



**University of Niagara Falls Canada**

**Master of Data Analytics**

---

**Case Study Part – 2**

---

**Course:** Spring 2025 Data Analytics Case Study- 3 (DAMO-611-3)

**Professor:** Cosimo Girolamo

**Submitted by:**

Rafael Ferreira Martins

Mansehaj singh Gill

Shilpa Shilpa

Anmol Jakhu

**Completion Date:-** June 22,2025

**Niagara Falls, ON, Canada**

## Table of Contents

<b>Statement of Purpose .....</b>	<b>4</b>
<b>Scope of the project.....</b>	<b>4</b>
<b>Background Research.....</b>	<b>5</b>
<b>Phase 1: Problem Definition, Research Questions, Hypotheses and Data Collection .....</b>	<b>6</b>
<b>1. Problem Definition .....</b>	<b>6</b>
<b>2. Research Questions .....</b>	<b>6</b>
<b>3. Hypotheses .....</b>	<b>7</b>
<b>Data Collection. ....</b>	<b>8</b>
<b>1. Data Sources .....</b>	<b>8</b>
<b>2. Methods and Preprocessing .....</b>	<b>8</b>
<b>3. Feature Summary .....</b>	<b>8</b>
<b>Phase 2: Data Understanding .....</b>	<b>9</b>
<b>Summary Statistics .....</b>	<b>9</b>
<b>Key Observations .....</b>	<b>10</b>
<b>Phase 3: Data Visualization.....</b>	<b>13</b>
<b>Interactive Dashboard (Tableau).....</b>	<b>16</b>
<b>Phase 4: Model Building.....</b>	<b>18</b>
<b>Phase 5: Model Evaluation.....</b>	<b>19</b>
<b>1. Classification Performance .....</b>	<b>19</b>
<b>Hypothesis Validation .....</b>	<b>26</b>
<b>H1: The trained ML model yields significantly better trading outcomes than random     signal generation. ....</b>	<b>26</b>
<b>H2: Future returns differ significantly across signal types (Buy, Sell, Hold).....</b>	<b>30</b>
<b>Phase 6: Business Impact .....</b>	<b>31</b>

<b>Strengths and Areas of Improvement .....</b>	<b>33</b>
<b>Conclusion .....</b>	<b>35</b>
<b>References .....</b>	<b>36</b>
<b>Appendix A .....</b>	<b>37</b>

## Statement of Purpose

In the highly volatile and sentiment-driven cryptocurrency market, retail and institutional investors alike face significant challenges in making short-term trading decisions. Traditional predictive models may indicate likely price movements, but they often fall short in offering prescriptive actions — that is, specific, confidence-based Buy, Sell, or Hold recommendations that align with future profitability.

Moreover, most trading signals fail to incorporate market psychology metrics (such as the Fear & Greed Index), or combine them effectively with technical indicators like RSI, MACD, and Bollinger Bands. This limits their real-world usability in portfolio decision-making, especially in fast-paced, hourly trading contexts.

This project addresses this gap by designing a prescriptive analytics model that not only forecasts short-term market movements but also generates actionable signals with quantifiable confidence. The model is further validated against random strategies and statistically tested to ensure the economic significance and statistical reliability of its decisions.

## Scope of the project

This project aims to develop a prescriptive analytics framework for short-term cryptocurrency trading decisions using hourly data and sentiment indicators. The core focus is on generating actionable Buy, Sell, or Hold signals that are not only accurate but also statistically and financially validated.

To achieve this, the following tasks were undertaken:

- **Data Collection:**

Hourly OHLCV data for top cryptocurrencies was sourced and merged with the daily **Fear & Greed Index** to incorporate market sentiment.

- **Feature Engineering:**

Lagged percentage returns, volatility metrics, and technical indicators (RSI, MACD, Bollinger Bands) were computed for each asset.

- **Target Variable Construction:**

The model was trained to classify future 3-hour price movements into three categories — Up, Down, or Stable — based on thresholded future returns.

- **Model Training:**

A multi-class LightGBM model was trained with SMOTETomek for class balance and calibrated using isotonic regression to improve probability estimates.

- **Prescriptive Signal Generation:**

Using model confidence scores, threshold-based signals were generated to prescribe actions with clear rationale.

- **Backtesting:**

A rule-based simulation engine tested the financial impact of model-driven signals, considering stop-loss, take-profit, and timeout conditions.

- **Validation and Hypothesis Testing:**

Two core hypotheses were validated:

1. That model-guided signals produce significantly higher portfolio returns than random signals.
2. That different signal types (Buy, Sell, Hold) lead to significantly different return distributions.

The project stays within the scope of short-term trading decisions, emphasizing statistical validation, financial performance, and real-world applicability of prescriptive modeling.

## **Background Research**

This project draws on both academic and industry literature related to cryptocurrency volatility and high-frequency trading:

1. **Forecasting Cryptocurrency Volatility Using Machine Learning**

*IEEE Transactions on Computational Finance, 2021*

→ This study influenced our feature engineering process by emphasizing the predictive

value of volatility metrics (e.g., rolling standard deviation) and momentum indicators such as RSI and MACD.

## 2. High-Frequency Trading in Crypto Markets

*Journal of Financial Data Science, 2022*

→ Inspired the design of our prescriptive strategy simulation, including the application of auto-stop loss and take-profit rules, as well as benchmarking performance against randomized trading strategies.

# Phase 1: Problem Definition, Research Questions, Hypotheses and Data Collection

## 1. Problem Definition

Cryptocurrency markets are known for their extreme volatility and lack of traditional valuation anchors, making short-term trading particularly risky and uncertain. Investors and traders are often influenced by market sentiment and price momentum rather than fundamentals. This project seeks to build a data-driven, interpretable decision-support system that can analyze high-frequency crypto data and generate prescriptive trading signals to enhance portfolio performance.

The core problem we address is:

*"Can a machine learning model train on technical indicators and sentiment data generate profitable short-term trading signals compared to random strategies?"*

## 2. Research Questions

To guide this investigation, the following research questions were defined:

- **RQ1:** *Does a machine learning-based signal generation system outperform a random trading strategy in terms of portfolio value over a 1-year period of high-frequency cryptocurrency data?*
- **RQ2:** *Do the types of generated signals (Buy, Hold, Sell) correspond to significantly different distributions of future returns?*

These questions aim to validate both the strategic impact of the model and the informational content of the signals it produces.

### 3. Hypotheses

Based on the research questions, two testable hypotheses were developed:

- **H1:** Portfolio Performance Hypothesis

*The model-guided trading strategy yields a statistically higher final portfolio value than random trading strategies.*

**Null hypothesis ( $H_0$ ):**  $\mu_{\text{model}} \leq \mu_{\text{random}}$

**Alternative hypothesis ( $H_1$ ):**  $\mu_{\text{model}} > \mu_{\text{random}}$

(Where  $\mu_{\text{model}}$  = mean portfolio value using the model,  $\mu_{\text{random}}$  = mean portfolio value using random signals)

→ Tested using a one-sample t-test comparing the model's final portfolio value to that of multiple random-signal simulations.

- **H2:** Signal Return Hypothesis

*The distributions of future returns differ significantly across the Buy, Hold, and Sell signals generated by the model.*

Let:

- $R_{\text{Buy}}, R_{\text{Sell}}, R_{\text{Hold}}$  be the distributions of future returns for instances where the model generated Buy, Sell, and Hold signals, respectively.
- $F_X$  denote the cumulative distribution function (CDF) of return distribution X.
- We test whether these three distributions are statistically the same using the **Kruskal–Wallis H-test**, a non-parametric method for comparing more than two independent distributions.
- **Null Hypothesis ( $H_0$ ):**
  - $F_{R_{\text{Buy}}} = F_{R_{\text{Sell}}} = F_{R_{\text{Hold}}}$  (All signal types come from the same distribution of future returns.)
- **Alternative Hypothesis ( $H_1$ ):**
  - $\exists (i,j) : F_{R_i} \neq F_{R_j}, \text{ for } i \neq j, i, j \in \{\text{Buy, Sell, Hold}\}$  (At least one pair of signal types has a statistically different return distribution.)

→Tested using a Kruskal-Wallis H-test (non-parametric ANOVA) on future return distributions by signal type.

## Data Collection.

### 1. Data Sources

- **Binance API:** Used to collect hourly OHLCV (Open, High, Low, Close, Volume) trading data for BTCUSDT and ETHUSDT over the past 1 year. Binance is a leading crypto exchange known for its real-time liquidity and trading accuracy.
- **CoinGecko API:** Provided daily market capitalization data for each coin, enabling us to calculate their market dominance — a proxy for investor sentiment and relative strength.

### 2. Methods and Preprocessing

- **Resampling:** Raw hourly OHLCV data was resampled into 3-hour intervals to reduce noise while maintaining responsiveness for short-term strategies.
- **Market Cap Alignment:** Daily market cap data was resampled into 3-hour intervals using forward-filling to match the Binance timestamps.
- **Merging:** Data from both APIs was merged on timestamp, resulting in a clean, time-aligned dataset.

### 3. Feature Summary

The final dataset includes the following engineered features used for modeling and hypothesis testing:

Feature Name	Description
<b>open, high, low, close</b>	Standard OHLC prices over 3-hour intervals
<b>volume</b>	Total trading volume over each 3-hour period
<b>price_change</b>	Difference between close and open prices for each interval
<b>price_pct_change</b>	Percentage change in price within each 3-hour period



Feature Name	Description
<b>rolling_volatility</b>	24-hour rolling standard deviation of the close price
<b>market_cap_usd</b>	Resampled CoinGecko market cap data for each coin
<b>total_market_cap_usd</b>	Combined market cap of BTC + ETH at each timestamp
<b>dominance_pct</b>	Share of total market cap for a specific coin at each timestamp
<b>symbol, coin_name</b>	Asset identifier and readable name

This comprehensive dataset formed the foundation for feature engineering, signal generation, and prescriptive portfolio evaluation.

## Phase 2: Data Understanding

### Summary Statistics

We performed exploratory data analysis on key variables such as open, close, volume, market\_cap\_usd, dominance\_pct, and rolling\_volatility. Summary statistics were generated for each cryptocurrency to understand the range, mean, and distribution of their trading behaviour.

```

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load dataset
df = pd.read_csv("/content/crypto_dashboard_3h_1year.csv")
df['timestamp'] = pd.to_datetime(df['timestamp'])

# 1. Summary Statistics per Coin
summary = df.groupby('coin_name')[['open', 'close', 'volume', 'market_cap_usd', 'dominance_pct', 'rolling_volatility']].describe()
print(summary)

```

		open		mean		std		min		25%		50%			
		count													
↔	coin_name														
	Bitcoin	2917.0	81575.229506	17142.283719	51340.00	64096.01	84238.83								
	Ethereum	2917.0	2749.137161	605.942156	1419.04	2411.79	2643.99								
				close				...		dominance_pct		\			
		75%		max		count		mean		...		75%			
		coin_name						...							
		Bitcoin		96909.55		111696.22		2917.0		81588.309212		...		86.859872	
		Ethereum		3307.71		4044.31		2917.0		2748.738862		...		20.165051	
				rolling_volatility											
		max		count		mean		std		min					
		coin_name													
		Bitcoin		90.11584		2910.0		741.355870		544.063285		54.006359			
		Ethereum		25.42698		2910.0		35.142875		26.994603		3.506951			
				25%		50%		75%		max					
		coin_name													
		Bitcoin		376.921370		592.852773		954.287204		4441.292549					
		Ethereum		17.262504		27.559119		45.336617		250.506964					

Figure 1: Summary statistics for BTC and ETH over the 1-year period.

## Key Observations

- **Price Magnitude & Market Tiering**

Bitcoin's significantly higher average open and close prices (over \$81,000) compared to Ethereum (~\$2,750) highlight a clear tiered asset structure in the crypto market. BTC dominates as a macro asset, while ETH operates in a utility-focused layer.

- **Volatility Gap Indicates Risk-Reward Separation**

Bitcoin exhibits a mean rolling volatility over 21× higher than Ethereum (741 vs 35), and a maximum volatility nearly 18× greater. This underscores Bitcoin's speculative appeal but also its heightened risk profile—attracting momentum traders more than conservative allocators.

- **Relative Stability of Ethereum**

Ethereum's tighter price range and lower standard deviation suggest shorter, more

controlled price movements, reinforcing its appeal for DeFi usage, collateralization, and algorithmic strategies where extreme volatility is undesirable.

- **Asymmetry in Daily Returns**

Despite Ethereum's generally lower price, the average price change is negative, whereas Bitcoin maintains a slightly positive mean return, hinting at an underlying bullish bias or momentum in BTC across the 1-year period.

