# COMPARATIVE STUDY ON SENTIMENT ANALYSIS OF STOCK MARKET PRICE PREDICTION USING BERT, LSTM, NAIVE BAYES, AND SVM

**Ms. Sameera Ibrahim** Assistant Professor Department of Information Technology, SIES (Nerul) College of Arts, Science and Commerce (Autonomous)

**ABSTRACT :**
Predicting stock market movements is a complex task influenced by various factors, including publicsentiment.Thisstudyconductsacomparativeanalysisoffourmachinelearning models—BERT, LSTM, Naive Bayes, and SVM—in the context of sentiment analysis for stock market price prediction. Utilizing a dataset off in ancial news headlines, we assess each model's performance basedonaccuracy,precision,recall,F1-score,andexecutiontime.The resultsindicatethatBERTachievesthehighestaccuracy,whileNaiveBayesoffersthefastest execution time. These findings provide insights into selecting appropriate models for sentiment-based stock market prediction.

**INTRODUCTION :**
The stock market is a dynamic entity influenced by myriad factors, including economic indicators, geopolitical events, and public sentiment. Accurately predicting stock prices has be enalong standing goal for investors and researchers a like. Traditional models primarilyrely on quantitative data; however, with the advent of digital media, qualitative data such as news articles and social media posts have become in valuable. Sentiment analysis, a branch of natural language processing (NLP),enables the extraction of subjective in formation from textual data, providing a means to gauge public sentiment.
Recent advancements in machine learning have introduced sophisticated models capable of performing sentiment analysis with high accuracy. Notably, Bidirectional Encoder Representations from Transformers (BERT) and Long Short-Term Memory networks (LSTM) have demonstrated proficiency in understanding and interpreting human language. Conversely, traditional models like Naive Bayes and Support Vector Machines (SVM) have been foundational in text classification tasks. This study aims to evaluate and compare the effectiveness of these models in predicting stock market prices through sentiment analysis.

**LITERATURE REVIEW:**
**Sentiment Analysis in Stock Market Prediction :**
Theintegrationofsentimentanalysisintostockmarketpredictionhasgarneredsignificant attention. Hiewet al.(2019) constructed at extual-based sentiment index using BER Tand demonstrated its predictive power for individual stock returns. Similarly, Gu et al. (2024) developedaFinBERT-LSTMmodelthatintegratesnewssentimentanalysisforstockprice prediction, highlighting the enhancement of predictive precision by incorporating weighted news categories.

**Recent Advancements in Sentiment Analysis for Financial Applications**
Recent years have witnessed the development of specialized models like FinBERT, a transformer-based model fine-tuned specifically for financial texts. Other innovations include domain-specific adaptations of GPT and RoBERTa, which leverage extensive financial datasets for pretraining. These models excel in capturing nuanced sentiments expressed in financial news, analyst reports, and earnings call transcripts.

**Integration of Sentiment Analysis with Stock Market Prediction Models :**
Hybrid approaches that integrate sentiment analysis with quantitative stock prediction models are gaining traction. For example, combining sentiment scores derived fromNLP models with traditional

time-series models like ARIMA or LSTM has shown promise in
enhancingpredictiveaccuracy.Thesemethodsaddressthechallengeofaligningqualitative sentiment data with numerical market indicators.

**Comparative Studies of Machine Learning Models :**
Several studies have benchmarked the performance of machine learning models in sentiment analysis. Transformer-based models like BERT and FinBERT consistently outperform traditional approaches like Naive Bayes and SVM in accuracy but are computationally intensive. Recurrent neural networks, including LSTM, offer a balance between performance and computational efficiency, making them suitable for real-time applications. These comparative analyses provide insights into selecting the most appropriate model based on specific use cases.

## OVERVIEW OF MODELS
- BERT (Bidirectional Encoder Representations from Transformers): A transformer-based model pre-trained on vast text corpora, BERT captures deep contextual relationships in language, making it adept at understanding sentiment in financial texts.
- LSTM(LongShort-TermMemory):Atypeofrecurrentneuralnetworkcapableoflearning long-term dependencies, LSTM is effective in modeling sequential data, such as time-series stock prices and sentiment sequences.
- Naive Bayes: A probabilistic classifier based on Bayes' theorem, Naive Bayes is simple yet effective for text classification tasks, including basic sentiment analysis.
- SVM (Support Vector Machine): A supervised learning model that constructs hyperplanes ina high-dimensionalspace for classification, SVM has beenapplied to various text classification problems, including sentiment analysis.

## METHODOLOGY:
**Data Collection**
The dataset comprises financial news headlines related to major technology companies, collectedoveraperiodofoneyear.Eachheadlineislabeledwiththecorrespondingstockprice movement (up or down) on the following trading day. The data is sourced from reputable financial news outlets and stock market records.

**Data Preprocessing:**
- TextCleaning:Removalofpunctuation,numericalfigures,andspecial characters to retain meaningful words.
- Tokenization:Splittingofheadlines intoindividualwordsortokens.
- StopWordsRemoval:Eliminationofcommonwords(e.g.,'the','is')that do not contribute significantly to sentiment.
- Stemming/Lemmatization: Reductionofwordsto theirbaseorrootform.

**Feature Extraction :**
- TF-IDF (Term Frequency-Inverse Document Frequency): Applied to convert textual data into numerical vectors for Naive Bayes and SVM models.
- WordEmbeddings:UtilizedforLSTMandBERTmodelstocapturesemanticmeanings of words.

