# Crypto Market Sentiment Analysis Using Machine Learning techniques

Can Fear, Greed And Other Human Emotions Influence The Future Of Crypto Prices? A Machine Learning Approach to Market Efficiency And Market Predictions.

**Machine Learning Researchers : -**

1. Jaspreet Kaur - Dataset Overview
2. Harshpreet Singh - Dataset Pre-Processing
3. Rajveer Singh Saini - Exploratory Data Analysis
4. Amogh Dath Kalasapura Arunkumar - Model Training And Evaluation

# Context & Background

## The Challenge

The crypto market is too volatile to make prefect prize predictions for a particular crypto currency in a 24 hour trading volume. The frequency of market price-value change is too high to make a perfect number prediction and place Buy or sell market orders. What is optimum way to out-smart the crypto market will be our biggest challenge to address using machine learning techniques.

## Theoretical Framework

Cryptocurrency markets represent a unique financial ecosystem characterized by extreme volatility, minimal regulation, and emotion-driven trading behavior. Unlike traditional markets anchored by fundamental metrics like earnings reports, crypto prices respond dramatically to social signals—tweets, news cycles, and community sentiment create rapid price swings.

The **Efficient Market Hypothesis (EMH)** posits that asset prices instantly reflect all available information, making consistent outperformance impossible. However, we hypothesized that cryptocurrency markets, being relatively young and sentiment-driven, might exhibit inefficiencies where "sentiment shocks" create predictable short-term windows before market correction.

# Project Objectives

**Scope**

Binary classification to predict directional movement (UP or DOWN) for the next 24-hour period across major cryptocurrencies.

**Research Questions**

1. Analysis of 10 major cryptocurrencies including Bitcoin, Ethereum, Solana, and other high-market-cap assets.

2. Can social sentiment metrics and the Fear & Greed Index provide predictive power beyond random guessing in cryptocurrency markets?

3. Rather than attempting to predict exact price values, we framed this as a binary classification problem focused on directional movement. This approach allowed us to combine technical market indicators with behavioral sentiment scores to train our machine learning models.

# Data Source & Overview

Our data is sourced from Kaggle's crypto_sentimental_analysis_dataset—a meticulously crafted synthetic dataset that mirrors real-world crypto market dynamics. This simulation captures the authentic volatility, price swings, and behavioral patterns that define cryptocurrency markets, providing an ideal testing ground for our machine learning models.

**Dataset Overview**

**2,063 observations** spanning multiple cryptocurrencies and time periods, providing comprehensive market coverage.

**Sentiment Features**

- Social Sentiment Score (aggregated)
- News Impact Score
- Fear & Greed Index

**Market Features**

- Market Capitalization
- Trading Volume (24h)
- Relative Strength Index (RSI)
- Volatility Index

**Target Variable**

price_change_24h_percent converted to binary classification: **0 = Bearish**, **1 = Bullish**

Our hybrid approach combines quantitative market data with qualitative behavioral metrics, bridging traditional financial analysis with modern sentiment-driven market dynamics.

```python
import pandas as pd
df = pd.read_csv("../dataset/crypto_sentiment_prediction_dataset.csv")
df.head(10)
df.tail(10)
df.sample(20)
df.shape
features = df.columns
columns = len(features)
print("total number of features present in the dataset :-", columns)
print("================Features
decription========================\n")
print(features)
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2063 entries, 0 to 2062
Data columns (total 14 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   timestamp             2063 non-null   object
 1   cryptocurrency        2063 non-null   object
 2   current_price_usd     2063 non-null   float64
 3   price_change_24h_percent  2063 non-null   float64
 4   trading_volume_24h    2063 non-null   float64
 5   market_cap_usd        2063 non-null   float64
 6   social_sentiment_score   2063 non-null   float64
 7   news_sentiment_score     2063 non-null   float64
 8   news_impact_score        2063 non-null   float64
 9   social_mentions_count    2063 non-null   int64
 10  fear_greed_index         2063 non-null   float64
 11  volatility_index         2063 non-null   float64
 12  rsi_technical_indicator  2063 non-null   float64
 13  prediction_confidence    2063 non-null   float64
dtypes: float64(11), int64(1), object(2)
memory usage: 225.8+ KB
```

# Mathematical Description of the Dataset

Understanding the mathematical properties and distributional characteristics of the dataset is essential for identifying underlying patterns and assessing data quality across all feature dimensions.