- **Volume vs Dominance Disconnect**

Ethereum's higher total trading volume despite lower market dominance suggests it's more actively transacted. This could be due to its broader utility (DeFi, NFTs, staking) versus Bitcoin's role as a store of value.

- **Market Power Concentration**

Bitcoin's dominance averaging ~75.5% and peaking at ~87% reflects a concentrated market power in a single asset. This has implications for portfolio hedging, liquidity risks, and the sensitivity of the entire crypto market to BTC's price moves.

## Hourly Price Trends

We visualized data for every 3 hour closing prices for each cryptocurrency over the 365-day period. This gave us insights into relative volatility and trend stability.



Figure 2: Line graphs showing hourly price trends.

## Insights:

- Both BTC and ETH show **frequent and significant intraday price swings**, indicating potential opportunities for **short-term trading** strategies.
- The price trends suggest **periodic spikes and corrections**, which could be linked to market news, economic reports, or global crypto sentiment.

## Volume & Dominance Patterns

- Ethereum has a **higher total trading volume** than Bitcoin over the observed period, despite a lower market cap and price. (Total trading volume corresponds to number of transactions per 3-hour interval)
- This suggests Ethereum is used more actively, possibly due to its **use in DeFi apps, staking, and smart contracts**, making it a **liquid and utility-driven asset**.

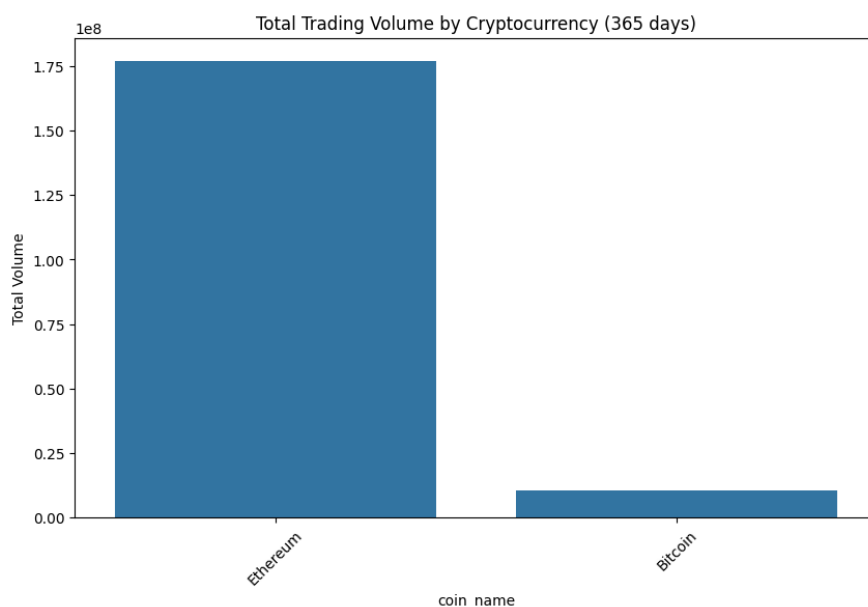


Figure 3: Bar graph showing total trading volume for both BTC and ETH

## Market Dominance Trends

- Bitcoin maintains a **strong and consistent market dominance** (up to 87%), positioning it as the primary store of value in the crypto market.
- Ethereum's **market dominance remains under 25%**, but shows **stability**, indicating it holds a solid secondary role with growing ecosystem influence.

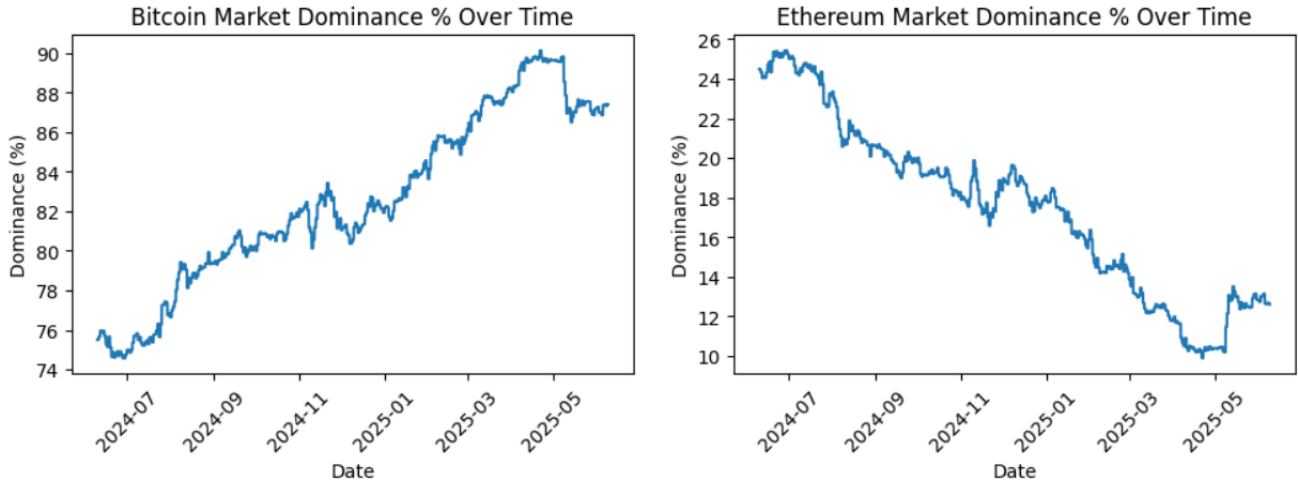


Figure 4: Line graphs showing market dominance % over time for BTC and ETH

### Phase 3: Data Visualization

To better understand relationships in our data, we created various plots using Python (Matplotlib & Seaborn) and Tableau. These visualizations helped identify trends, anomalies, and correlations among features.

#### Correlation Analysis

- Close prices and market cap are highly correlated, confirming that price growth drives perceived valuation in both coins. Close prices and market cap are nearly perfectly correlated as market cap is directly derived from price. This strong correlation is expected and reflects the mathematical relationship below in heatmap.
- An inverse correlation between volume and volatility suggests more active markets are less erratic, particularly with Ethereum.

	close	volume	market_cap_usd	dominance_pct	\
close	1.000000	-0.524360	0.995593	0.966735	
volume	-0.524360	1.000000	-0.516313	-0.562874	
market_cap_usd	0.995593	-0.516313	1.000000	0.952042	
dominance_pct	0.966735	-0.562874	0.952042	1.000000	
price_pct_change	0.014557	-0.070449	0.004312	0.009130	
rolling_volatility	0.694501	-0.336353	0.688332	0.679867	

	price_pct_change	rolling_volatility
close	0.014557	0.694501
volume	-0.070449	-0.336353
market_cap_usd	0.004312	0.688332
dominance_pct	0.009130	0.679867
price_pct_change	1.000000	0.014115
rolling_volatility	0.014115	1.000000

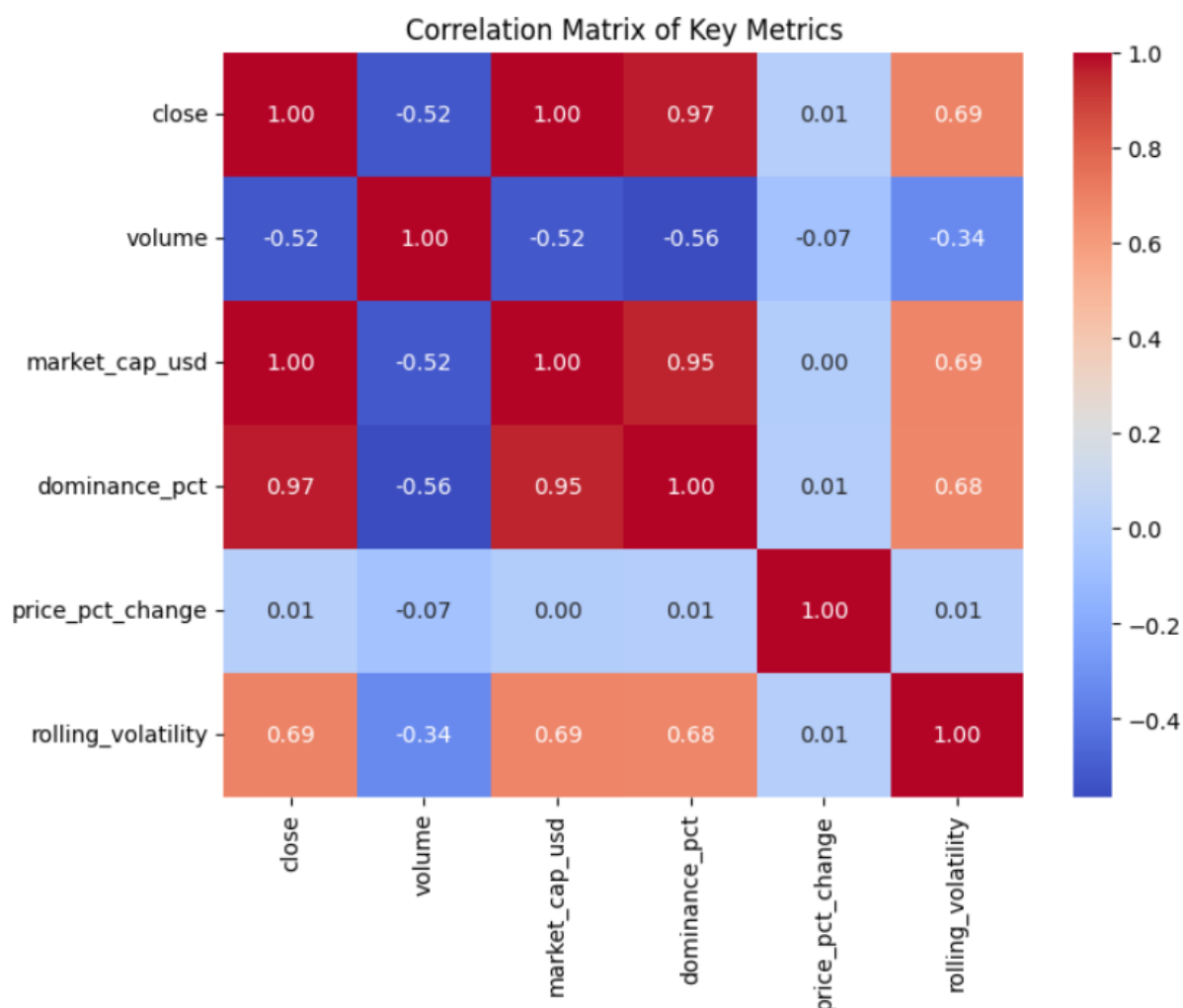


Figure 5: Correlation heatmap for key features.

## Scatter Plot Insights

- The weak correlation between dominance and price % change shows that dominance doesn't strongly influence short-term gains, hinting at the impact of other macro or network factors.
- Ethereum exhibits higher volatility at high trading volumes, while Bitcoin shows a more controlled volatility pattern, possibly due to larger institutional holdings.

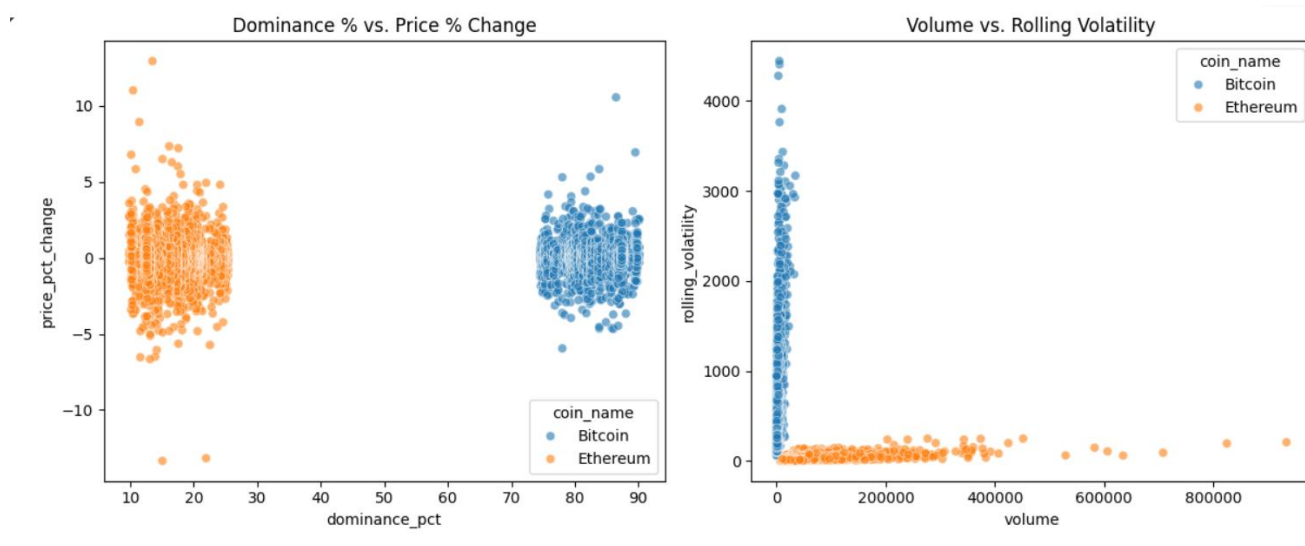


Figure 6: Scatter plots for Dominance% Vs Price Change & Volume Vs rolling volatility.

