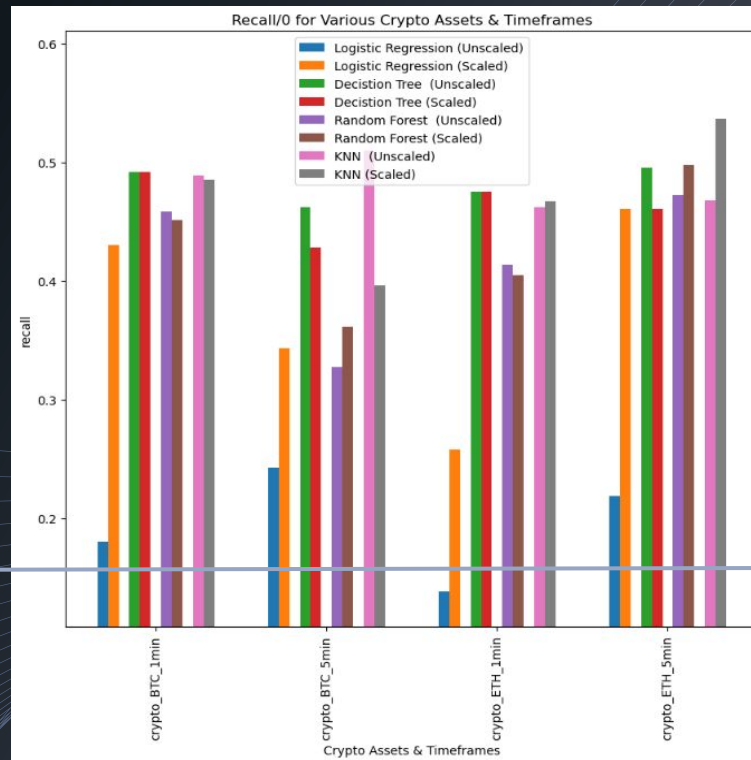
# ML Supervised Learning Model Results

## Recall-1

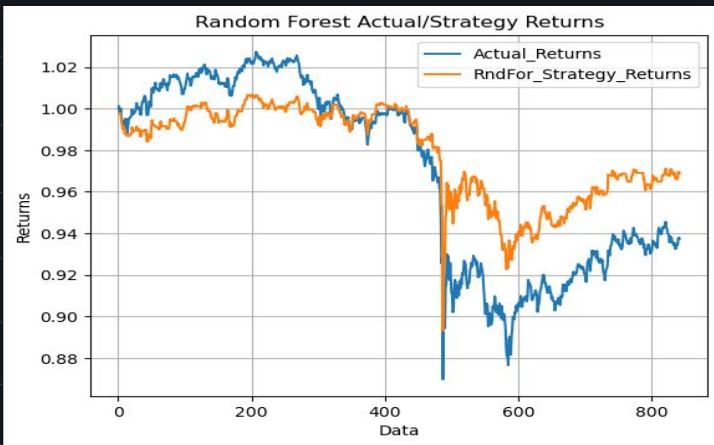
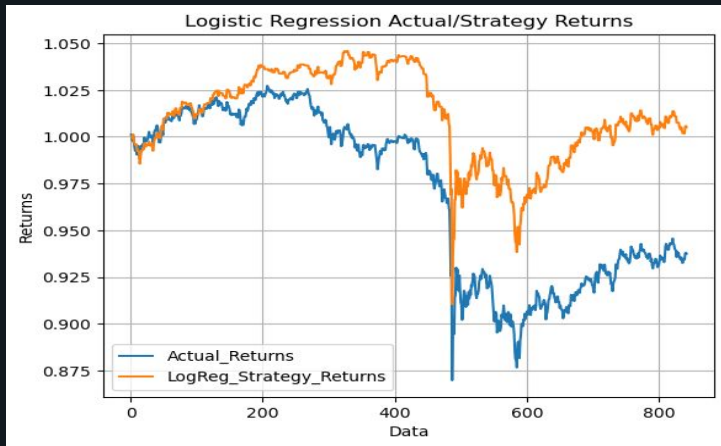