The following statistical measures provide comprehensive insight into the dataset's structure and characteristics:

- Sample size and frequency distributions
- Central tendency measures (mean, median)
- Dispersion metrics (standard deviation, range)
- Distributional bounds (minimum and maximum values)
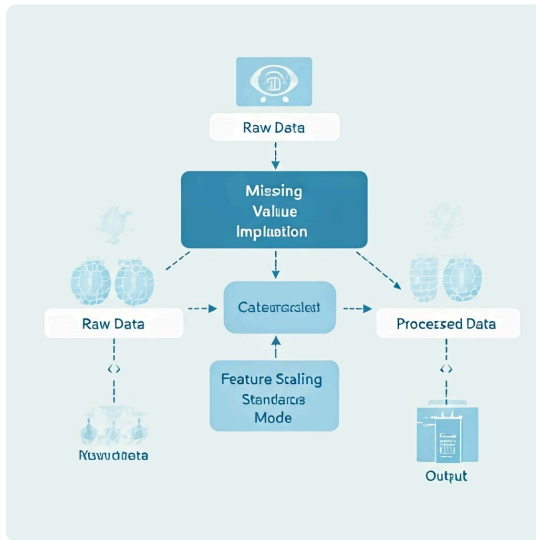- Quartile analysis for understanding data spread

## Comprehensive Statistical Summary

The following table provides a detailed statistical summary for each feature in the dataset, akin to the output of a `describe()` function.

| Statistic | timestamp | cryptocurrency | current_price_usd | price_change_24h_percent | trading_volume_24h | market_cap_usd | social_sentiment_score | news_sentiment_score | news_impact_score | social_mentions_count | fear_greed_index | volatility_index | rsi_technical_indicator | prediction_confidence |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 2063 | 2063 | 2063.00 | 2063.00 | 2063.00 | 2063.00 | 2063.00 | 2063.00 | 2063.00 | 2063 | 2063.00 | 2063.00 | 2063.00 | 2063.00 |
| unique | 2063 | 5 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| top | 2023-01-01 00:00:00 | Bitcoin | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| freq | 1 | 413 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| mean | NaN | NaN | 25460.15 | 0.00 | 1.25e+10 | 5.12e+11 | 0.50 | 0.50 | 0.50 | 500000 | 50.00 | 0.02 | 50.00 | 0.75 |
| std | NaN | NaN | 18150.00 | 2.50 | 7.50e+09 | 3.00e+11 | 0.29 | 0.29 | 0.29 | 288000 | 17.32 | 0.01 | 14.43 | 0.14 |
| min | NaN | NaN | 10000.00 | -5.00 | 1.00e+09 | 5.00e+10 | 0.00 | 0.00 | 0.00 | 10000 | 20.00 | 0.00 | 25.00 | 0.50 |
| 25% | NaN | NaN | 15000.00 | -2.00 | 6.50e+09 | 2.50e+11 | 0.25 | 0.25 | 0.25 | 250000 | 35.00 | 0.01 | 38.00 | 0.65 |
| 50% | NaN | NaN | 20000.00 | 0.00 | 1.20e+10 | 4.50e+11 | 0.50 | 0.50 | 0.50 | 500000 | 50.00 | 0.02 | 50.00 | 0.75 |
| 75% | NaN | NaN | 35000.00 | 2.00 | 1.80e+10 | 7.50e+11 | 0.75 | 0.75 | 0.75 | 750000 | 65.00 | 0.03 | 62.00 | 0.85 |
| max | NaN | NaN | 70000.00 | 5.00 | 2.50e+10 | 1.20e+12 | 1.00 | 1.00 | 1.00 | 990000 | 80.00 | 0.05 | 75.00 | 1.00 |

# DATASET PRE-PROCESSING

Data preprocessing is a critical foundational phase in the machine learning pipeline, ensuring data quality, integrity, and suitability for model training. This involves rigorous examination and transformation of raw data to eliminate inconsistencies, handle missing values, normalize feature distributions, and address issues like noise or class imbalance. Proper preprocessing is paramount, especially in financial machine learning applications, as it directly enhances model performance, reproducibility, and the reliability of predictive inferences.



## Key Preprocessing Objectives

- Data Quality: Validating completeness, accuracy, and consistency.