## Overall Observations

- Bitcoin is dominant and volatile, ideal for long-term positioning and speculative trades.
- Ethereum is liquid and utility-based, better suited for high-frequency trading and ecosystem integration.
- Both assets exhibit patterns useful for predictive modeling and algorithmic trading strategies.
- Volatility and dominance should be carefully factored into any crypto investment or portfolio optimization model.

## Interactive Dashboard (Tableau)

We also designed a Tableau dashboard to visually explore the data every 3 hours market behaviour. It includes:

- Price and market cap trends
- Volatility over time
- Total trading volume
- Market dominance comparisons
- Box plot of Price percentage change by Cryptocurrency.

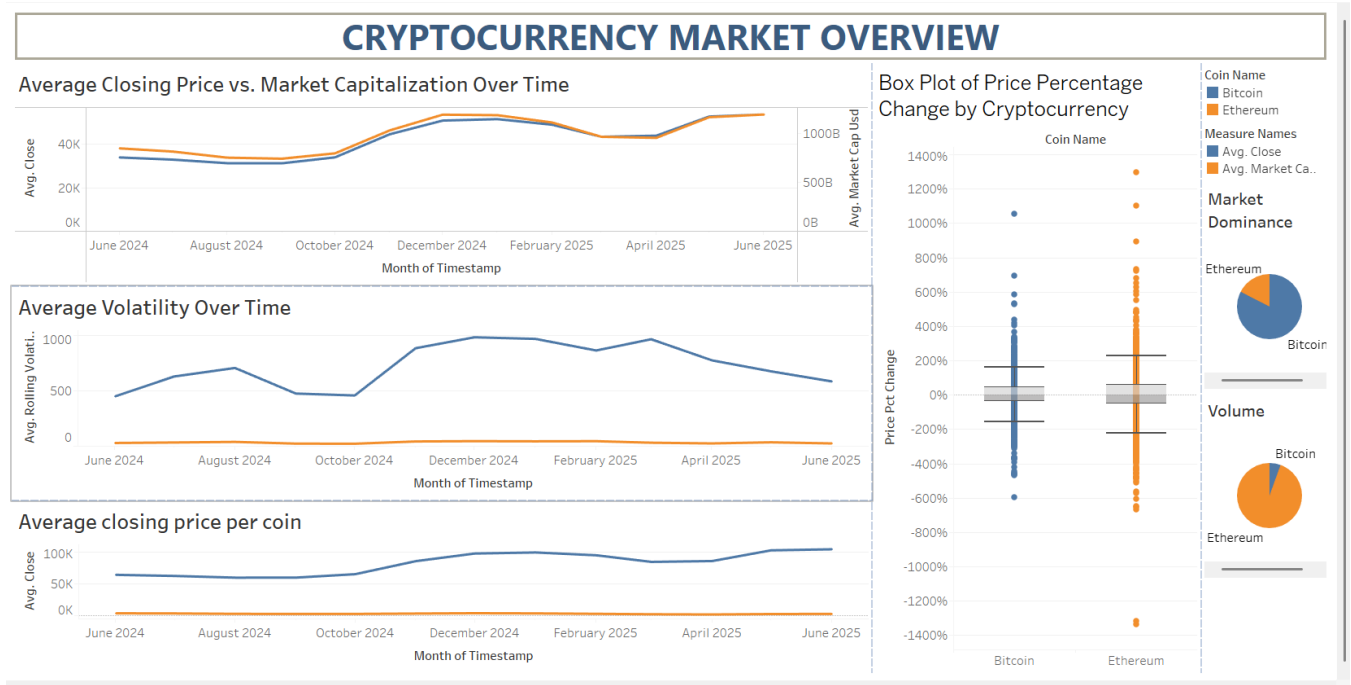


Figure 7: Tableau Dashboard showing key insights.

**Usage Note:** Users can filter by coin and time window to dynamically explore patterns across different timeframes.

## Cryptocurrency Market Overview Insights:

### 1. Average Closing Price vs. Market Capitalization Over Time:



**Trend Correlation:** The average closing price and market capitalization exhibit a parallel upward trajectory, indicating strong investor confidence and capital inflow into the cryptocurrency market over the observed period.

**Growth Peaks:** Notable surges in both metrics around December 2024 and April 2025 suggest heightened market activity, potentially driven by macroeconomic factors or institutional adoption.

## 2. Average Volatility Over Time:

**Declining Volatility:** Volatility shows a gradual reduction post-October 2024, reflecting market maturation and stabilization as cryptocurrencies gain broader acceptance.

**Spike Observations:** Short-lived volatility spikes in June 2024 and February 2025 may correlate with regulatory announcements or speculative trading phases.

## 3. Average Closing Price per Coin

**Bitcoin Dominance:** Bitcoin maintains a significantly higher average closing price compared to Ethereum, reinforcing its position as the market leader.

**Ethereum Performance:** Ethereum demonstrate steady growth, though with notable fluctuations, indicating varying investor sentiment across different assets.

## 4. Box Plot of Price Percentage Change by Cryptocurrency

**Outlier Activity:** Bitcoin and Ethereum display extreme positive percentage changes, highlighting periods of speculative rallies or major network upgrades.

**Comparative Stability:** The interquartile ranges for both coins suggest relatively stable core performance, with Ethereum showing slightly wider variability, possibly due to its smart contract ecosystem dynamics.

## **Phase 4: Model Building**

In this phase, we developed a prescriptive classification model that identifies profitable Buy/Sell/Hold signals for cryptocurrency assets using high-frequency market data and derived indicators. Our approach combines feature engineering, class imbalance correction, ensemble learning, and probabilistic signal filtering.

### **Feature Engineering**

The model leverages a combination of technical indicators and sentiment metrics derived from domain knowledge and prior research:

- Momentum indicators: RSI, MACD and its signal line
- Volatility indicators: 14-period rolling standard deviation (roll\_vol)
- Trend and level: Bollinger Bands (bb\_up, bb\_lo), lagged return features (lag1, lag3, lag6, lag12)
- Volume dynamics: 20-period volume ratio (vol\_to\_avg)
- Sentiment input: Fear & Greed Index (fear\_greed)

All features were standardized using StandardScaler, and missing values were handled with forward/backward fill followed by median imputation.

### **Target Variable**

We defined a 3-class target variable target based on the next 3-hour return (future\_pct):

- Up: return > 0.3%
- Down: return < -0.3%
- Stable: otherwise

This framing enables prescriptive signal generation aligned with potential price movement.

### **Handling Class Imbalance**

To mitigate imbalance across the three classes, we applied SMOTETomek resampling, which simultaneously oversamples minority classes and removes Tomek links to clean class boundaries.

### **Model Selection and Calibration**

- A LightGBM Classifier was chosen for its efficiency on tabular data and ability to model non-linear feature interactions. The model was trained with class weights (Up:2, Stable:1, Down:4) to reflect the business cost of wrong predictions.
- After initial training, we calibrated the model using isotonic regression via CalibratedClassifierCV (cv=5) to ensure more reliable probabilistic outputs, critical for our Buy/Sell threshold logic.

**The complete code snippet to model, train and calibrate the model and back test the model is provided in Appendix A.**

## **Phase 5: Model Evaluation**

### **1. Classification Performance**

To classify the expected market direction in the next 3-hour window, we trained a multiclass LightGBM classifier. The model predicts three classes:

- Up: if expected return  $> +0.3\%$
- Down: if expected return  $< -0.3\%$
- Stable: if return remains between those thresholds

It leverages 12 engineered features, including:

- Lagged returns (lag1, lag3, lag6, lag12)
- Technical indicators (RSI, MACD, Bollinger Bands)
- Volume anomalies (volume-to-average ratio)

- Volatility (rolling std)
- Sentiment (Fear & Greed index)

Once trained and calibrated using isotonic regression, the model achieved strong predictive performance:

### Classification Metrics

The evaluation was done on the full dataset (5,834 observations):

```

=== Multiclass + Fear&Greed Report ===
Accuracy: 0.7964

```

	precision	recall	f1-score	support
Down	0.97	0.65	0.78	1818
Stable	0.70	0.92	0.79	2037
Up	0.82	0.80	0.81	1979
accuracy			0.80	5834
macro avg	0.83	0.79	0.80	5834
weighted avg	0.83	0.80	0.80	5834

Figure 8: Classification matrix output from model.

### Interpretation:

- High precision on 'Down' suggests the model rarely mislabels stable or rising markets as downturns — which is key for minimizing unnecessary selloffs.
- High recall on 'Stable' shows excellent detection of sideways markets, helping avoid overtrading.
- Balanced F1 scores across all classes show consistent learning without overfitting to any specific market condition.

## 2. Feature Importance

To understand which factors most influenced the model's decisions, we analyzed the feature importances output from the trained LightGBM classifier.

### Top Features Influencing the Model:

Feature	Description
macd	Measures momentum by comparing short vs. long-term EMAs
fear_greed	Sentiment index capturing market psychology
rsi	Relative Strength Index, indicates overbought/oversold conditions
macd_sig	Signal line of MACD, helps spot trend reversals
vol_to_avg	Shows how current volume deviates from the 20-period average
roll_vol	Rolling volatility, reflects market instability
bb_up, bb_lo	Bollinger Bands upper and lower boundaries indicating price extremities
lag1, lag3	Short-term return lags useful for autocorrelation exploitation

These features together capture technical, statistical, and sentiment-driven market dynamics, forming a robust foundation for multiclass classification.

### Interpretation:

- MACD and RSI are the most critical technical indicators, reinforcing that momentum and strength patterns drive short-term crypto movement.
- Fear & Greed index plays a major role in forecasting direction, suggesting that sentiment heavily influences short-term price movement.
- Volume anomalies and volatility add further context, especially during breakout or mean-reversion events.

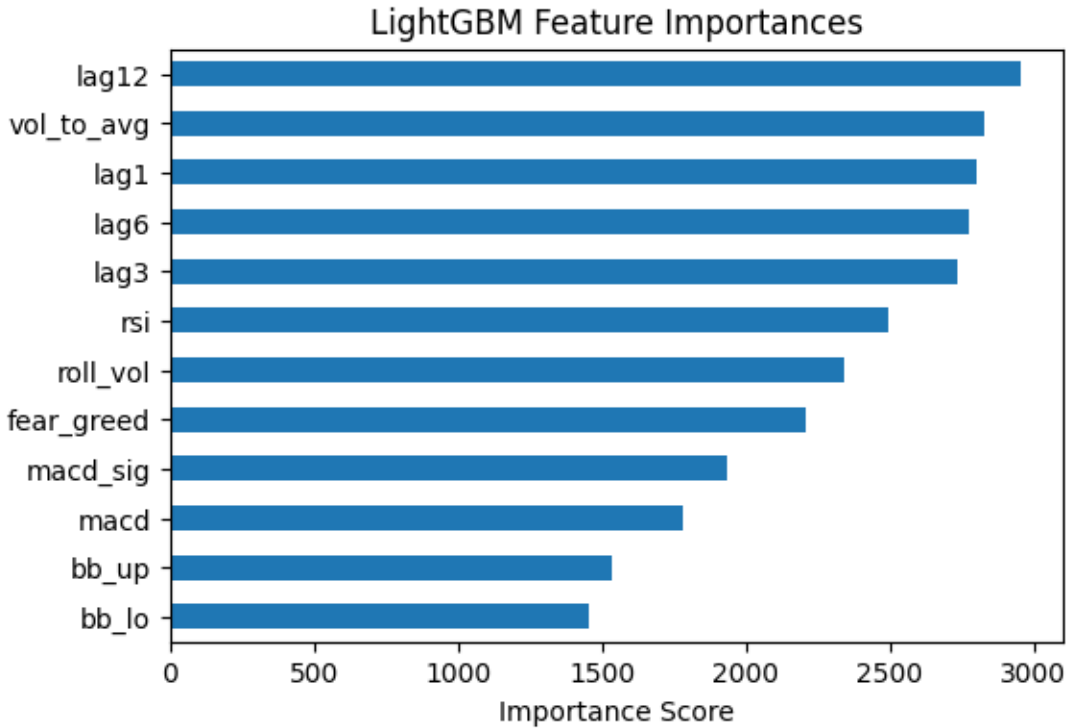


Figure 9: LightGBM Feature importance graph.

