# Domain Background

## Introduction to Stock Market Prediction

### Background and Importance

Stock market prediction has been a cornerstone of financial analysis and investment strategy for decades. Accurate forecasting of stock prices enables investors, hedge funds, and financial institutions to make informed decisions, maximize returns, and mitigate risks. The ability to predict market behavior, even within a narrow margin, provides a competitive edge in a field driven by rapid, high-stakes decision-making.

### Historical Developments

Historically, stock market prediction relied on traditional methods such as technical analysis, which focuses on historical price patterns, and fundamental analysis, which evaluates a company's financial health and industry position. Over the past few decades, the integration of computing power and the availability of vast amounts of financial data have transformed these practices.

With the rise of machine learning (ML), models have become capable of uncovering complex patterns in market data. Techniques like regression analysis, time series modeling, and neural networks are now widely used. Machine learning can process large datasets, identify non-linear relationships, and adapt to changing market conditions, making it a valuable tool for stock prediction.

Growth of Machine Learning in Stock Prediction

The application of machine learning in stock prediction accelerated in the early 2000s with advancements in data collection and computing power. Academic studies like those by Chen et al. (2003) on support vector machines and Hiransha et al. (2018) on deep learning for stock price prediction demonstrated the potential of ML to outperform traditional methods.

### Notable examples include:

- LSTM (Long Short-Term Memory) models: These are widely used for time series forecasting due to their ability to capture temporal dependencies. Studies have shown LSTM's effectiveness in predicting stock prices and trends.
- Quantitative Hedge Funds: Firms like Renaissance Technologies utilize algorithmic models and machine learning to trade successfully, underscoring the impact of ML in finance.

### Personal Motivation

This project is driven by a strong interest in the intersection of finance and machine learning. With the global market becoming increasingly data-driven, I am inspired to explore the potential of ML algorithms to make informed predictions and contribute to the growing body of work in this field. This project also provides an opportunity to apply advanced ML techniques to real-world problems, bridging theoretical knowledge with practical implementation.

## The Problem

The Problem

The goal is to develop a machine learning model that predicts the Adjusted Close price of stocks based on daily trading data over a specific date range. The Adjusted Close price accounts for stock splits and dividends, making it a reliable metric for analyzing stock performance and trends over time. The model will use features such as Open, High, Low, Volume, and Adjusted Close from historical data to generate predictions for future dates.

The Need for This Solution

The stock market is a complex and volatile system influenced by numerous factors, including economic trends, market sentiment, and geopolitical events. Predicting stock prices is inherently challenging, but accurate predictions can significantly impact investment strategies. By leveraging machine learning, this project aims to address key pain points in financial decision-making:

1. Informed Investment Decisions: Accurate stock price predictions enable investors to identify opportunities, optimize entry and exit points, and align their strategies with market trends.
2. Risk Reduction: Machine learning models can analyze vast amounts of historical data and detect patterns that human analysis might miss. This reduces the likelihood of poor investment choices based on incomplete or biased information.
3. Profit Maximization: By identifying trends and anticipating price movements, investors can capitalize on market fluctuations, maximizing returns on their investments.
4. Scalability and Efficiency: Traditional stock analysis is time-intensive and subjective. Machine learning models provide scalable, automated, and consistent predictions, allowing traders and firms to make decisions faster and more efficiently.

This solution is particularly relevant for individual day traders and small investment firms who lack access to the advanced tools and resources available to larger institutions. By making predictive tools more accessible, the project contributes to leveling the playing field in financial markets.

The implementation of this model not only meets the immediate need for actionable insights but also lays the foundation for advanced trading strategies, including portfolio optimization and risk analysis.

## Proposed Solution

The solution involves creating an API-based system for predicting the Adjusted Close price of stocks using a machine learning model. The model will be trained on historical stock price data and leverage features such as:

- Open price
- High price
- Low price
- Volume
- Adjusted Close price

The API will serve as the interface for querying stock price predictions, making it accessible for integration into various applications or workflows.

API Functionality

The API will handle POST requests and provide the following core functionality:

1. Training Interface:

- Input: A JSON payload containing a date range (start_date, end_date) and a list of ticker symbols (e.g., AAPL, GOOG).
- Action: The API will train the machine learning model using the specified historical stock price data.
- Output: A response indicating the training status (e.g., "Model trained successfully").

2. Prediction Interface:

- Input: A JSON payload with query dates and ticker symbols.
- Action: The API will use the trained model to predict the Adjusted Close prices for the specified stocks on the given dates.
- Output: A JSON response containing the predicted prices for each stock on each query date.