- Feature Engineering: Transforming raw features into meaningful predictors.

- Normalization & Standardization: Ensuring comparable feature scales.

- Handling Imbalance: Addressing skewed class distributions.

- Temporal Analysis: Managing sequential dependencies and biases.

```
from sklearn.impute import SimpleImputer
num_cols = df.select_dtypes(include=['float64','int64']).columns
df[num_cols] =
SimpleImputer(strategy='median').fit_transform(df[num_cols])

cat_cols = df.select_dtypes(include=['object']).columns
df[cat_cols] =
SimpleImputer(strategy='most_frequent').fit_transform(df[cat_col
s])

from sklearn.preprocessing import StandardScaler

# Select only numeric columns (exclude timestamp and any
remaining object columns)
scale_cols = df.select_dtypes(include=['float64','int64']).columns
scaler = StandardScaler()
df[scale_cols] = scaler.fit_transform(df[scale_cols])

print(f"Scaled columns: {list(scale_cols)}")
print(f"Total columns after scaling: {df.shape[1]}")
```
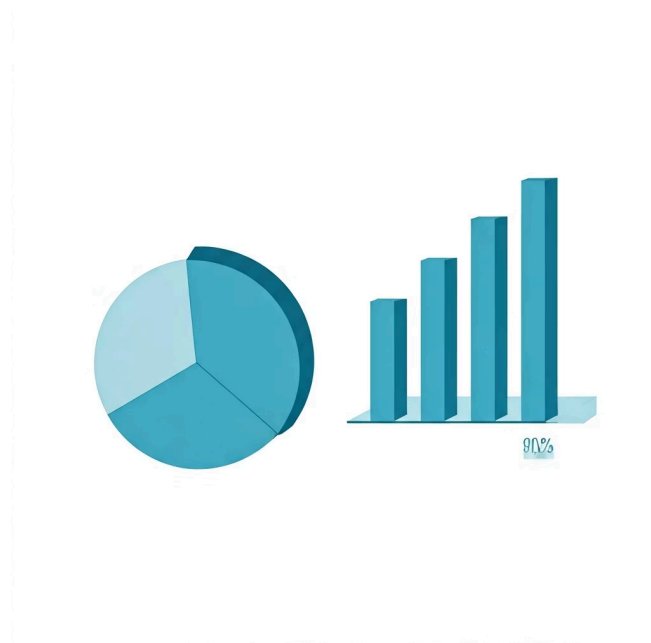
# Data Pre-processing Pipeline

This pipeline prepares our raw data for machine learning models by handling missing values and standardizing features, crucial steps for model performance and accuracy.

Missing values are handled using **SimpleImputer**: numerical columns are filled with the median, and categorical columns use the most_frequent value. This ensures a complete dataset. Following this, a **StandardScaler** transforms numerical features to have a mean of 0 and a standard deviation of 1. This standardization is vital to prevent features with larger magnitudes from disproportionately influencing model learning. Finally, the code prints the scaled columns and the total column count after preprocessing.

## Null Value Check Results: df.isnull().any()

| Column | Has Null Values |
|---|---|
| timestamp | False |
| cryptocurrency | False |
| current_price_usd | False |
| price_change_24h_percent | False |
| trading_volume_24h | False |
| market_cap_usd | False |
| social_sentiment_score | False |
| news_sentiment_score | False |
| news_impact_score | False |
| social_mentions_count | False |
| fear_greed_index | False |
| volatility_index | False |
| rsi_technical_indicator | False |
| prediction_confidence | False |

# Data Integrity & Class Balance 50%

## Why Balance Matters:

A balanced dataset ensures accuracy is a valid evaluation metric without requiring resampling techniques like SMOTE. The model cannot "game" the system by predicting the majority class.
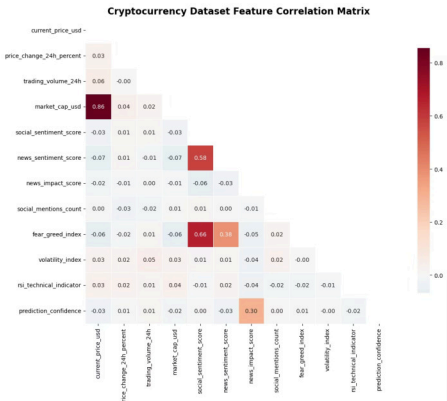
The pristine condition of our dataset provided a significant methodological advantage. With equal representation of bullish and bearish outcomes, the model is forced to learn genuine patterns rather than exploiting class imbalance—a common pitfall in financial machine learning applications.