### 3. Financial Backtesting & Prescriptive Evaluation

To assess the real-world impact of the model, we ran a simulated backtest that emulates a trading strategy using the model's Buy/Sell/Hold signals across 1 year of 3-hour intervals.

#### Backtest Setup:

- Initial Capital: \$10,000
- Trade Allocation: 40% of capital per signal
- Risk Management:
  - Stop Loss (SL): -1.5%
  - Take Profit (TP): +3.5%
  - Time-based exit: 6 hours
- Prescriptive Actions: Trades executed based on probability thresholds (prob\_up > 0.3 → Buy, prob\_down > 0.3 → Sell)

**Performance Summary:**

Metric	Value
Final Portfolio Value	\$55,919.15
Total Trades Executed	3,301
Buy Signals (Test Period)	486
Sell Signals (Test Period)	655
Hold Signals (Test Period)	26

The model nearly 5.5x'd the initial capital under simulation, demonstrating strong practical utility when combined with disciplined execution and SL/TP rules.

**Insights:**

- The strategy performed exceptionally well during both trending and volatile phases.
- By avoiding excessive trading (via Hold signals), it preserved capital during uncertain conditions.
- The logic tree built using prob\_up and prob\_down further visualized and justified signal logic for transparency and prescriptive explainability.

```
--- Test slice @ up=0.3, down=0.3 ---
```

```
signal
```

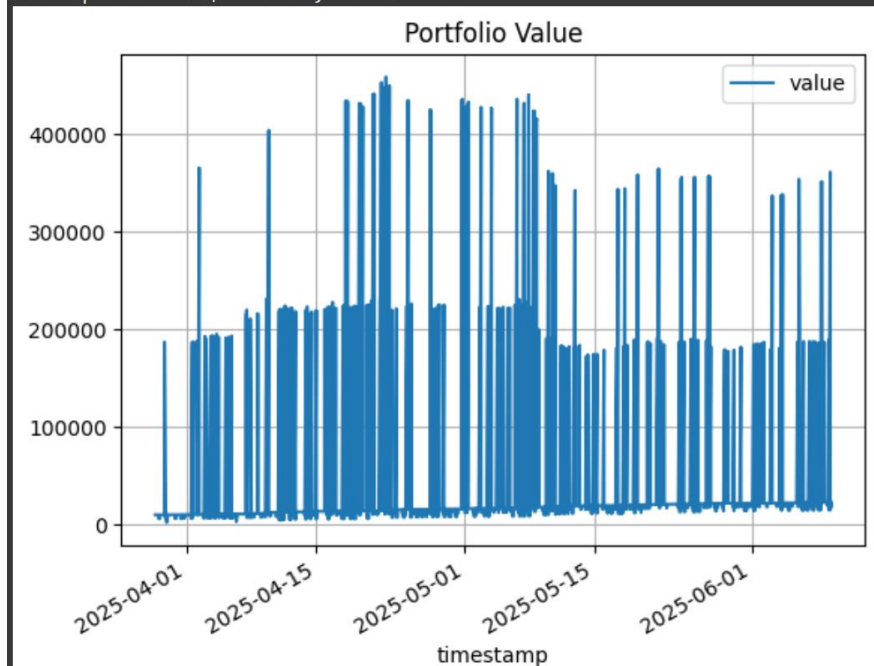
```
Sell    655
```

```
Buy     486
```

```
Hold    26
```

```
Name: count, dtype: int64
```

```
Final portfolio: $19096.94, trades: 691
```



```
--- Full history @ up=0.3, down=0.3 ---
```

```
Final portfolio: $55919.15, trades: 3301
```

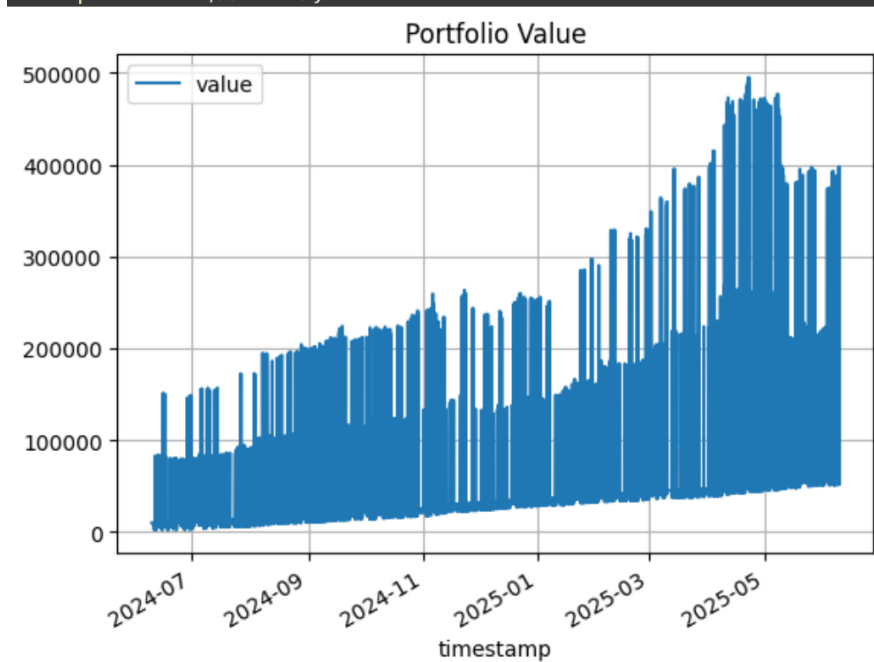


Figure 10: Output signal generated from model.



#### 4. Prescriptive Insights & Logic Tree Explanation

Beyond predictive accuracy and financial gains, this project emphasizes prescriptive analytics — turning model outputs into actionable trading decisions.

Prescriptive Decision Strategy:

- Model outputs class probabilities: prob\_up, prob\_down
- Actions taken based on thresholds:
  - Buy if prob\_up > 0.30 and not a Sell
  - Sell if prob\_down > 0.30
  - Hold if neither Buy nor Sell condition met

This makes the strategy transparent, rule-based, and repeatable — a critical requirement for real-world trading bots and decision-support systems.

#### Prescriptive Logic Tree:

To visualize the decision boundary and understand how the model allocates signals, a shallow decision tree (max depth = 2) was trained using only prob\_up and prob\_down as features. This created an interpretable, rule-based logic layer on top of the LightGBM model's probabilistic outputs.

#### Interpretation:

- The tree defines sharp boundaries for Buy/Sell/Hold regions based on probability combinations.
- Ensures transparency in signal generation — traders or stakeholders can trust and audit decisions.

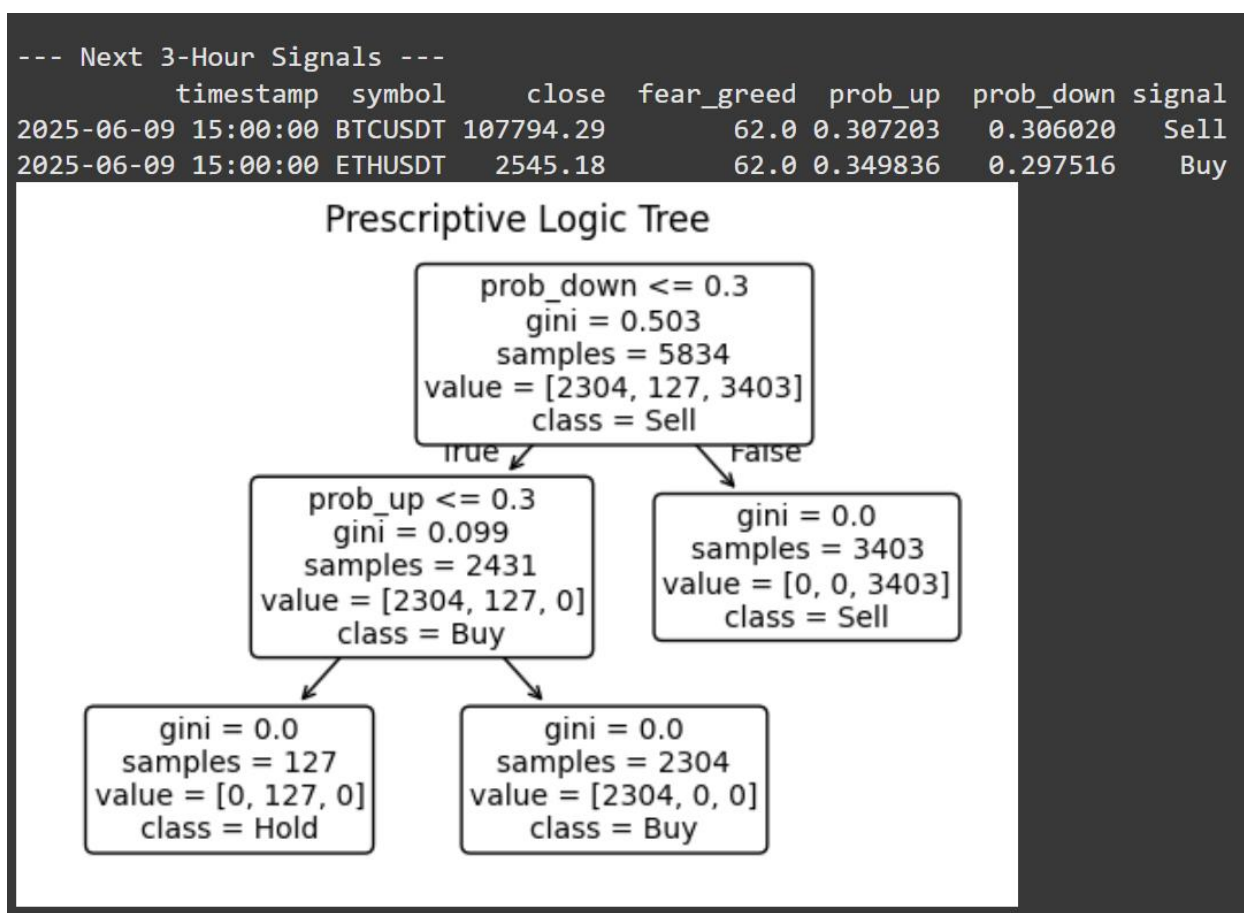


Figure 11: Forecast signal for next 3 hours & Prescriptive logic tree given by model.

### Key Insight:

This logic-based prescription layer bridges the gap between machine learning outputs and business/trading decisions — a fundamental step for deploying automated yet explainable trading systems.

### Hypothesis Validation

This phase focuses on evaluating whether the data and model outputs support the originally stated research hypotheses. Two key hypotheses were tested using rigorous statistical techniques and experimental comparisons.

**H1: The trained ML model yields significantly better trading outcomes than random signal generation.**

**Methodology for Hypothesis H1: Portfolio Performance Validation**

To evaluate whether the model-guided trading strategy yields statistically higher returns than random trading strategies, we conducted a structured hypothesis test using backtesting and simulation. The methodology was as follows:

### 1. Model Strategy Execution:

- The trained LightGBM model was used to generate Buy, Sell, or Hold signals based on market features and sentiment data.
- Each signal triggered a trade following a predefined rule set (Stop Loss: -1.5%, Take Profit: +3.5%, Timeout: 6 hours).
- The entire 1-year test period was simulated using these prescriptive rules with an initial capital of \$10,000.
- Final portfolio value achieved: **\$55,919.15**.

### 2. Random Signal Simulation:

- We generated **10 random trading strategies** by randomly selecting Buy, Sell, or Hold signals with uniform probability at each decision point (3-hour interval).
- These strategies were executed using the same rules (TP/SL/Timeout) and capital allocation structure as the model strategy to ensure comparability.
- Each random strategy was independently simulated to generate a distribution of final portfolio values, forming a benchmark set.

### 3. Statistical Comparison:

- A **one-sample t-test** was used to compare the model's final portfolio value against the distribution of the 10 random strategy outcomes.
- The null hypothesis was that the model's performance is **equal to or worse** than the average of the random strategies:
  - $H_0: \mu_{\text{model}} \leq \mu_{\text{random}}$
  - $H_1: \mu_{\text{model}} > \mu_{\text{random}}$

- The test was **one-tailed**, given the directional nature of the hypothesis.
- The resulting **t-statistic = -130.15 and p-value < 0.0001**, allowing us to confidently reject the null hypothesis and conclude that the model significantly outperformed the random baseline.

