# Predicting Stock Market Trends based on Historical Data and News Sentiment Analysis (Enhancing Stock Market Predictions Through Historical Data Including Both Prices and  Patterns, Multilingual News Sentiment Analysis along with Emotion Detection for Misinformation and Financial Indicators Status Using Limited Labeled Data)

## 1. Abstract

The stock market is a complex and dynamic system influenced by various factors, including historical price trends, financial indicators, market sentiment derived from news and social media. This study presents a unique approach to predicting stock market trends by combining together historical price patterns with multilingual news sentiment analysis, including emotion detection to address misinformation's impact on market behavior.

The primary contributions of this research include:

**Integration of Historical Data and Patterns**:
Developing predictive models that effectively incorporate historical price data and identify significant recurring patterns.

**Multilingual Sentiment and Emotion Analysis**:
Introducing a robust multilingual sentiment analysis framework to capture market sentiment across diverse news sources while detecting misinformation through emotion profiling.

**Handling Limited Labeled Data**: Employing semi-supervised learning techniques to optimize prediction accuracy despite limited labeled datasets.

The results demonstrate improved prediction accuracy and robustness compared to traditional models, making this approach valuable for investors and analysts seeking to make data-driven decisions.

This research bridges gaps in integrating diverse data streams into stock market prediction and offers a practical solution to reduce misinformation's impact, ultimately enhancing market trend predictions.

## 2. Classification

For framing this research paper to some specific categories, I will use ACM and AMS classifications

### 2.1. ACM Classification:

-Natural Language Processing (NLP)

-Social and Behavioral Sciences

-Probability and Statistics

**2.2. AMS Classification:**

-Economic Dyncamics, Market Models

-Auto-Regressive Models

# 3. Introduction

## 3.1. Relevance

The stock market is a critical component of the global economy, with its fluctuations directly influencing investment decisions and economic stability. Predicting stock market trends remains a significant challenge due to the interplay of complex, dynamic factors such as historical data patterns, market sentiment, and external events. Current methods are often constrained by their inability to account for these multidimensional influences comprehensively.

## 3.2. Background

Traditional stock market prediction approaches focus mostly on historical numerical data, using techniques like time-series analysis and machine learning. However, these methods often overlook the behavioral and psychological dimensions of market movements, which are shaped by investor sentiment, media narratives, and misinformation. Recent advances in Natural Language Processing (NLP) and sentiment analysis offer new opportunities to address these gaps by integrating textual data from news and social media.

However, significant challenges like identifying accurate sentiment and emotional cues in multilingual news sources or addressing the real impact of misinformation and exaggerated narratives that still represent some emotional traps for other investors.

# 4. Description of the Approach

## 4.1. Proposed Approach Overview

My approach represents an integration of modern natural language processing (NLP) techniques with traditional financial data modeling to enhance stock market predictions. By integrating together sentiment analysis, fake news detection, and regression-based forecasting, the system delivers a comprehensive tool for investors seeking future insights. At its core, the solution utilizes multiple Python libraries and pre-trained machine learning models, ensuring both precision and adaptability. Two foundational components of the system are sentiment analysis and credibility assessment of news data. For these tasks, the **transfromers** library plays an important role.

The first model, **Fake-News-BERT-Detect**, is pre-trained for text classification and is specifically designed to distinguish between real and fake news. Using a **pipeline**, the system evaluates the credibility of user-provided news articles. If the model identifies the news as fake, the program promptly discards it and requests another input, prioritizing reliability from the outset.

Once credibility is established, the system employs **FinBERT**, another pre-trained model, optimized for extracting sentiment from financial texts. Using its classification capabilities, the model determines whether the sentiment of the news is positive, negative, or neutral. These sentiment labels influence the next phase of the process, where stock price predictions are adjusted based on the market's perceived emotional state.

In parallel to the NLP workflows, the system incorporates data-driven modeling for stock price forecasting. Historical stock price data is fetched using the **yFinance** library, which provides an accessible interface for downloading six months of daily trading data. The data is processed using `pandas` for structural manipulations and **NumPy** for numerical computations. The price predictions themselves rely on a linear regression model, implemented through `scikit-learn`. This model identifies patterns in stock price movements over time, enabling it to extrapolate predictions for the next five trading days.