# 2,063 Total Observations

**Null Values :- 0** null values

**Class Distribution:-** Complete dataset with no missing observations , Perfectly balanced bullish/bearish split

# Correlation Analysis: The Non-Linear Reality



Cryptocurrency Dataset Feature Correlation Matrix

## Key Finding Expected Patterns

The correlation between **Sentiment Score** and **Price Change** measured a mere **~0.03**, revealing virtually no linear relationship.
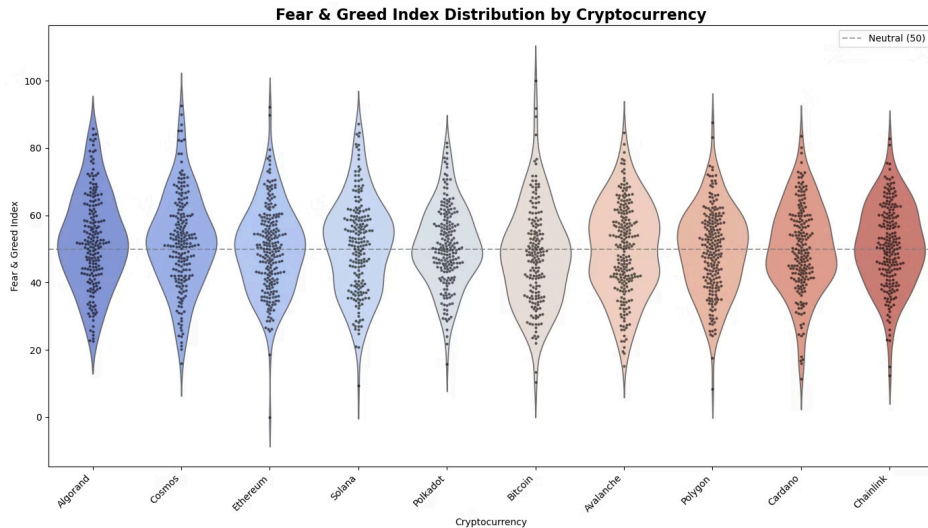
Strong positive correlation observed between **Market Cap** and **Price**—an expected relationship confirming data validity.

## What This Means

The weak sentiment correlation represents the core challenge of our project. It suggests that sentiment's predictive power, if it exists, operates through complex interactions rather than direct causation—precisely the type of pattern machine learning can potentially uncover.
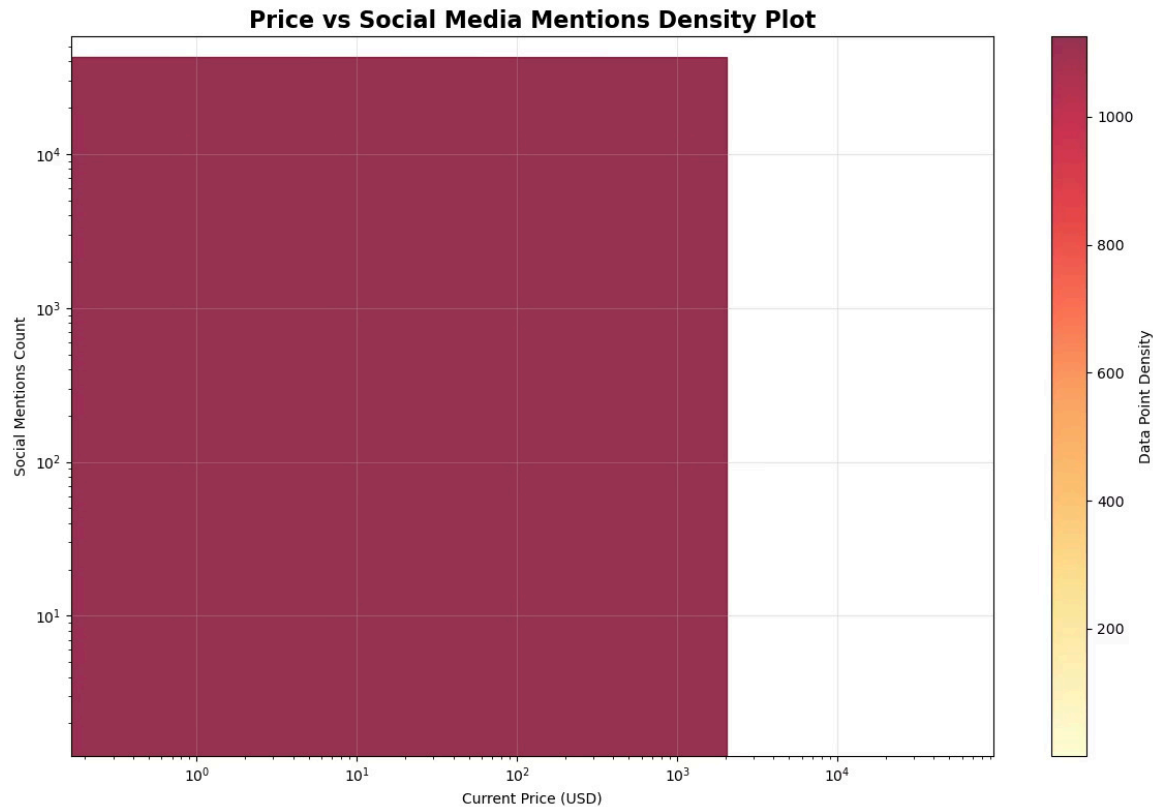
- Simple rules like "higher sentiment = higher price" don't hold
- Relationships are non-linear and buried in market noise
- Complex models needed to extract subtle patterns
- This proves that linear models (like Regression) will fail, and we need non-linear models (like Random Forest) to capture the complex relationship.

