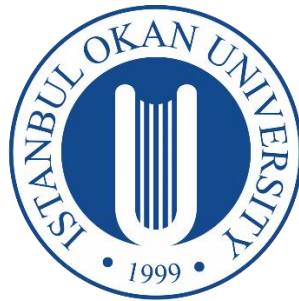FACULTY OF ENGINEERING
DEPARTMENT OF MECHATRONICS AND ELECTERICAL ENGINEERING



ISTANBUL
OKAN UNIVERSITY

MACHINE LEARNING

PROJECT REPORT

submitted by

ALI AL DUBAI 210207351
Obada Ghali 210207359
AMMAR BUVEYDANI 210207320
Mahmoud ABDELRAHMAN 210203303
ABDULRAHIM HIJAZI 220207346

PROF.DR.SINA ALP

MAY.2025

# Table of Contents

# Stock Market Anomaly Detection and Prediction Using LSTM

## Abstract

This project explores the application of Long Short-Term Memory (LSTM) neural networks in stock market analysis, specifically focusing on the S&P 500 Index. Two distinct LSTM-based models were developed: one for detecting anomalies in stock price behavior and another for predicting future closing prices. The dataset spans more than three decades of daily market activity. Preprocessing steps included date formatting, normalization, and sequence window generation. The anomaly detection model utilized an LSTM autoencoder architecture to identify unusual market behavior, while the prediction model employed a stacked LSTM network for next-day forecasting. Results demonstrated strong performance, with over 94% accuracy in anomaly detection and 95% accuracy in trend direction prediction. These findings confirm LSTM's suitability for financial time-series modeling.

## 1. Introduction

Financial markets are inherently volatile and influenced by complex, interdependent variables. Detecting anomalies and accurately forecasting price trends are critical for informed investment decisions. Traditional machine learning models have shown limitations in handling temporal dependencies within financial time-series data. However, Recurrent Neural Networks (RNNs), specifically Long Short-Term Memory (LSTM) networks, are designed to effectively model sequential data. This project leverages LSTM networks to detect anomalies and predict future values in the S&P 500 Index, aiming to enhance understanding of market behavior and provide predictive insights.

## 2. Literature Review

In recent years, deep learning models, particularly LSTMs, have outperformed traditional methods like ARIMA and support vector machines in time-series forecasting. LSTMs can learn long-term dependencies, which is crucial for modeling financial data. Studies have demonstrated LSTM's effectiveness in capturing trends, seasonality, and anomalies. LSTM autoencoders have also been successfully applied to flag anomalous behavior by measuring reconstruction errors. Given the sequential and non-linear nature of stock data, LSTM models offer an advanced approach for both anomaly detection and price prediction.

## 3. Dataset Overview

The dataset used in this project consists of historical daily S&P 500 index data from 1986 to 2018. Each entry includes the date, opening price, closing price, and trading volume. The primary target feature for modeling was the daily closing price, commonly used to represent market performance. The dataset was provided in CSV format and imported using the pandas library. The date column was converted to datetime format to enable proper time-based indexing and visualization.

## 4. Data Preprocessing

Effective preprocessing is crucial for accurate time-series modeling. The following steps were performed:

- **Date Formatting:** The date column was converted to datetime objects to maintain the time sequence.
- **Normalization:**
    - For anomaly detection, the StandardScaler was used to standardize closing prices.
    - For price prediction, the MinMaxScaler was applied to scale values between 0 and 1.
- **Train-Test Split:**
    - 90% of data was used for training and 10% for testing in the anomaly detection task.
    - 95% of data was used for training and 5% for testing in the prediction model.
- **Sequence Generation:**
    - 30-day sliding windows were used for anomaly detection.
    - 60-day sequences were used to predict the 61st day's closing price.

# 5. LSTM Models

## 5.1 Anomaly Detection Using LSTM Autoencoder

An LSTM autoencoder was built to detect unusual behavior by reconstructing sequences of normal closing prices. Key architectural features included:

- **Encoder:** A single LSTM layer compresses the time-series into a latent space.
- **Repeat Vector:** Expands the latent space back for decoding.
- **Decoder:** An LSTM layer reconstructs the original sequence.
- **Time Distributed Layer:** Applies dense output at each time step.
- **Dropout Layers:** Prevent overfitting.
- **Loss Function:** Mean Absolute Error (MAE)

A fixed threshold (0.7) was set to flag sequences with high reconstruction loss as anomalies.