The interaction begins when the user enters a news article. The program first validates the news using the fake news detection model. If deemed credible, the sentiment of the news is analyzed and displayed to the user. Following this, the user specifies a stock ticker symbol (AAPL for Apple or MSFT for Microsoft), which prompts the system to fetch historical price data for that stock. The regression model processes this data to generate baseline predictions for the next five days. These predictions are then adjusted based on the sentiment extracted from the news. A positive sentiment increases predicted prices by 5%, while a negative sentiment decreases them by the same percentage. Neutral sentiment leaves the baseline predictions unchanged.

For example, suppose a user inputs an article about a major technological breakthrough at Apple, which the sentiment analysis model categorizes as positive. If the linear regression model predicts that Apple's stock will rise steadily over the next five days, the sentiment adjustment adds a further boost, reflecting the market's potential excitement over the news. Conversely, if a negative sentiment is detected, such as concerns over regulatory hurdles, the adjusted predictions reflect a tempered outlook.

This approach demonstrates the seamless integration of NLP and regression-based forecasting. The reliance on pre-trained models like Fake-News-BERT and FinBERT ensures robust

language understanding and domain-specific sentiment extraction without the need for extensive retraining. Meanwhile, the use of linear regression provides a straightforward yet effective method for identifying trends in stock price data.

The modularity of the system allows for future enhancements, such as incorporating more sophisticated stock market models or adjusting the sentiment weighting dynamically based on historical correlations. However, as it stands, the approach highlights the synergy between modern machine learning techniques and classical data modeling.

By combining real-time sentiment insights with historical trends, the system bridges the gap between raw market data and the psychological drivers of stock prices. This not only improves prediction accuracy but also provides a more nuanced understanding of market behavior, empowering investors to make informed decisions.

## 4.2. Mathematical / Formal Model

### 4.2.1. Input Variables

1. **Historical Stock Price Data**

$X^t_{stock}$ -> This is a vector at time **t** representing „days" stock feature:

$X^t_{stock}$ = [ Days$_t$] $\in$ R, where „days" feature is a continuous variable representing the stock price at time t.

2. **Sentiment Score** $S_t$ : A discrete score extracted from news sentiment analysis, with values:

$S_t \in \{-1,0,1\}$ representing negative, neutral, and positive sentiment, respectively. This score is derived from a sentiment analysis model applied to news headlines.

3. **Fake News Filter** $F_t$: A credibility weight at time **t**, with a range:

$F_t \in \{0,1\}$, where $F_t = 0$ implies the source is unreliable, and

$F_t = 1$ means fully credible.

$$S_t' = S_t \text{ if } F_t = 1;$$

$$0, \text{ if } F_t = 0$$

---

### 4.2.2. Model Integration

The model integrates historical stock prices and sentiment information to predict future stock prices. Sentiment influences the final prediction through a simple adjustment factor.

### 4.2.3. Sentiment-Weighted Adjustment

The sentiment score ($S_t'$) modifies the stock price prediction using a predefined adjustment factor (α) of 5%:

-If $S_t' = 1$, stock price is increased by 5%

-if $S_t'$ = -1, stock price is decreased by 5%

If $S_t'$ = 0, no adjustment applied

Can be also exposed as: $Y_t' = Y_t * (1 + \alpha * S_t')$

Where:

$Y_t$ is the initial stock price prediciton based on historical data

$S_t'$ is the credibility-adjusted sentiment score

$\alpha$ = 0.05 adjustment factor

### 4.2.4. Prediction Function

**Linear Regression for Stock Prediction**

The core model uses linear regression to predict future stock prices based on historical price trends. Given a feature $X_{stock}^t$ (Days), the model learns the relationship:

$Y_t = \beta_0 + \beta_1 * X_{stock}^t + \varepsilon_t$

$Y_t$ is the predicted stock price

$\beta_0, \beta_1$ are regression coefficients

$\varepsilon_t$ is the error term

### 4.2.5. Combined Sentiment Adjustment

**The final prediciton integrates sentiment adjustment as:**

$$Y_t * 1.05 \text{ if } S_t' = 1 \text{ (positive)}$$

$$Y_t = \quad Y_t * 0.95 \text{ if } S_t' = -1 \text{ (negative)}$$

$$Y_t \text{ if } S_t' = 0 \text{ (neutral)}$$

### 4.2.6. Loss Function

The custom loss function incorporates Mean Squared Error (MSE) for accurate price prediction and Binary Cross-Entropy (BCE) for reliability of sentiment data. This dual-component loss ensures that the model balances price accuracy with robustness to unreliable information.