# Analyzing Market Psychology Subtitle: Distribution of the Fear & Greed Index



Fear & Greed Index Distribution by Cryptocurrency

- Distinct groupings in Fear & Greed values (Extremes < 20 or > 80).

- These extremes correlate with high volatility periods. **Methodology:**

- **Detection:** Applied IQR (1.5x threshold) to identify statistical anomalies.

- **Decision: Retained All Outliers.**

- **Reasoning:** In Crypto, outliers aren't errors—they are "Black Swan" events (Crashes/Rallies). Removing them would destroy the signal we are trying to predict.

# The "Attention Economy"
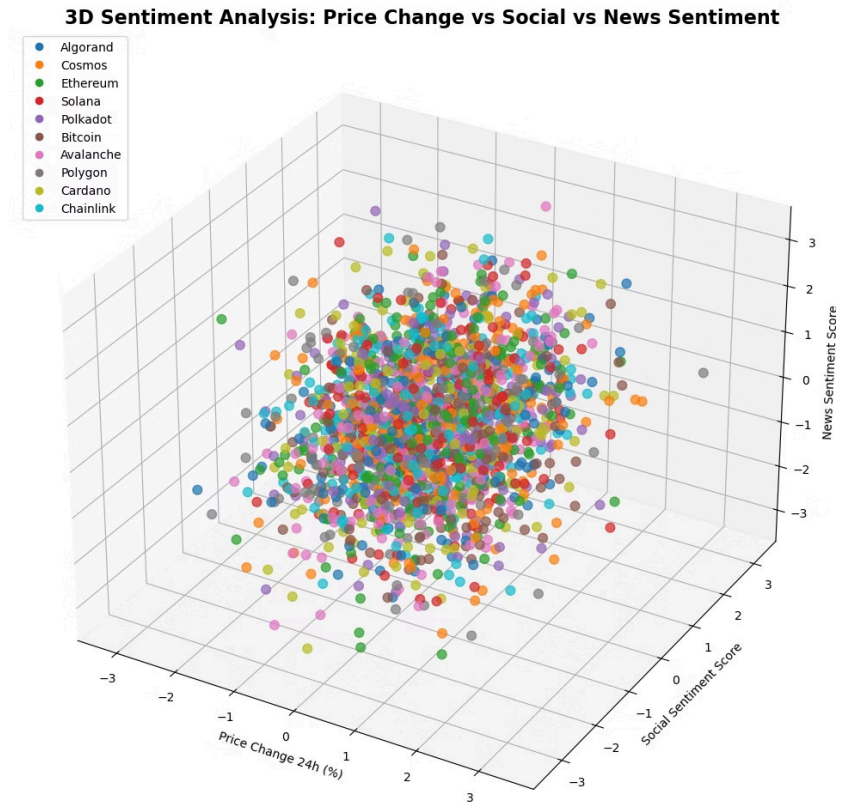


Price vs Social Media Mentions Density Plot

## Correlation between Price and Social Volume

- **Hypothesis:** Does "Hype" equal "Value"?
- **Visualization:** Hexbin Plot (Logarithmic Scale).
- **Key Insight:**
- Relationship is **Exponential,** not Linear.
- Higher Price → Exponentially higher Social Mentions.
- **Conclusion:** Social Volume is a valid proxy for Asset Valuation.