## Dataset Characteristics

- Fields Used:
- Open Price: The stock's price at the start of the trading day.
- High Price: The highest price the stock reached during the trading day.
- Low Price: The lowest price the stock reached during the trading day.
- Volume: The total number of shares traded during the trading day.
- Adjusted Close Price: The closing price adjusted for stock splits and dividends (target variable).
- Historical Data:
- The dataset will include historical daily trading data for multiple stocks over a user-specified date range. For instance, users can query stock data spanning several months or years to train the prediction model.
- Availability, Reliability, and Limitations:
- Availability: Yahoo Finance is widely accessible through unofficial means, and the yfinance library simplifies data retrieval.
- Reliability: Yahoo Finance provides robust and consistent financial data, but the unofficial nature of yfinance introduces potential risks, such as API rate limiting or website structure changes.
- Limitations:
- No guarantees of long-term API support due to its unofficial nature.
- Data may lack certain fields or granularity compared to premium providers (e.g., second-level trade data).

## Selected Benchmark Model

Moving Average Model The Moving Average model is selected as the benchmark for this project. It predicts the next day's Adjusted Close price by calculating the average of the Adjusted Close prices over a fixed window of past days (e.g., 5, 10, or 20 days).

**Why?** The Moving Average model is chosen because:

1. Simplicity: It is widely used in financial analysis and easy to implement.
2. Applicability: The model leverages historical stock prices, aligning with the project's dataset and goals.
3. Baseline for Performance: Provides a realistic baseline for evaluating the accuracy and effectiveness of advanced machine learning models.

## Evaluation Metrics

To evaluate the performance of both the benchmark model and the advanced machine learning model, the following metrics will be used:

1. Mean Absolute Error (MAE) MAE measures the average magnitude of errors in the predicted stock prices without considering their direction. It provides an intuitive understanding of the average prediction error in dollar terms.
2. Root Mean Squared Error (RMSE) RMSE emphasizes larger errors by squaring them before averaging. This is particularly relevant for stock price prediction, where large deviations from actual values can have significant financial implications.
3. Mean Absolute Percentage Error (MAPE) MAPE expresses errors as a percentage, allowing performance comparison across stocks with varying price ranges. This is useful when the dataset includes stocks with both low and high prices.
4. R-squared (Coefficient of Determination) R-squared indicates how well the model explains the variability in stock prices. It provides a measure of goodness-of-fit, with a value close to 1 indicating a strong correlation between predicted and actual prices.

**Why These Metrics?**

These metrics were chosen to provide a comprehensive evaluation of model performance. While MAE and RMSE quantify prediction accuracy in terms of dollar values, MAPE adds a percentage-based perspective. R-squared provides insight into the model's ability to explain variability in the data.

## Workflow

The project workflow is divided into the following stages:

1. Data Collection

- Input:
- Fetch historical stock price data using the yfinance Python library.
- Features include Open, High, Low, Volume, and Adjusted Close prices.
- Process:
- Query a dataset based on a user-specified date range and stock ticker symbols.
- Handle missing values and ensure data quality.
- Output:
- A clean and preprocessed dataset ready for model training.

2. Data Preprocessing

- Normalization:

- Scale numerical features (e.g., Open, High, Low, Volume) to a standard range (e.g., 0 to 1) for model compatibility.
- Feature Engineering:
- Add derived features, such as moving averages, volatility, or percentage changes.
- Train-Test Split:
- Split the dataset into training and testing sets based on dates to prevent data leakage (e.g., train on 2020–2022, test on 2023).
- Time-Series Considerations:
- Ensure sequential integrity by not shuffling data.

3. Model Development

- Model Selection:
- Start with basic models such as Linear Regression and Random Forest.
- Progress to advanced models like LSTM (Long Short-Term Memory) networks to capture temporal dependencies.
- Training:
- Use the training dataset to fit the model parameters.
- Implement cross-validation to assess robustness.
- Hyperparameter Tuning:
- Optimize parameters using techniques like grid search or Bayesian optimization.
- Benchmark Comparison:
- Compare the machine learning model's performance with the benchmark (e.g., Moving Average or Naive Prediction).

4. Model Evaluation

- Evaluate the trained model on the test set using metrics such as:
- Mean Absolute Error (MAE)
- Root Mean Squared Error (RMSE)
- Mean Absolute Percentage Error (MAPE)
- R-squared
- Visualize predictions against actual prices to assess trends and accuracy.

5. API Development

- API Features:
- Training Endpoint:
- Accepts POST requests with stock ticker symbols and date ranges to train the model.
- Prediction Endpoint:
- Accepts POST requests with query dates and ticker symbols.
- Returns predicted Adjusted Close prices in JSON format.
- Implementation:
- Ensure proper error handling, logging, and validation of input data.

6. Deployment

- Infrastructure:
- Deploy the API on a cloud platform such as AWS (API Gateway + Lambda).
- Scalability:
  - Allow concurrency
- Monitoring:
  - Implement logging and monitoring to track API performance and usage.

**Deliverables**

1. A functional API that provides predictions for Adjusted Close prices via POST requests.
2. Documentation for API usage, including input and output formats.
3. Performance Report comparing the benchmark and advanced models.
4. Codebase with clear modularity for data processing, model training, and inference.