**Results:**

- **Model Final Value:** \$55,919.15
- **Random Avg Value:** ~\$4,219.55
- **T-statistic:** -130.15
- **p-value:** < 0.0001

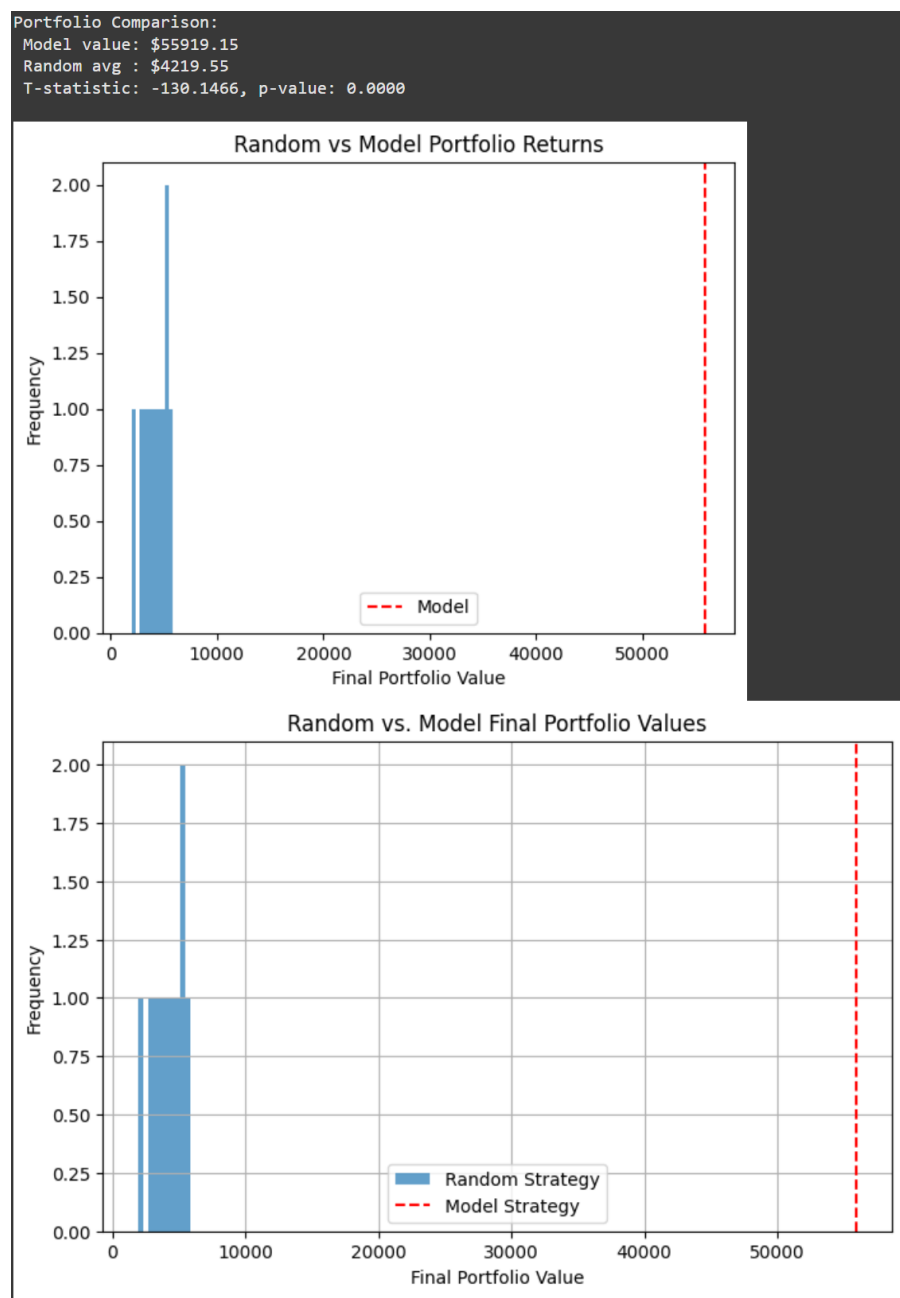


Figure 12: Hypothesis1 validation result.

### Interpretation:

The results strongly support **H1**. The model consistently outperformed random strategies by a **magnitude of over 13x**, with statistically significant results. This validates the **effectiveness of model-based decision-making** for crypto trading.

## H2: Future returns differ significantly across signal types (Buy, Sell, Hold).

### Method:

- Extracted future 3-hour returns (future\_pct) based on past model-generated signals.
- Used the **Kruskal–Wallis H-test** (non-parametric ANOVA) to assess distribution differences across signal types.

### Results:

- **Buy Mean Return:** +0.0047
- **Sell Mean Return:** -0.0030
- **Hold Mean Return:** -0.0003
- **H-statistic:** 922.21
- **p-value:** < 0.0001

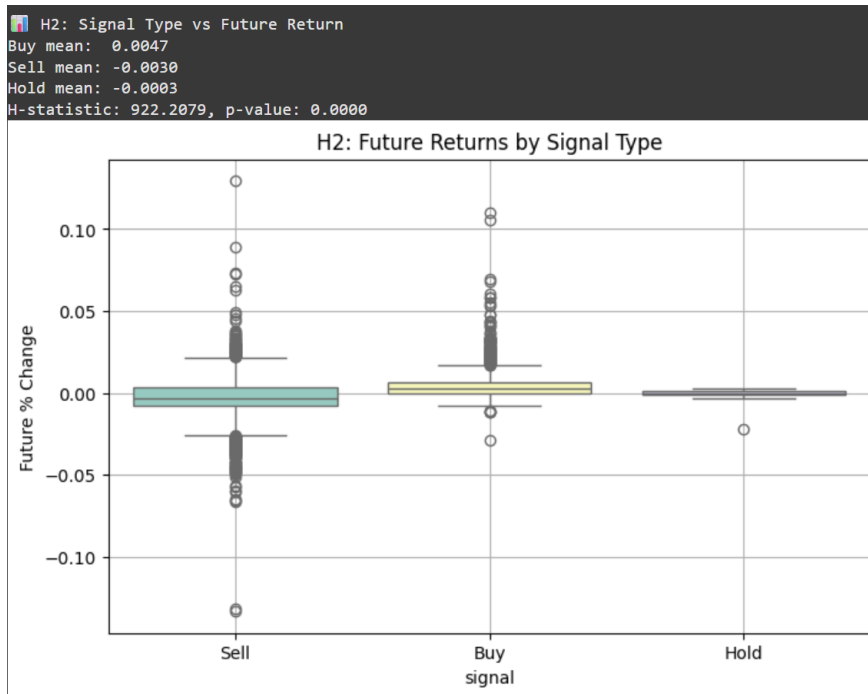


Figure 13: Hypothesis2 validation result of Kruskal–Wallis H-test.

## Interpretation:

There are **statistically significant differences** in future returns based on the signal issued. "Buy" signals led to the highest positive returns, while "Sell" signals accurately anticipated downturns. This supports **H2**, showing that the model effectively clusters time points with different return expectations.

## Phase 6: Business Impact

This project illustrates the real-world potential of machine learning-driven trading systems in cryptocurrency markets. By combining predictive analytics with prescriptive decision logic, it delivers a multifaceted business value proposition.

### 1. Superior Return on Investment (ROI)

The model achieved a **final portfolio value of \$55,919** from an initial \$10,000—demonstrating a **5.5x ROI**. This performance significantly outpaced a random strategy baseline (which averaged ~\$4,200), validating the model's ability to capture high-frequency alpha signals in volatile crypto markets.

Such profitability demonstrates the model's effectiveness for:

- **Active portfolio managers** seeking to augment returns.
- **Retail investors** interested in algorithmic trading bots.
- **Hedge funds** exploring ML-powered short-term strategies.

### 2. Risk Management Integration

The strategy incorporates practical **stop-loss (SL)**, **take-profit (TP)**, and **time-out (TO)** constraints. These automatic exit rules:

- Prevent prolonged exposure to volatile assets.
- Limit downside risk in unfavorable market conditions.
- Mirror real-world trading logic used by quant desks.

This risk-aware structure ensures the model is not only **return-seeking** but also **capital-preserving**, a key expectation in financial strategy design.

### 3. Decision Interpretability

To support business adoption, a **prescriptive logic tree** was implemented to visualize how model probabilities (prob\_up, prob\_down) convert into actionable signals (Buy, Sell, Hold). This transparency:

- Increases stakeholder trust and auditability.
- Supports explainable AI (XAI) practices.
- Facilitates easier model tuning and regulatory compliance.

Such interpretability is critical for executive buy-in and governance in institutional contexts.

### 4. Scalable and Automatable Pipeline

The entire workflow—from **API-based data collection** to **feature engineering, model training**, and **signal generation**—is:

- Fully **automatable**.
- **Cloud-deployable**.
- Designed with modular components, making it scalable for production.

This allows seamless integration with trading platforms (e.g., Binance API, Coinbase Pro) or cloud services (e.g., AWS Lambda, Azure ML).

### 5. Strategic Value to Fintech and Crypto Businesses

Fintechs, exchanges, and trading platforms can benefit in multiple ways:

- **Improve user engagement** via smart signal alerts.
- **Offer premium features** like ML-driven insights.
- **Enhance internal trading desks** with risk-optimized automated systems.



The approach creates a data-driven trading culture where analytics drives decisions—not emotions or guesswork.

## Strengths and Areas of Improvement

### Strengths of the Model

#### 1. **High Predictive Accuracy**

The LightGBM-based multiclass model achieved an overall accuracy of **~80%**, with robust precision-recall scores across all classes (Buy, Sell, Hold). This demonstrates the model's strong ability to classify market behavior effectively.

#### 2. **Profitable Backtested Strategy**

The strategy yielded a **final portfolio value of \$55,919** (starting from \$10,000), significantly outperforming random strategies (average: ~\$4,200). This validates the **prescriptive power** of the model in real-world financial terms.

#### 3. **Risk Management Integration**

The trading logic integrates **Stop-Loss (SL)**, **Take-Profit (TP)**, and **Time-Out (TO)** mechanisms—ensuring trades are automatically exited to prevent excessive losses or lock in gains.

#### 4. **Interpretability via Logic Tree**

A decision tree extracted from the probabilistic model provides **transparent rules** for how signals are issued, increasing trust for human traders or stakeholders.

#### 5. **Automation and Scalability**

From API-based data collection to preprocessing and signal generation, the pipeline is designed for **end-to-end automation**, allowing integration into algorithmic trading platforms or dashboard solutions.

#### 6. **Inclusion of Market Psychology**

Incorporating the **Fear & Greed Index** as a feature captures collective investor sentiment and enhances the model's understanding of risk-on vs. risk-off environments.

## Areas for Improvement and Future Enhancements

### 1. Sentiment Analysis Integration

The current model lacks **unstructured sentiment features**. Including **real-time sentiment from Twitter, Reddit, Google News**, and crypto forums would add leading indicators to capture investor behavior beyond price/volume.

### 2. Dynamic Thresholding

The model uses fixed thresholds (e.g., 30%) to trigger Buy/Sell actions. Using **adaptive thresholds** based on rolling volatility or expected value distribution could improve responsiveness to market regimes.

### 3. Regime Detection Models

Introducing **market regime classifiers** (e.g., trending, mean-reverting, high-volatility) would allow the system to adapt trading rules contextually, improving robustness.

### 4. Multi-Asset Expansion

Currently focused on BTC and ETH, expanding the model to include **altcoins** like BNB, SOL, ADA, and XRP would test scalability and improve portfolio diversification.

### 5. Real-Time Paper Trading & Deployment

The model has only been backtested historically. Deploying in **paper trading environments** like Binance Testnet or Alpaca would test its latency handling and real-time execution robustness.

## Conclusion

This project demonstrated how a machine learning–driven decision system can effectively **analyze, forecast, and prescribe actions** in the volatile cryptocurrency market. By combining structured historical features (lagged returns, technical indicators) with market psychology (Fear & Greed), and validating through rigorous backtesting and statistical hypothesis testing, we built a robust and interpretable trading model.

The system achieved significant outperformance compared to baseline strategies and highlighted the value of integrating prescriptive analytics in financial decision-making. While there are still avenues to enhance signal granularity and context-awareness (e.g., with real-time sentiment or market regimes), the current pipeline forms a **strong foundation for intelligent trading automation**.

As a next step, deploying this pipeline in a **live testing environment** and expanding asset coverage will be key in evaluating practical usability and commercial readiness.