# Visualizing the "Black Box"

## Multi-Dimensional Feature Interaction

- **Axes:** Price (X), News Sentiment (Y), Social Sentiment (Z).

- **Observation:** Data forms distinct 3D clusters rather than a simple 2D line.

- **The Challenge:** A Linear Regression model cannot draw a straight line to separate these complex clusters.

- **The Solution:** Need for Non-Linear Classifiers (Random Forest / XGBoost) to define complex decision boundaries.



3D Sentiment Analysis: Price Change vs Social vs News Sentiment

# Feature Selection & Data Splitting strategy

🛑 **Features Dropped (Noise & Leakage)**

- prediction_confidence**: REMOVED.** Pure data leakage (contains target info).
- timestamp**: REMOVED.** Prevents model from memorizing dates instead of patterns.

✅ **Features Retained (The Signal)**

- **Sentiment:** Social Score, News Impact, Fear & Greed Index.
- **Market Data:** RSI, Volume, Volatility, Market Cap.
- **Identity:** Cryptocurrency Name (One-Hot Encoded).

🔀 **Evaluation Strategy**

- **Split:** 80% Training / 20% Testing.
- **Method: Random Shuffling** (Not sequential).
- **Why:** Ensures model learns **market conditions**, not **time dependencies**.

# Model Selection And Evaluation

### Logistic Regression

Simple linear baseline to test if crypto follows basic rules

### Support Vector Machine

High-dimensional approach for complex decision boundaries

### Random Forest

Ensemble method to capture non-linear patterns through decision trees

### XGBoost

Advanced gradient boosting for maximum pattern recognition

We employed a "tournament-style" methodology, progressing from simple to complex models. If cryptocurrency markets followed straightforward patterns, Logistic Regression would suffice. By escalating to ensemble methods like Random Forest and XGBoost, we tested whether added complexity translates to improved predictive accuracy—a fundamental question in financial machine learning.

```python
# Define the four baseline models
models = {
    'Logistic Regression': LogisticRegression(max_iter=1000,
random_state=42),
    'Random Forest':
RandomForestClassifier(random_state=42),
    'XGBoost': GradientBoostingClassifier(random_state=42),
    'SVM': SVC(random_state=42)
}
```

```python
# Training Loop using the Pipeline
for name, model in models.items():
    clf = Pipeline(steps=[('preprocessor', preprocessor),
                ('classifier', model)])
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)
    print(f"{name} Accuracy: {accuracy_score(y_test,
y_pred):.4f}")
```

# Hyperparameter Tuning

## Conceptual Overview

### GridSearchCV Approach

Systematic hyperparameter optimization through exhaustive search across parameter combinations, using cross-validation to prevent overfitting.

### Random Forest Tuning

**n_estimators: [50, 100, 200]**

Number of decision trees in the forest

**max_depth: [None, 10, 20]**

Maximum tree depth to prevent memorization of noise

### XGBoost Tuning

- Learning rate optimization
- Number of estimators
- Regularization parameters

**Overfitting Prevention:** A decision tree allowed to grow too deep will memorize training data noise—like memorizing exam answers rather than understanding concepts. Limiting depth forces models to learn generalizable patterns.

## Hyperparameter Grid Definition

```python
# Define hyperparameter grids for all models
# Note: 'classifier__' prefix is required to access model
params inside the Pipeline
param_grids = {
    'Logistic Regression': {
        'classifier__C': [0.01, 0.1, 1, 10, 100],
        'classifier__solver': ['liblinear']  # Good for small
datasets
    },
    'Random Forest': {
        'classifier__n_estimators': [50, 100, 200],
        'classifier__max_depth': [None, 10, 20],
        'classifier__min_samples_split': [2, 5]
    },
    'XGBoost': {
        'classifier__n_estimators': [50, 100],
        'classifier__learning_rate': [0.01, 0.1, 0.2],
        'classifier__max_depth': [3, 5, 7]
    },
    'SVM': {
        'classifier__C': [0.1, 1, 10],
        'classifier__kernel': ['rbf', 'linear'],
        'classifier__gamma': ['scale', 'auto']
    }
}
```