50%

## Recall-0

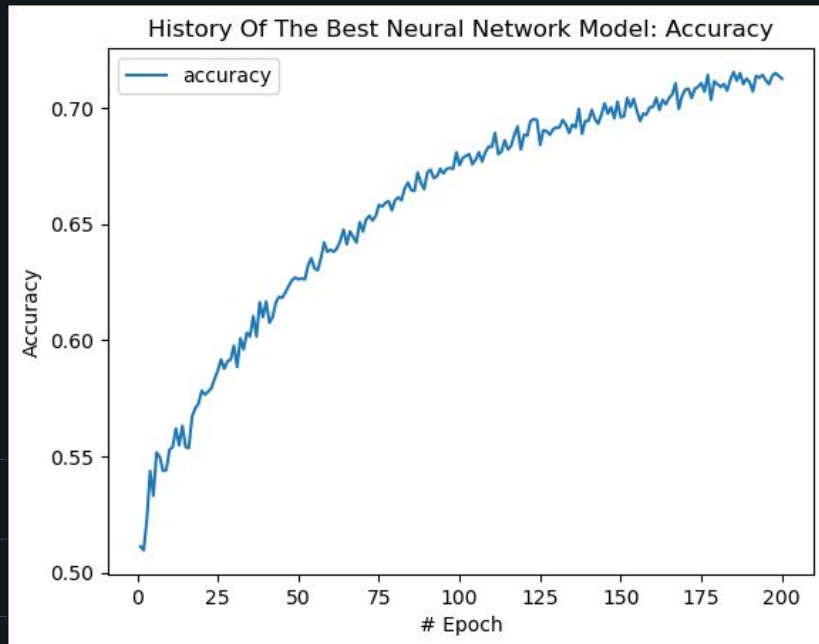


# ML Supervised Learning Strategy Bot Results

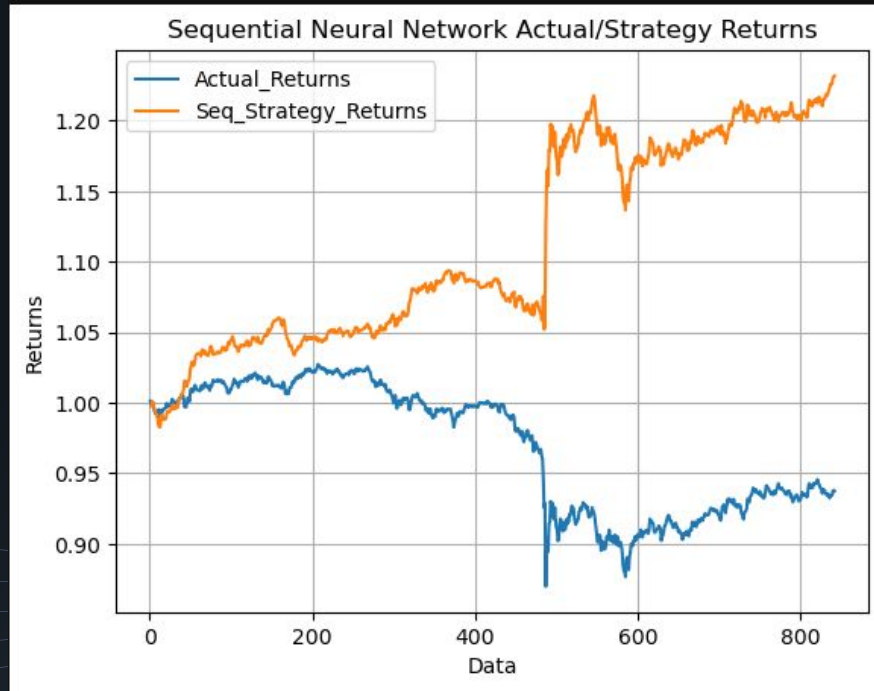




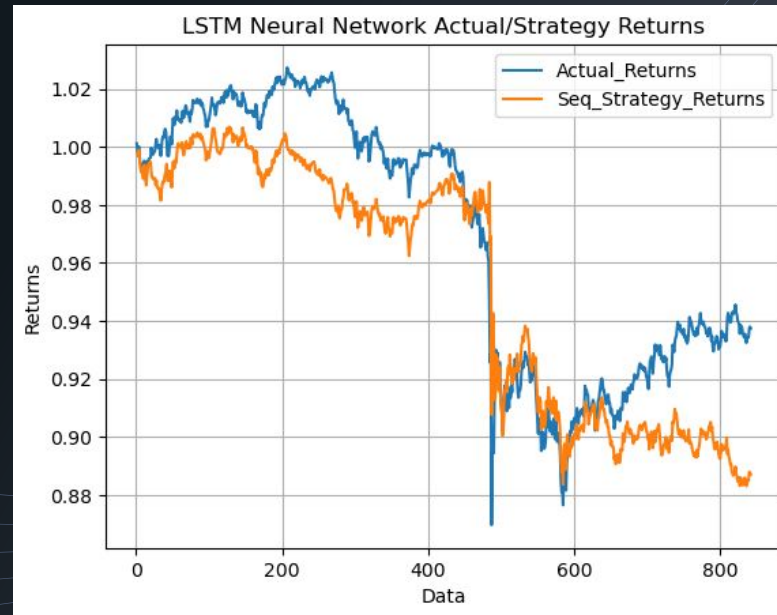
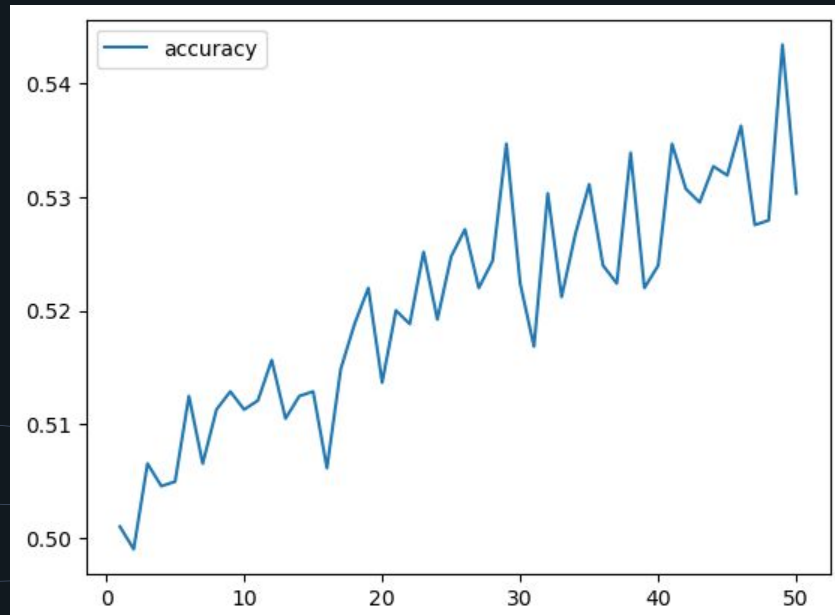
# ML Neural Network Dense Strategy Bot Results



\*\*\* Final Best Results: Permutation:20 ==> Accuracy: 0.5397390127182007 \*\*\*



# ML Neural Network LSTM Strategy Bot Results



Epoch 50/50

79/79 [=====] - 1s 12ms/step - loss: 0.6890 - accuracy: 0.5303

# 6. Next Steps

# Hurdles

- Data acquisition
- Model accuracy
- Binary  $\pm 1$  target signals  $\longrightarrow$  binary 1/0 target signals.

# Further Analysis

- Complete neural networks “price prediction” for project submission.
- Additional ML models
- Additional feature engineering.



# Thank you

<https://github.com/lisetlopez/project2>