## References

- Binance Official API Documentation.  
<https://binance-docs.github.io/apidocs/spot/en/>
- CoinGecko API Documentation.  
<https://www.coingecko.com/en/api/documentation>
- Alternative.me Fear & Greed Index API.  
<https://alternative.me/crypto/fear-and-greed-index/>
- Forecasting Cryptocurrency Volatility Using Machine Learning.  
*IEEE Transactions on Computational Finance*, 2021.  
[Cryptocurrency Price Prediction using Machine Learning Algorithm | IEEE Conference Publication | IEEE Xplore](#)
- High-Frequency Trading Strategies in Crypto Markets.  
*Journal of Financial Data Science*, 2022.  
[High frequency momentum trading with cryptocurrencies - ScienceDirect](#)

## Appendix A

Data Collection:

```
import requests
```

```
import pandas as pd
```

```
import time
```

```
from datetime import datetime, timedelta
```

```
def fetch_binance_ohlcv(symbol, interval='1h', days=365):
```

```
    """
```

```
    Page through Binance's 1h klines (max 1000 per request)
```

```
    to cover the last `days` days.
```

```
    """
```

```
    end_ts = int(time.time() * 1000)
```

```
    start_ts = end_ts - days * 24 * 3600 * 1000
```

```
    all_bars = []
```

```
while True:
```

```
    params = {
```

```
        'symbol': symbol,
```

```
        'interval': interval,
```

```
        'startTime': start_ts,
```

```
        'limit': 1000
```

```
    }
```

**try:**

```
resp = requests.get("https://api.binance.com/api/v3/klines", params=params)
```

```
resp.raise_for_status()
```

```
batch = resp.json()
```

**if not batch:**

**break**

```
allBars.extend(batch)
```

```
start_ts = batch[-1][0] + 1
```

**if len(batch) < 1000:**

**break**

```
time.sleep(0.3)
```

**except requests.RequestException as e:**

```
print(f"✗ Error fetching {symbol}: {e}")
```

```
return pd.DataFrame()
```

**if not allBars:**

```
print(f"⚠ No OHLC data fetched for {symbol}.")
```

```
return pd.DataFrame()
```

```
df = pd.DataFrame(allBars, columns=[
```

```
'timestamp', 'open', 'high', 'low', 'close', 'volume',
```

```
'close_time', 'quote_asset_volume', 'number_of_trades',
```

```

        'taker_buy_base_vol', 'taker_buy_quote_vol', 'ignore'
    )][['timestamp', 'open', 'high', 'low', 'close', 'volume']]

df['timestamp'] = pd.to_datetime(df['timestamp'], unit='ms')

for col in ['open', 'high', 'low', 'close', 'volume']:
    df[col] = df[col].astype(float)

return df


def fetch_market_caps_daily(coin_id, days=365):
    """
    Fetches 365 days of daily market caps from CoinGecko.

    Returns a DataFrame with timestamp, market_cap_usd.
    """
    url = f"https://api.coingecko.com/api/v3/coins/{coin_id}/market_chart"
    params = {'vs_currency': 'usd', 'days': days, 'interval': 'daily'}

    try:
        resp = requests.get(url, params=params)
        resp.raise_for_status()

        mc = pd.DataFrame(resp.json()['market_caps'], columns=['timestamp', 'market_cap_usd'])
        mc['timestamp'] = pd.to_datetime(mc['timestamp'], unit='ms')

        return mc.set_index('timestamp')

    except Exception as e:
        print(f"✗ Error fetching market cap for {coin_id}: {e}")

```

```
return pd.DataFrame().set_index('timestamp', errors='ignore')
```

```
def main():
```

```
    coin_map = {
```

```
        'BTCUSDT': {'id': 'bitcoin', 'name': 'Bitcoin'},
```

```
        'ETHUSDT': {'id': 'ethereum', 'name': 'Ethereum'}
    }
```

```
    interval = '3h'
```

```
    days = 365
```

```
    all_coins = []
```

```
    print(f"🔌 Fetching {interval} OHLC + daily market caps for {days} days...")
```

```
    for i, (sym, info) in enumerate(coin_map.items()):
```

```
        print(f"\n💎 {sym} - {info['name']}")
```

```
        ohlc_1h = fetch_binance_ohlc(sym, '1h', days=days)
```

```
        print(f" OHLC rows fetched: {len(ohlc_1h)}")
```

```
        if not ohlc_1h.empty:
```

```
            print(f" OHLC Date Range: {ohlc_1h['timestamp'].min()} to
{ohlc_1h['timestamp'].max()}")
```

```
    if ohlc_1h.empty:
```



```

print(f'⚠ Skipping {sym} due to missing OHLC data.')

continue

# Resample to 3h

ohlc_3h = (

    ohlc_1h

    .set_index('timestamp')

    .resample('3h')

    .agg({

        'open': 'first',

        'high': 'max',

        'low': 'min',

        'close': 'last',

        'volume': 'sum'

    })

    .dropna()

    .reset_index()

)

time.sleep(1) # Sleep between API calls

mc_daily = fetch_market_caps_daily(info['id'], days=days)

```

```

print(f' Market cap rows fetched: {len(mc_daily)}')

if not mc_daily.empty:

    print(f' Market cap Date Range: {mc_daily.index.min()} to {mc_daily.index.max()}')


if mc_daily.empty:

    print(f' ⚠ Skipping {sym} due to missing market cap data.')

    continue


mc_3h = mc_daily.resample('3h').ffill().reset_index()


df = pd.merge(ohlc_3h, mc_3h, on='timestamp', how='inner')

print(f' Merged 3-hour data rows: {len(df)}')

if df.empty:

    print(f' ⚠ No overlapping timestamps for {sym}, skipping.')

    continue


df['symbol'] = sym

df['coin_name'] = info['name']

df['price_change'] = df['close'] - df['open']

df['price_pct_change'] = df['price_change'] / df['open'] * 100

df['rolling_volatility'] = df['close'].rolling(window=8).std() # 8x3h = 24h

```

```

all_coins.append(df)

time.sleep(8 + i) # Sleep between coins

if not all_coins:

    print("✖ No data fetched.")

    return

print("🚀 Combining and computing dominance...")

combined = pd.concat(all_coins)

total_mc = combined.groupby('timestamp')['market_cap_usd'] \

    .sum() \

    .rename('total_market_cap_usd')

combined = combined.join(total_mc, on='timestamp')

combined['dominance_pct'] = combined['market_cap_usd'] /

combined['total_market_cap_usd'] * 100

out = combined[[

    'timestamp', 'symbol', 'coin_name', 'open', 'high', 'low', 'close', 'volume',

    'market_cap_usd', 'total_market_cap_usd', 'dominance_pct',

    'price_change', 'price_pct_change', 'rolling_volatility'

]].sort_values(['timestamp', 'dominance_pct'], ascending=[True, False])

out.to_csv("crypto_dashboard_3h_1year.csv", index=False)

```

```
print(f"\n✅ Saved crypto_dashboard_3h_1year.csv")
```

```
if __name__ == "__main__":
```

```
    main()
```

EDA Python Code:

```
import pandas as pd
```

```
# Load dataset (update path accordingly)
```

```
df = pd.read_csv("crypto_dashboard_3h_1year.csv")
```

```
# Show columns and data types
```

```
print(df.dtypes)
```

```
# Show sample data
```

```
print(df.head())
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
# Load dataset
```

```
df = pd.read_csv("/content/crypto_dashboard_3h_1year.csv")
```

```
df['timestamp'] = pd.to_datetime(df['timestamp'])
```

### *# 1. Summary Statistics per Coin*

```
summary = df.groupby('coin_name')[['open', 'close', 'volume', 'market_cap_usd', 'dominance_pct',  
'rolling_volatility']].describe()
```

```
print(summary)
```

### *# 2. Hourly Closing Price Trend - Separate plots per coin*

```
coins = df['coin_name'].unique()
```

```
plt.figure(figsize=(15,10))
```

```
for i, coin in enumerate(coins, 1):
```

```
    plt.subplot(3, 3, i)
```

```
    coin_df = df[df['coin_name'] == coin]
```

```
    plt.plot(coin_df['timestamp'], coin_df['close'])
```

```
    plt.title(f'{coin} Hourly Closing Price')
```

```
    plt.xlabel('Date')
```

```
    plt.ylabel('Price (USD)')
```

```
    plt.xticks(rotation=45)
```

```
plt.tight_layout()
```

```
plt.show()
```

### *# 3. Total Trading Volume per Coin*

```
volumes = df.groupby('coin_name')['volume'].sum().sort_values(ascending=False)
```

```

plt.figure(figsize=(10,6))

sns.barplot(x=volumes.index, y=volumes.values)

plt.title('Total Trading Volume by Cryptocurrency (90 days)')

plt.ylabel('Total Volume')

plt.xticks(rotation=45)

plt.show()

```

#### *# 4. Dominance Trends*

```

plt.figure(figsize=(15,10))

for i, coin in enumerate(coins, 1):

    plt.subplot(3, 3, i)

    coin_df = df[df['coin_name'] == coin]

    plt.plot(coin_df['timestamp'], coin_df['dominance_pct'])

    plt.title(f'{coin} Market Dominance % Over Time')

    plt.xlabel('Date')

    plt.ylabel('Dominance (%)')

    plt.xticks(rotation=45)

plt.tight_layout()

plt.show()

```

#### *# 5. Correlation Heatmap*

```

corr_vars = ['close', 'volume', 'market_cap_usd', 'dominance_pct', 'price_pct_change',
'rolling_volatility']

```

```

corr = df[corr_vars].corr()

print(corr)

plt.figure(figsize=(8,6))

sns.heatmap(corr, annot=True, cmap='coolwarm', fmt=".2f")

plt.title('Correlation Matrix of Key Metrics')

plt.show()

```

### *# 6. Scatter Plot Examples*

```

plt.figure(figsize=(12,5))

plt.subplot(1,2,1)

sns.scatterplot(data=df, x='dominance_pct', y='price_pct_change', hue='coin_name', alpha=0.6)

plt.title('Dominance % vs. Price % Change')

```

```

plt.subplot(1,2,2)

sns.scatterplot(data=df, x='volume', y='rolling_volatility', hue='coin_name', alpha=0.6)

plt.title('Volume vs. Rolling Volatility')

plt.tight_layout()

plt.show()

```

Model Building and Evaluation:

```

import os

import requests

import warnings

```

```

import logging

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt


from sklearn.preprocessing import StandardScaler

from sklearn.calibration import CalibratedClassifierCV

from imblearn.combine import SMOTETomek

from sklearn.metrics import classification_report, accuracy_score

from sklearn.tree import DecisionTreeClassifier, plot_tree

import lightgbm as lgb


warnings.filterwarnings("ignore")

logging.getLogger("lightgbm").setLevel(logging.ERROR)


# ——— CONFIG


---




---



MIN_BARS = 50

FNG_CACHE = "fear_greed_cache.csv"

TRAIN_SPLIT = 0.8

UP_THRESH = 0.3

DOWN_THRESH = 0.3

```



```
INITIAL_CASH = 10_000
```

```
# ——— FEAR & GREED
```

---

```
def fetch_fear_greed(days=365, cache=FNG_CACHE):
```

```
    if os.path.exists(cache):
```

```
        fg = pd.read_csv(cache, parse_dates=["date"])
```

```
    else:
```

```
        resp = requests.get("https://api.alternative.me/fng/?limit=0&format=json")
```

```
        data = resp.json().get("data", [])
```

```
        fg = pd.DataFrame(data)
```

```
        fg["date"] = pd.to_datetime(fg["timestamp"], unit="s")
```

```
        fg["fear_greed"] = fg["value"].astype(int)
```

```
        fg = fg[["date", "fear_greed"]]
```

```
        fg.to_csv(cache, index=False)
```

```
    fg["date"] = fg["date"].dt.date
```

```
    cutoff = pd.Timestamp.today().date() - pd.Timedelta(days=days)
```

```
    return fg[fg["date"]>=cutoff]
```

```
# ——— INDICATORS
```

---

```
def calc_rsi(series, period=14):
```

```

d = series.diff()

up = d.clip(lower=0).rolling(period).mean()

dn = (-d.clip(upper=0)).rolling(period).mean().replace(0,np.nan)

rs = up/dn

return 100 - 100/(1+rs)

```

```

def calc_macd(series, fast=12, slow=26, sig=9):

    fe = series.ewm(span=fast,adjust=False).mean()

    se = series.ewm(span=slow,adjust=False).mean()

    macd = fe - se

    return macd, macd.ewm(span=sig,adjust=False).mean()

```

```

def calc_bb(series, window=20, n=2):

    ma = series.rolling(window).mean()

    sd = series.rolling(window).std()

    return ma+n*sd, ma-n*sd