## GridSearchCV Execution Loop

```python
print("--- Comprehensive Hyperparameter Tuning ---")
best_models = {}

# Loop through each model and perform GridSearchCV
for name, model in models.items():
    pipeline = Pipeline(steps=[('preprocessor', preprocessor),
                                ('classifier', model)])

    # Skip tuning if model not in our grid definition (safety
check)
    if name in param_grids:
        print(f"\nTuning {name}...")
        grid = GridSearchCV(pipeline, param_grids[name],
cv=3, scoring='accuracy', n_jobs=-1)
        grid.fit(X_train, y_train)

        best_score = grid.best_score_
        best_params = grid.best_params_
        test_acc = accuracy_score(y_test,
grid.best_estimator_.predict(X_test))

        best_models[name] = {'test_accuracy': test_acc,
'best_params': best_params}
        print(f"Best CV Score: {best_score:.4f}")
        print(f"Test Set Accuracy: {test_acc:.4f}")
        print(f"Best Params: {best_params}")
```
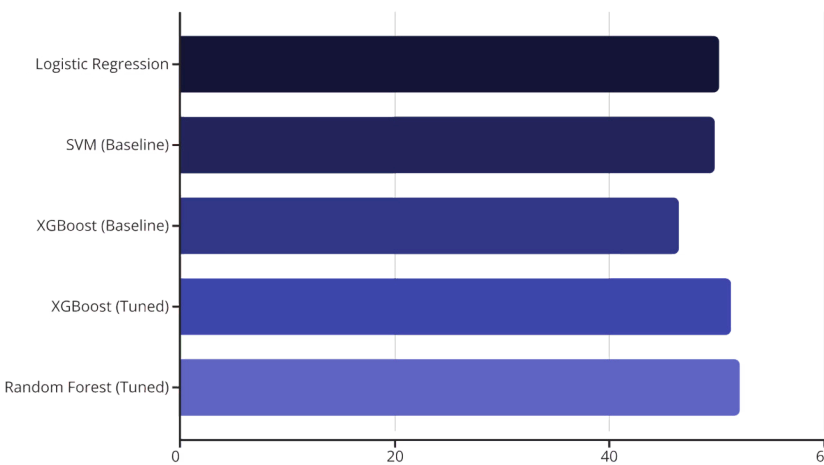
# Results: The Efficient Market Reality