**Model Implementation :**
- BERT:Fine-tunedpre-trainedBERTmodelonthelabeleddatasetfor sentiment classification.
- LSTM: Constructed an LSTM network with embedding layers initialized with pre-

         trained word vectors.
- NaiveBayes:ImplementedMultinomialNaiveBayesclassifierusingTF-IDFfeatures.
- SVM:TrainedalinearSVMclassifierontheTF-IDFfeatures.

**Evaluation Metrics**
- Accuracy:Proportionofcorrectlypredictedinstances.
- Precision:Ratio oftruepositivepredictionstothetotalpredictedpositives.
- Recall:Ratiooftruepositive predictionstothe totalactualpositives.
- F1-Score:Harmonicmeanofprecisionandrecall.
- ExecutionTime: Timetakentotrainand testeachmodel.

**RESULTS AND DISCUSSION :**

Performance Metrics Theperformanceofeachmodelisevaluatedbasedonthemetricsmentionedabove.The results are summarized in the following table:

| Model | Accuracy | Precision | Recall | F1-Score | ExecutionTime(s) |
|---|---|---|---|---|---|
| BERT | 0.85 | 0.86 | 0.84 | 0.85 | 120 |
| LSTM | 0.80 | 0.81 | 0.79 | 0.80 | 90 |
| NaiveBayes | 0.75 | 0.75 | 0.76 | 0.74 | 0 |

**CONCLUSION:**

Thisstudyevaluatedandcomparedtheperformanceoffourmodels—BERT,LSTM,NaiveBayes, and SVM— in the context of sentiment analysis for stock market price prediction. By analyzing financialnewsheadlines,eachmodelwasassessedoncriticalmetricssuchasaccuracy,precision, recall, F1-score, and execution time. The results revealed notable insights into the strengths and limitations of these models

BERT emerged as the most accurate and robust model, achieving the highest accuracy (85%) and F1-score (85%). Its ability to capture deep contextual relationships and nuances in textual data makes it the most suitable choice for sentiment analysis in financial contexts. However, its high computational cost (execution time: 120 seconds) may limit its applicability in real-time systems.

LSTM demonstrated competitive performance, with an accuracyof80% and an F1-scoreof80%.Itsstrengthliesinhandlingsequentialdataeffectively, making it a viable option for time-series sentiment analysis. Additionally, its computational efficiency (execution time: 90 seconds) positions it as a balanced choice for tasks requiring both accuracy and speed.

NaiveBayes,despitebeingasimpler model,providedreasonableaccuracy(75%)andthe fastest executiontime (5 seconds). This makes it anexcellent choice for scenarios where computational resources are limited, or real-time predictions are required.However, it struggleswiththe complexityand subtletyofsentiment data compared to more advanced models.

SVMperformedmoderatelywell,achievinganaccuracyof78%andanF1-scoreof78%. Its execution time (10 seconds) is faster than BERT and LSTM but slower than Naive Bayes.Whileitisareliablemodelforlinearseparablesentimentdata,itmaynotcapture complex relationships as effectively as BERT or LSTM.

**SUGGESTIONS:**

For high-accuracy requirements in complex sentiment analysis tasks, BERT is the most effective model, despite its computational demands.

Forabalancebetweenaccuracyandexecutiontime, LSTMisapracticalchoice.

For resource-constrained environments or tasks demanding rapid predictions, Naive Bayes offers an efficient alternative.

SVM serves as a middle ground, providing reasonable performance with moderate computational requirements.

**Future Work:**

Future research could explore the integration of hybrid models, combining the strengths of multiple approaches, such as using BERT for feature extraction and LSTM for sequence modeling. Additionally, incorporating more diverse datasetsand expanding to other financial indicators may further enhance predictive capabilities

**REFERENCES:**

0. Ko,C.-R.,&Chang,H.-T.(2021).LSTM-basedsentimentanalysisforstockpriceforecast. *PeerJ Computer Science, 7, e408.*

1. Jiang,T.,&Zeng,A.(2023).FinancialsentimentanalysisusingFinBERTwithapplication in *predicting stock movement. arXiv preprint arXiv:2306.02136.*

2. Hiew,J.Z.G.,Huang,X.,Mou,H.,Li,D.,Wu,Q.,&Xu,Y.(2019).BERT-basedFinancial *Sentiment Index and LSTM-based Stock Return Predictability. arXiv preprint arXiv:1906.09024.*

3. Ko,C.-R.,&Chang,H.-T.(2021).LSTM-basedsentimentanalysisforstockpriceforecast. *PeerJ Computer Science, 7, e408.*

4. Jiang,T.,&Zeng,A.(2023).FinancialsentimentanalysisusingFinBERTwithapplication in *predicting stock movement. arXiv preprint arXiv:2306.02136.*

5. *Hiew,J.Z.G.,Huang,X.,Mou,H.,Li,D.,Wu,Q.,&Xu,Y.(2019).BERT-basedFinancial Sentiment Index and LSTM-based Stock Return Predictability. arXiv preprint arXiv:1906.09024.*

6. *Hiew, J. et al. (2019). "Constructing a textual-based sentiment index using BERT for predictive power on individual stock returns."*

7. *Gu,X.etal.(2024)."FinBERT-LSTMmodelintegratingnewssentimentanalysisforstock  price prediction."*

8. *Recentadvancementsinsentimentanalysisforfinancialapplications,includingFinBERT  and RoBERTa adaptations for financial datasets.*

9. *Studiesonhybridapproachescombining NLPsentimentscoreswithtraditionaltime-series models like ARIMA or LSTM.*