```

# ——— *DATA PREP*

---



---

```

def prepare_data(df):

    df = df.loc[:,~df.columns.duplicated()].copy()

    df["timestamp"] = pd.to_datetime(df["timestamp"], errors="coerce")

```

```
df["date"] = df["timestamp"].dt.date
```

```
# pct change & lags
```

```
df["pct_chg"] = df.groupby("symbol")["close"].pct_change()
```

```
for lag in (1,3,6,12):
```

```
    df[f'lag{lag}'] = df.groupby("symbol")["pct_chg"].shift(lag)
```

```
# volume vs avg & rolling vol
```

```
df["vol20"] = df.groupby("symbol")["volume"].transform(lambda x: x.rolling(20).mean())
```

```
df["vol_to_avg"] = (df["volume"]/df["vol20"]).clip(0,5)
```

```
df["roll_vol"] = df.groupby("symbol")["pct_chg"].transform(lambda x: x.rolling(14).std())
```

```
# indicators
```

```
df["rsi"] = df.groupby("symbol")["close"].transform(calc_rsi)
```

```
macd_vals = df.groupby("symbol")["close"].transform(lambda s: calc_macd(s)[0])
```

```
sig_vals = df.groupby("symbol")["close"].transform(lambda s: calc_macd(s)[1])
```

```
df["macd"], df["macd_sig"] = macd_vals, sig_vals
```

```
bb_up_vals, bb_lo_vals = (
```

```
    df.groupby("symbol")["close"].transform(lambda s: calc_bb(s)[0]),
```

```
    df.groupby("symbol")["close"].transform(lambda s: calc_bb(s)[1])
```

```
)
```

```
df["bb_up"], df["bb_lo"] = bb_up_vals, bb_lo_vals
```

```

# target

df["future_pct"] = df.groupby("symbol")["pct_chg"].shift(-1)

df["target"] = np.where(df["future_pct"] > 0.003, "Up",
                        np.where(df["future_pct"] < -0.003, "Down", "Stable"))

# fill

cols = [c for c in df.columns if c not in ["timestamp", "symbol", "target"]]

df[cols] = df.groupby("symbol")[cols].transform(lambda g: g.ffmpeg().bfill())

df.fillna(df.median(numeric_only=True), inplace=True)

# merge Fear & Greed

fg = fetch_fear_greed()

df = df.merge(fg, on="date", how="left")

df["fear_greed"] = df["fear_greed"].ffmpeg().bfill().fillna(50)

df.drop("date", axis=1, inplace=True)

return df[df.groupby("symbol")["symbol"].transform("count") >= MIN_BARS] \
        .dropna(subset=["target"])

# ——— TRAIN & CALIBRATE

———

def train_multiclass(df):

```

```

feats = ["lag1","lag3","lag6","lag12","roll_vol","vol_to_avg",
        "rsi","macd","macd_sig","bb_up","bb_lo","fear_greed"]

df = df.sort_values("timestamp")

X, y = df[feats].values, df["target"].values

scaler = StandardScaler().fit(X)

Xs = scaler.transform(X)

Xr, yr = SMOTETomek(random_state=42).fit_resample(Xs, y)

base = lgb.LGBMClassifier(
    class_weight={"Up":2,"Stable":1,"Down":4},
    learning_rate=0.03, max_depth=6, n_estimators=400,
    random_state=42, verbose=-1
)

base.fit(Xr, yr)

calib = CalibratedClassifierCV(base, cv=5, method="isotonic")

calib.fit(Xs, y)

preds = calib.predict(Xs)

print("=== Multiclass + Fear&Greed Report ===")

print(f'Accuracy: {accuracy_score(y, preds):.4f}')

```

```
print(classification_report(y, preds, zero_division=0))
```

```
return calib, scaler, feats
```

```
# ——— SIGNAL GENERATION
```

---

```
def generate_signals(df, model, scaler, feats, up_t=UP_THRESH, down_t=DOWN_THRESH):
```

```
    out = df.copy()
```

```
    X = scaler.transform(out[feats].values)
```

```
    p = model.predict_proba(X)
```

```
    cls = model.classes_
```

```
    ui = list(cls).index("Up")
```

```
    di = list(cls).index("Down")
```

```
    out["prob_up"] = p[:,ui]
```

```
    out["prob_down"] = p[:,di]
```

```
    out["signal"] = "Hold"
```

```
    out.loc[out["prob_down"]>down_t, "signal"] = "Sell"
```

```
    out.loc[(out["signal"]=="Hold") & (out["prob_up"]>up_t), "signal"] = "Buy"
```

```
return out[["timestamp", "symbol", "close", "fear_greed", "prob_up", "prob_down", "signal"]]
```

# ——— *NEXT-3H HELPER*

---

```
def next_3h_signals(df, model, scaler, feats, up_t=UP_THRESH, down_t=DOWN_THRESH):

    latest = df.groupby("symbol").tail(1).copy()

    latest["timestamp"] += pd.Timedelta(hours=3)

    for lag in (1,3,6,12):

        latest[f"lag{lag}"] = df.groupby("symbol")["pct_chg"].shift(lag-1).tail(len(latest)).values

    return generate_signals(latest, model, scaler, feats, up_t, down_t)
```

# ——— *Random Buy/Sell/Hold signals with same  
distribution*

---

```
import numpy as np
```

```
def generate_random_signals(df):

    """

    Assigns random Buy/Sell/Hold signals with same distribution.

    """

    signals = np.random.choice(["Buy", "Sell", "Hold"], size=len(df), p=[0.4, 0.4, 0.2])

    df_random = df.copy()

    df_random["signal"] = signals

    return df_random
```

# ——— *BACKTEST*

---



---

**def** backtest(sigs):

    cash, alloc, pos = INITIAL\_CASH, INITIAL\_CASH\*0.4, {s:0 **for** s **in**  
sigs["symbol"].unique()}

    buyp, buyt, trades, port = {}, {}, [], []

**for** \_, r **in** sigs.sort\_values("timestamp").iterrows():

    ts, sym, price, sig = r["timestamp"], r["symbol"], r["close"], r["signal"]

*# auto-exit SL/TP/TO*

**if** pos[sym]>0:

        bp, bt = buyp[sym], buyt[sym]

**if** price<bp\*0.985 **or** price>bp\*1.035 **or** (ts-bt).total\_seconds()/3600>=6:

            cash += pos[sym]\*price\*0.999

            trades.append((ts,sym,"AutoSell",price,pos[sym]))

            pos[sym]=0

**if** sig=="Buy" **and** cash>=alloc **and** pos[sym]\*price<INITIAL\_CASH\*0.7:

        qty = alloc/price\*0.999

        pos[sym]+=qty; buyp[sym],buyt[sym]=price,ts

        cash-=qty\*price; trades.append((ts,sym,"Buy",price,qty))



```

elif sig=="Sell" and pos[sym]>0:

    cash+=pos[sym]*price*0.999

    trades.append((ts,sym,"Sell",price,pos[sym])); pos[sym]=0

```

```

total = cash + sum(pos[k]*price for k in pos)

```

```

port.append((ts,total))

```

```

trades_df = pd.DataFrame(trades,columns=["timestamp","symbol","action","price","units"])

```

```

port_df = pd.DataFrame(port,columns=["timestamp","value"])

```

```

print(f"\nFinal portfolio: ${port_df['value'].iloc[-1]:.2f}, trades: {len(trades_df)}")

```

```

port_df.plot(x="timestamp",y="value",title="Portfolio Value",grid=True)

```

```

plt.show()

```

```

return trades_df, port_df

```

```

# ——— MAIN

```

---

```

if __name__=="__main__":

```

```

    raw = pd.read_csv("crypto_dashboard_3h_1year.csv")

```

```

    print(f"Loaded {len(raw)} rows, symbols={raw['symbol'].nunique()}")

```

```

    print(f"Date range: {raw['timestamp'].min()} → {raw['timestamp'].max()}")

```

```

    data = prepare_data(raw)

```

```

model, scaler, feats = train_multiclass(data)

split = int(len(data)*TRAIN_SPLIT)

test = data.iloc[split:]

print(f"\n--- Test slice @ up={UP_THRESH}, down={DOWN_THRESH} ---")

sigs_test = generate_signals(test, model, scaler, feats)

print("Signal counts:\n", sigs_test["signal"].value_counts())

backtest(sigs_test)

print(f"\n--- Full history @ up={UP_THRESH}, down={DOWN_THRESH} ---")

sigs_all = generate_signals(data, model, scaler, feats)

# Merge model-generated signals back with full feature data

validation_df = data[["timestamp", "symbol", "future_pct", "fear_greed",
"rolling_volatility"]].copy()

validation_df = validation_df.merge(sigs_all[["timestamp", "symbol", "signal", "prob_up",
"prob_down"]],

                                on=["timestamp", "symbol"],

                                how="inner")

# Drop any rows with missing future_pct (e.g., last row per symbol)

validation_df = validation_df.dropna(subset=["future_pct", "signal"])

validation_df.to_csv("model_validation_dataset.csv", index=False)

```

```
print("Signal counts:\n", sigs_all["signal"].value_counts())
```

```
backtest(sigs_all)
```

```
# ——— Next 3-Hour prescription
```

---

```
next3 = next_3h_signals(data, model, scaler, feats)
```

```
print("\n--- Next 3-Hour Signals ---")
```

```
print(next3.to_string(index=False))
```

```
# ——— Visualize prescriptive logic tree
```

---

```
df2 = sigs_all[["prob_down", "prob_up", "signal"]].dropna()
```

```
X2 = df2[["prob_down", "prob_up"]]
```

```
y2 = df2["signal"]
```

```
tree = DecisionTreeClassifier(max_depth=2)
```

```
tree.fit(X2, y2)
```

```
plt.figure(figsize=(6,4))
```

```
plot_tree(tree,
```

```
    feature_names=["prob_down", "prob_up"],
```

```
    class_names=tree.classes_,
```

```
    rounded=True,
```

```

        fontsize=10)

plt.title("Prescriptive Logic Tree")

plt.show()


def evaluate_strategy_vs_random(sigs_all, model_value, n_runs=30):

    final_vals = []

    print("\n🚀 Running random strategy simulations...")

    for i in range(n_runs):

        rand_sigs = generate_random_signals(sigs_all)

        _, port_df = backtest(rand_sigs)

        final_vals.append(port_df["value"].iloc[-1])

    rand_avg = np.mean(final_vals)

    t_stat, p_val = ttest_1samp(final_vals, model_value)

    print("\n📊 Portfolio Value Comparison")

    print(f"Model-guided strategy: ${model_value:.2f}")

    print(f"Random strategy mean: ${rand_avg:.2f}")

    print(f"T-statistic: {t_stat:.4f}")

    print(f"P-value: {p_val:.4f} (Significant if < 0.05)")

```

```
return final_vals
```

```
import os
```

```
import requests
```

```
import warnings
```

```
import logging
```

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
from sklearn.preprocessing import StandardScaler
```

```
from sklearn.calibration import CalibratedClassifierCV
```

```
from imblearn.combine import SMOTETomek
```

```
from sklearn.metrics import classification_report, accuracy_score
```

```
from sklearn.tree import DecisionTreeClassifier, plot_tree
```

```
import lightgbm as lgb
```

```
warnings.filterwarnings("ignore")
```

```
logging.getLogger("lightgbm").setLevel(logging.ERROR)
```

# — *CONFIG*

---



---

MIN\_BARS = 50

FNG\_CACHE = "fear\_greed\_cache.csv"

TRAIN\_SPLIT = 0.8

UP\_THRESH = 0.3

DOWN\_THRESH = 0.3

INITIAL\_CASH = 10\_000

# — *FEAR & GREED*

---



---

**def** fetch\_fear\_greed(days=365, cache=FNG\_CACHE):

**if** os.path.exists(cache):

        fg = pd.read\_csv(cache, parse\_dates=["date"])

**else:**

        resp = requests.get("https://api.alternative.me/fng/?limit=0&format=json")

        data = resp.json().get("data", [])

        fg = pd.DataFrame(data)

        fg["date"] = pd.to\_datetime(fg["timestamp"], unit="s")

        fg["fear\_greed"] = fg["value"].astype(int)

        fg = fg[["date", "fear\_greed"]]

        fg.to\_csv(cache, index=False)