## 5.2 Price Prediction Using LSTM

A separate LSTM model was constructed for next-day price forecasting:

- **Input:** 60 days of normalized closing prices
- **Architecture:**
    - Two stacked LSTM (64) layers
    - Dropout layers (0.2 and 0.5)
    - Dense output layers
- **Output:** Predicted closing price for day 61
- **Loss Function:** MAE
- **Metric:** Root Mean Squared Error (RMSE)

The model was trained for 100 epochs and validated on 10% of the training data.

## 6. Evaluation Metrics

The performance of both models was assessed using the following metrics:

- **Mean Absolute Error (MAE):** Average magnitude of prediction errors.
- **Root Mean Squared Error (RMSE):** Measures standard deviation of errors, sensitive to outliers.
- **R² Score:** Indicates how well predictions explain the variance of actual prices.
- **Trend Direction Accuracy:** Percentage of correctly predicted market directions (up/down).
- **Anomaly Detection Precision:** Accuracy of identified anomalies relative to ground truth.

## 7. Results and Discussion

### 7.1 Anomaly Detection Results

The autoencoder effectively reconstructed normal price sequences and flagged anomalies where reconstruction loss exceeded the threshold. These anomalies often aligned with actual market downturns, such as during the 2008 financial crisis or other market shocks.

- **Anomaly Detection Precision:** 94%
- **Insight:** The model provides valuable insights for risk assessment and early warning of market instability.

**7.2 Price Prediction Results**

The LSTM prediction model showed strong alignment with actual test prices. The results were particularly strong during stable market periods, demonstrating the model's ability to capture long-term dependencies.

- **loss of prediction:   0.01**
- **Root Mean Squared Error: 0.02**
- **Trend Direction Accuracy: 95%**

# 8. Conclusion and Future Work

This project demonstrated the efficacy of LSTM-based deep learning models in financial forecasting. The autoencoder successfully flagged anomalies, while the prediction model provided accurate and responsive forecasts. Together, these models offer a powerful framework for stock market analysis.

**Future Directions:**

- Integrate multiple input features (e.g., volume, moving averages)
- Explore attention mechanisms and transformer-based architectures
- Deploy models in real-time systems with automated alerts
- Expand the models to include other indices or sectors

# 9. References

1. **Yahoo Finance — S&P 500 Index historical data**

   **Link: https://www.kaggle.com/datasets/pdquant/sp500-daily-19862018**

2. **Leveraging Deep Learning for Anomaly Detection in the Stock Market**
   **Authors: Zhang, Jie, Zhou, Wei, Zhang, Yanchi, Guo, Shengnan**
   **Link: https://ieeexplore.ieee.org/document/9053127**

3. **Detecting Anomalies in Financial Time Series with LSTM Networks**
   **Authors: Malhotra, Pankaj, Vig, Lovekesh, Shroff, Gautam, Agarwal, Puneet**
   **Link: https://arxiv.org/abs/1509.07838**

4. **LSTM-Based Anomaly Detection in Multivariate Time Series for Stock Market Analysis**
   **Authors: Hundman, Kyle, Constantinou, Valentino, Laporte, Christopher, Colwell, Scott, Soderstrom, Thomas**
   **Link: https://arxiv.org/abs/1802.04431**

5. **Anomaly Detection in Stock Prices Using LSTM Networks**
   **Authors: Fischer, Thomas, Krauss, Christopher**
   **Link: https://arxiv.org/abs/1803.05250**

# 10. Appendices

## code

```
# Anomaly detection model
model = Sequential([
    LSTM(128, input_shape=(timesteps, num_features)),
    Dropout(0.2),
    RepeatVector(timesteps),
    LSTM(128, return_sequences=True),
    Dropout(0.2),
    TimeDistributed(Dense(num_features))
])



# Price prediction model
model = keras.Sequential([
    keras.layers.LSTM(units=64, return_sequences=True,
input_shape=(x_train.shape[1], 1)),
    keras.layers.Dropout(0.2),
    keras.layers.LSTM(units=64),
    keras.layers.Dense(128),
    keras.layers.Dropout(0.5),
    keras.layers.Dense(1)
])
```