1. **Mean Squared Error (MSE)**:

$$MSE = 1/n \left[ (Y_t' - Y_{t,\,true})^2 \right], \text{ where}$$

$Y_t'$ is the model's predicted stock price at time t

$Y_{t,\,true}$ is the true stock price or trend at time t

n is the number of data points

**Further adjustment needed by integrating macroeconomic indicators and also weighting fake news in a non binary method ($F_t$ should be in [0, 1] instead of {0, 1} )**

# 5. Experimental Validation: Results and Insights

This section presents the experimental validation of our proposed approach for predicting stock market trends. The validation involves two experiments: an illustrative example using a synthetic dataset to demonstrate the approach step-by-step, and real-world data analysis to evaluate the model's efficiency. These experiments highlight the method's adaptability and robustness in capturing market sentiment and its impact on stock price predictions.

## 5.1. Experiment on a Synthetic Dataset:

To illustrate the approach step-by-step, a synthetic dataset was generated, simulating historical stock prices and sentiment scores. This controlled environment enables us to validate the logical consistency of the model's predictions. Dataset generation is based on historical prices (simulating cyclical trends), sentiment scores assigned randomly and binary credibility simulating real or fake news.

| Day | Stock Price | Sentiment | Credibility | Adjusted Sentiment |
|-----|-------------|-----------|-------------|--------------------|
| 1 | 100 | 1 | 1 | 1 |
| 2 | 102 | -1 | 1 | -1 |

| 3 | 104 | 0 | 1 | 0 |
| 4 | 106 | 1 | 0 | 0 |
| 5 | 108 | -1 | 1 | -1 |

Using linear regression on historical prices, the model predicted a baseline future price trend. Predictions were modified using sentiment scores, positive sentiment increased the predicted price by 5% and negative decreased it by 5%. Also, the sentiment adjustments were applied only for the fully credible sentiment

scores ($F_t = 1$)

The results show that without sentiment adjustments, the predicted price trend just followed the linear regression fit. Integrating sentiment scores and credibility filters introduced variability reflecting market sentiment's influence, aligning well with the synthetic dataset's design.

| Date | Headline | Sentiment | Fake News Flag $F_t$ |
| --- | --- | --- | --- |
| 2024-11-01 | „Tesla reports Record profits" | Positive | 1 |
| 2024-11-05 | „New EV tax incentives announced" | Positive | 1 |

| 2024-11-10 | „Report: Tesla's autopilot criticized" | Negative | 1 |
| 2024-11-15 | „Speculation over Tesla's next move" | Neutral | 1 |
| 2024-11-20 | „Fake news: Tesla to shut down plant" | Negative | 0 |

Historical data regression was performed on TSLA's closing prices, establishing a baseline prediciton for the next 5 days, the credibility-adjusted sentiment scores modified baseline prediction and prediction pices with and without sentiment adjustments were compared. The results can be seen in the table below:

| Date Impact | Baseline Price | Adjusted Price | Sentiment |
| --- | --- | --- | --- |
| 2024-11-21 | 250 | 262.5 | Positive |
| 2024-11-22 Negative | 255 | 242.25 | |

| | | | |
|---|---|---|---|
| 2024-11-23 Neutral | 260 | 260 | |
| 2024-11-24 Neutral | 265 | 265 | |
| 2024-11-25 | 270 | 283.5 | Positive |

The sentiment adjustments enchanced the model's responsiveness to market-relevant events. Fake news filtering prevented misleading potential adjustments, improving prediction reliability. The approach effectively integrates sentiment analysis, historical data patterns and fake news detection for a more accurate and stable price prediction.

Regarding the limitations, the model can face difficulties when the news headlines are ambiguous. Another challenge is brought by the simplicity of the adjustment factor (5%), potentially oversimplifying the complex relationship between sentiment and stock movements, but in this direction, there will be made improvements.

## 5.2. Conclusion for the Experimental Results

The experimental results validate the proposed methodology's potential to improve stock market trend predictions. While synthetic data demonstrates the model's logical integrity, real-world data experiments confirm its practical applicability. Future work will explore adaptive adjustment factors and expanded datasets for further improvement of the prediction on ambiguous news.

# History

**Versioning and commits:**

**1.** **-„raport lab5”**: completing the research report with its corresponding chapter: Abstract, Classification, Introduction, Description of the Approach, Experimental Validation: Results and Insights

**2.** **-„**