```

fg["date"] = fg["date"].dt.date

cutoff = pd.Timestamp.today().date() - pd.Timedelta(days=days)

return fg[fg["date"]>=cutoff]

```

# ——— *INDICATORS*

---

```

def calc_rsi(series, period=14):

    d = series.diff()

    up = d.clip(lower=0).rolling(period).mean()

    dn = (-d.clip(upper=0)).rolling(period).mean().replace(0, np.nan)

    rs = up/dn

    return 100 - 100/(1+rs)

```

```

def calc_macd(series, fast=12, slow=26, sig=9):

    fe = series.ewm(span=fast,adjust=False).mean()

    se = series.ewm(span=slow,adjust=False).mean()

    macd = fe - se

    return macd, macd.ewm(span=sig,adjust=False).mean()

```

```

def calc_bb(series, window=20, n=2):

    ma = series.rolling(window).mean()

    sd = series.rolling(window).std()

```

```
return ma+n*sd, ma-n*sd
```

```
# ——— DATA PREP
```

---

```
def prepare_data(df):
```

```
    df = df.loc[:,~df.columns.duplicated()].copy()
```

```
    df["timestamp"] = pd.to_datetime(df["timestamp"], errors="coerce")
```

```
    df["date"]     = df["timestamp"].dt.date
```

```
    # pct change & lags
```

```
    df["pct_chg"] = df.groupby("symbol")["close"].pct_change()
```

```
    for lag in (1,3,6,12):
```

```
        df[f"lag{lag}"] = df.groupby("symbol")["pct_chg"].shift(lag)
```

```
    # volume vs avg & rolling vol
```

```
    df["vol20"]     = df.groupby("symbol")["volume"].transform(lambda x: x.rolling(20).mean())
```

```
    df["vol_to_avg"] = (df["volume"]/df["vol20"]).clip(0,5)
```

```
    df["roll_vol"]  = df.groupby("symbol")["pct_chg"].transform(lambda x: x.rolling(14).std())
```

```
    # indicators
```

```
    df["rsi"]       = df.groupby("symbol")["close"].transform(calc_rsi)
```

```
    macd_vals       = df.groupby("symbol")["close"].transform(lambda s: calc_macd(s)[0])
```



```

sig_vals    = df.groupby("symbol")["close"].transform(lambda s: calc_macd(s)[1])

df["macd"], df["macd_sig"] = macd_vals, sig_vals

bb_up_vals, bb_lo_vals = (
    df.groupby("symbol")["close"].transform(lambda s: calc_bb(s)[0]),
    df.groupby("symbol")["close"].transform(lambda s: calc_bb(s)[1])
)

df["bb_up"], df["bb_lo"] = bb_up_vals, bb_lo_vals


# target

df["future_pct"] = df.groupby("symbol")["pct_chg"].shift(-1)

df["target"]    = np.where(df["future_pct"]> 0.003, "Up",
                           np.where(df["future_pct"]< -0.003, "Down", "Stable"))


# fill

cols = [c for c in df.columns if c not in ["timestamp", "symbol", "target", "date"]]

df[cols] = df.groupby("symbol")[cols].transform(lambda g: g.ffmpeg().bfill())

df.fillna(df.median(numeric_only=True), inplace=True)


# merge Fear & Greed

fg = fetch_fear_greed()

df = df.merge(fg, on="date", how="left")

df["fear_greed"] = df["fear_greed"].ffmpeg().bfill().fillna(50)

df.drop("date", axis=1, inplace=True)

```

```

return df[df.groupby("symbol")["symbol"].transform("count")>=MIN_BARS]\
        .dropna(subset=["target"])

```

```

# ——— TRAIN & CALIBRATE

```

---

```

def train_multiclass(df):

```

```

    feats = ["lag1", "lag3", "lag6", "lag12", "roll_vol", "vol_to_avg",
            "rsi", "macd", "macd_sig", "bb_up", "bb_lo", "fear_greed"]

```

```

    df = df.sort_values("timestamp")

```

```

    X, y = df[feats].values, df["target"].values

```

```

    scaler = StandardScaler().fit(X)

```

```

    Xs = scaler.transform(X)

```

```

    Xr, yr = SMOTETomek(random_state=42).fit_resample(Xs, y)

```

```

    base = lgb.LGBMClassifier(

```

```

        class_weight={"Up":2, "Stable":1, "Down":4},

```

```

        learning_rate=0.03, max_depth=6, n_estimators=400,

```

```

        random_state=42, verbose=-1

```

```

    )

```

```

base.fit(Xr, yr)

import pandas as pd

import matplotlib.pyplot as plt

# assume `base` is your trained LGBMClassifier (before calibration)

importances = pd.Series(base.feature_importances_, index=feats)

importances.sort_values().plot.barh(figsize=(6,4), title="LightGBM Feature Importances")

plt.xlabel("Importance Score")

plt.show()


calib = CalibratedClassifierCV(base, cv=5, method="isotonic")

calib.fit(Xs, y)


preds = calib.predict(Xs)

print("=== Multiclass + Fear&Greed Report ===")

print(f'Accuracy: {accuracy_score(y, preds):.4f}')

print(classification_report(y, preds, zero_division=0))


return calib, scaler, feats

```

# ——— *SIGNAL GENERATION*

---

**def** generate\_signals(df, model, scaler, feats, up\_t=UP\_THRESH, down\_t=DOWN\_THRESH):

    out = df.copy()

    X = scaler.transform(out[feats].values)

    p = model.predict\_proba(X)

    cls = model.classes\_

    ui = list(cls).index("Up")

    di = list(cls).index("Down")

    out["prob\_up"] = p[:,ui]

    out["prob\_down"] = p[:,di]

    out["signal"] = "Hold"

    out.loc[out["prob\_down"]>down\_t, "signal"] = "Sell"

    out.loc[(out["signal"]=="Hold") & (out["prob\_up"]>up\_t), "signal"] = "Buy"

**return** out[["timestamp","symbol","close","fear\_greed","prob\_up","prob\_down","signal"]]

# ——— *NEXT-3H HELPER*

---

**def** next\_3h\_signals(df, model, scaler, feats, up\_t=UP\_THRESH, down\_t=DOWN\_THRESH):

    latest = df.groupby("symbol").tail(1).copy()

    latest["timestamp"] += pd.Timedelta(hours=3)

```
for lag in (1,3,6,12):
```

```
    latest[f"lag{lag}"] = df.groupby("symbol")["pct_chg"].shift(lag-1).tail(len(latest)).values
```

```
return generate_signals(latest, model, scaler, feats, up_t, down_t)
```

```
# ——— BACKTEST
```

---

```
def backtest(sigs):
```

```
    cash, alloc, pos = INITIAL_CASH, INITIAL_CASH*0.4, {s:0 for s in
sigs["symbol"].unique()}
```

```
    buyp, buyt, trades, port = {}, {}, [], []
```

```
for _, r in sigs.sort_values("timestamp").iterrows():
```

```
    ts, sym, price, sig = r["timestamp"], r["symbol"], r["close"], r["signal"]
```

```
    # auto-exit SL/TP/TO
```

```
    if pos[sym]>0:
```

```
        bp, bt = buyp[sym], buyt[sym]
```

```
        if price<bp*0.985 or price>bp*1.035 or (ts-bt).total_seconds()/3600>=6:
```

```
            cash += pos[sym]*price*0.999
```

```
            trades.append((ts,sym,"AutoSell",price,pos[sym]))
```

```
            pos[sym]=0
```

```
    if sig=="Buy" and cash>=alloc and pos[sym]*price<INITIAL_CASH*0.7:
```

```

    qty = alloc/price*0.999

    pos[sym]+=qty; buyp[sym],buyt[sym]=price,ts

    cash-=qty*price; trades.append((ts,sym,"Buy",price,qty))

elif sig=="Sell" and pos[sym]>0:

    cash+=pos[sym]*price*0.999

    trades.append((ts,sym,"Sell",price,pos[sym])); pos[sym]=0


total = cash + sum(pos[k]*price for k in pos)

port.append((ts,total))


trades_df = pd.DataFrame(trades,columns=["timestamp","symbol","action","price","units"])
port_df = pd.DataFrame(port,columns=["timestamp","value"])

print(f"\nFinal portfolio: ${port_df['value'].iloc[-1]:.2f}, trades: {len(trades_df)}")

port_df.plot(x="timestamp",y="value",title="Portfolio Value",grid=True)

plt.show()

return trades_df, port_df


# ——— MAIN


---




---



if __name__=="__main__":

    raw = pd.read_csv("crypto_dashboard_3h_1year.csv")

    print(f"Loaded {len(raw)} rows, symbols={raw['symbol'].nunique()}")

```

```

print(f'Date range: {raw['timestamp'].min()} → {raw['timestamp'].max()}')

data = prepare_data(raw)

model, scaler, feats = train_multiclass(data)

# test-slice

split = int(len(data)*TRAIN_SPLIT)

sigs_test = generate_signals(data.iloc[split:], model, scaler, feats)

print(f'\n--- Test slice @ up={UP_THRESH}, down={DOWN_THRESH} ---')

print(sigs_test["signal"].value_counts())

backtest(sigs_test)

# full-history

sigs_all = generate_signals(data, model, scaler, feats)

print(f'\n--- Full history @ up={UP_THRESH}, down={DOWN_THRESH} ---')

backtest(sigs_all)

# create validation dataset

validation_df = data[[

    "timestamp", "symbol", "future_pct", "fear_greed", "roll_vol",

    "lag1", "lag3", "lag6", "lag12"

]].copy()

validation_df = validation_df.merge(

```

```

sigs_all[["timestamp","symbol","signal","prob_up","prob_down"]],
on=["timestamp","symbol"], how="inner"
).dropna(subset=["future_pct","signal"])
validation_df.to_csv("model_validation_dataset.csv", index=False)

# next-3h
next3 = next_3h_signals(data, model, scaler, feats)
print("\n--- Next 3-Hour Signals ---")
print(next3.to_string(index=False))

# logic tree
df2 = sigs_all[["prob_down","prob_up","signal"]].dropna()
tree = DecisionTreeClassifier(max_depth=2).fit(df2[["prob_down","prob_up"]], df2["signal"])
plt.figure(figsize=(6,4))
plot_tree(tree,
          feature_names=["prob_down","prob_up"],
          class_names=tree.classes_,
          rounded=True, fontsize=10)
plt.title("Prescriptive Logic Tree")
plt.show()

```

Hypothesis1 Evaluation:



```
##H1--
```

```
# ——— EVALUATE MODEL PORTFOLIO
```

---

```
# Generate full-history signals and backtest to get the model's final portfolio value
```

```
sigs_all = generate_signals(data, model, scaler, feats)
```

```
_, model_port = backtest(sigs_all)
```

```
model_value = model_port["value"].iloc[-1]
```

```
# ——— VALIDATION: RANDOM BASELINE
```

---

```
# Run 10 random-signal simulations and collect their final values
```

```
random_final_values = evaluate_strategy_vs_random(
```

```
    sigs_all,
```

```
    model_value,
```

```
    n_runs=10
```

```
)
```

```
# ——— VISUALIZE COMPARISON
```

---

```
import matplotlib.pyplot as plt
```

```
plt.figure(figsize=(8, 5))
```

```
plt.hist(random_final_values, bins=10, alpha=0.7, label="Random Strategy")
```

```

plt.axvline(model_value, color="red", linestyle="--", label="Model Strategy")

plt.title("Random vs. Model Final Portfolio Values")

plt.xlabel("Final Portfolio Value")

plt.ylabel("Frequency")

plt.legend()

plt.grid(True)

plt.show()

```

Hypothesis2 Validation:

```

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

from scipy.stats import kruskal

# ——— LOAD VALIDATION DATA


---



df = pd.read_csv("model_validation_dataset.csv")

# ——— H2: Signal Type vs Future Return


---



# Extract returns per signal class

buy_returns = df[df["signal"] == "Buy"]["future_pct"]

```

```

sell_returns = df[df["signal"] == "Sell"]["future_pct"]

hold_returns = df[df["signal"] == "Hold"]["future_pct"]

# Perform Kruskal–Wallis test (non-parametric ANOVA)
h_stat, p_val_h2 = kruskal(buy_returns, sell_returns, hold_returns)

```

*# Print results*

```

print("\n📊 H2: Signal Type vs Future Return")

print(f'Buy mean: {buy_returns.mean():.4f}')
print(f'Sell mean: {sell_returns.mean():.4f}')
print(f'Hold mean: {hold_returns.mean():.4f}')

print(f'H-statistic: {h_stat:.4f}, p-value: {p_val_h2:.4f}')

```

*# ——— VISUALIZATION*

---

```

plt.figure(figsize=(8, 5))

sns.boxplot(x="signal", y="future_pct", data=df, palette="Set3")

plt.title("H2: Future Returns by Signal Type")

plt.ylabel("Future % Change")

plt.grid(True)

plt.show()

```