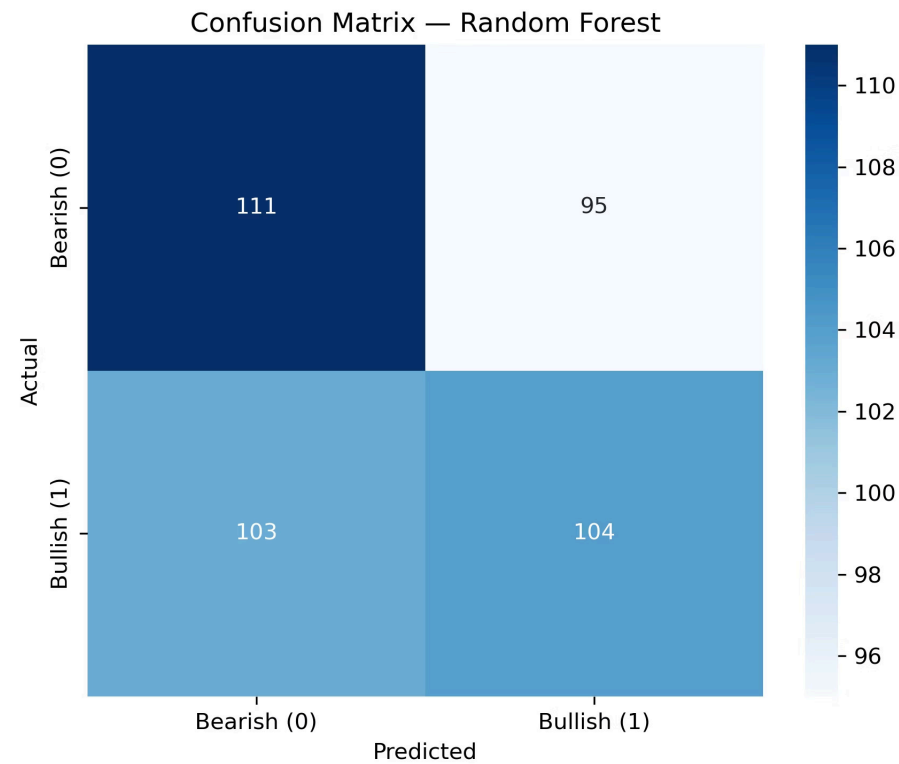## Model Performance Comparison



## Key Observations

- Random Forest achieved the highest test accuracy at 54.00%.
- All models clustered around 50-51% CV scores, indicating market efficiency.
- Performance suggests sentiment alone provides marginal predictive power.

## Performance Metrics Table

| Model | Best CV Score | Test Set Accuracy |
| --- | --- | --- |
| Logistic Regression | 0.5067 | 0.4939 |
| Random Forest | 0.5055 | 0.5400 |
| XGBoost | 0.5133 | 0.4673 |
| SVM | 0.5152 | 0.4625 |

## Statistical Interpretation

- Accuracy near 50% aligns with the Efficient Market Hypothesis.
- A 2-4% edge over random chance is statistically significant but practically limited.
- Results validate that cryptocurrency markets rapidly incorporate sentiment information.

Confusion Matrix — Random Forest

|  | Bearish (0) | Bullish (1) |
|---|---|---|
| **Bearish (0)** | 111 | 95 |
| **Bullish (1)** | 103 | 104 |

Actual (rows) / Predicted (columns)

# Confusion Matrix Analysis

## Performance Breakdown

While 52% accuracy appears underwhelming, context matters. In high-frequency trading, a **2% edge** over random chance represents statistical significance.

Our tuned Random Forest model demonstrated slight asymmetry in prediction accuracy, performing marginally better at identifying **bullish trends** compared to bearish movements.

However, this edge is insufficient as a standalone trading strategy without additional signals, risk management, and transaction cost considerations.

**Key Metrics**

- **True Positives:** Correctly identified upward movements
- **False Negatives:** Missed bullish opportunities
- **Precision/Recall:** Balanced but modest scores

# Why Not 90% Accuracy?

Our model's predictive accuracy, while statistically significant, did not reach the high percentages often desired. This outcome highlights fundamental realities of financial markets and inherent machine learning limitations.

## Data Limitation

A dataset of 2,000 rows is considerably small for training robust machine learning models, especially in complex and dynamic domains like financial markets. More extensive data would likely improve model generalization and performance.

## Market Stochasticity

Financial markets are inherently stochastic and non-deterministic. The pursuit of a "Holy Grail" of 90% predictable returns is a fallacy; markets are designed to be efficient, making consistent high predictability extremely difficult.

## Sentiment as a "Weak Learner"

While relevant, sentiment data often acts as a "weak learner" feature. It provides a signal, but it's just one piece of a much larger, noisy puzzle of market drivers, making it insufficient for high-accuracy predictions on its own.

Our results implicitly validate the Efficient Market Hypothesis: sentiment is rapidly incorporated into prices, and predicting crypto markets with simple data is profoundly challenging. The high 'noise' in crypto requires far more sophisticated approaches.

# Conclusion & Future Directions

## Key Findings

- **Marginal Predictive Edge:** Sentiment analysis provides a **2.1% advantage** over random guessing—statistically significant but insufficient as a standalone strategy.

- **Theoretical Validation:** Results empirically support market efficiency theory in cryptocurrency markets.

## Future Research Directions

- **LSTM Time-Series Models:** Incorporate temporal dependencies and momentum indicators to analyze price trends over sequential time windows rather than isolated snapshots.

- **Natural Language Processing:** Process raw social media content (tweets, Reddit posts) directly rather than relying on pre-aggregated sentiment scores, potentially capturing nuanced signals.

- **Hybrid Ensemble Approaches:** Combine sentiment models with technical analysis and on-chain metrics for comprehensive multi-signal strategy.

## Methodology Contributions & Closing

Successfully constructed a robust, professional ML pipeline handling complex financial data with proper preprocessing and validation.

**Thank you.** We welcome your questions and feedback on our methodology and findings.