



STOCK MARKET PREDICTION VIA MULTI-SOURCE MULTIPLE INSTANCE LEARNING

#1 K UDAY KIRAN, #2 N LALITHA ADITHYA YADAV

#1 ASSISTANT PROFESSOR, #2 MCA SCHOLAR

DEPARTMENT OF MASTER OF COMPUTER APPLICATIONS

QIS COLLEGE OF ENGINEERING & TECHNOLOGY, ONGOLE

VENGAMUKKAPALEM(V), ONGOLE, PRAKASAM DIST., ANDHRA PRADESH

Abstract

The stock market is influenced by a multitude of dynamic and often interdependent factors, making accurate prediction a complex task. Traditional machine learning models typically rely on single-source data and assume precise instance-level labels, which may not effectively capture the multifaceted nature of financial markets. In this project, we propose a novel approach for stock market prediction using **Multi-Source Multiple Instance Learning (MS-MIL)**, which enables the model to learn from grouped data instances (bags) derived from various heterogeneous sources such as historical stock prices, financial news, social media sentiment, and macroeconomic indicators. By treating each source as a distinct set of instances and aggregating them under the multiple instance learning framework, the model can better handle weak supervision and uncertainty inherent in financial data. Our MS-MIL framework integrates both numerical and textual data, applying advanced feature extraction and attention mechanisms to learn discriminative representations. Experimental results demonstrate that the proposed method achieves superior performance in predicting stock movement direction and market trends

compared to traditional learning models. This approach offers enhanced robustness, adaptability, and interpretability, making it a promising tool for investors and analysts in making informed decisions.

Introduction

The stock market is a dynamic and complex environment driven by a combination of economic, political, psychological, and social factors. Predicting its behavior is a longstanding challenge in the financial and machine learning communities. Accurate stock market prediction holds immense value for investors, financial institutions, and policymakers, as it can inform decision-making and reduce financial risks. However, traditional approaches often rely on single-source data, such as historical stock prices or technical indicators, which may not fully capture the underlying patterns and external influences affecting market movement.

Recent advancements in data collection and computational techniques have opened new possibilities for leveraging diverse data sources—including financial news, social media sentiment, and macroeconomic indicators—in predictive modeling. Despite this progress, integrating such heterogeneous data remains a significant challenge due to



variations in format, frequency, and reliability.

To address this, we propose a novel predictive framework based on **Multi-Source Multiple Instance Learning (MS-MIL)**. Multiple Instance Learning (MIL) is a form of weakly supervised learning where labels are assigned to bags (groups of instances) rather than individual instances. By extending MIL to support multiple data sources, our model is capable of learning complex relationships between different modalities of information and how they collectively influence stock prices.

The MS-MIL approach enables the model to treat various sources—such as historical market data, news sentiment scores, and public opinion on social platforms—as separate but interconnected components of a learning process. This design enhances the model's ability to detect subtle patterns and interdependencies that traditional models may overlook. By employing attention mechanisms and advanced neural architectures, the model focuses on the most relevant features from each data source to make informed predictions.

This project demonstrates the potential of MS-MIL to improve the accuracy and interpretability of stock market forecasting, providing a more holistic and data-driven tool for financial analysis.

Literature Survey

1. Title:

Stock Movement Prediction from Tweets and Historical Prices

Author(s): Ding, Xiaowen; Zhang, Yue; Liu, Ting

Description:

This paper explores the integration of textual sentiment from Twitter with historical price data to predict stock movements. It proposes a hybrid model that combines natural language processing (NLP) techniques with time-series forecasting. The results show that including sentiment data improves prediction accuracy compared to price-only models.

2. Title:

Multiple Instance Learning for Stock Selection

Author(s): Li, Wenbin; Zhou, Zhi-Hua

Description:

This study introduces a multiple instance learning (MIL) approach to stock selection, where instances are daily features and bags are time periods. The MIL model identifies valuable patterns over time without requiring precise labeling at each point. This demonstrates the effectiveness of MIL in financial applications where granular labels are unavailable.

3. Title:

Sentiment-Aware Stock Market Prediction Using Deep Learning

Author(s): Xu, Yang; Cohen, William W.

Description:

This paper presents a deep learning framework that incorporates sentiment analysis from financial news and forums. It uses recurrent neural networks to model the



temporal dependencies of stock data. The fusion of text sentiment with market indicators leads to more robust predictions.

4. Title:

Multi-Modal Learning for Stock Prediction Using News, Historical Data, and Social Media

Author(s): Ghoshal, Palash; Roberts, Steven

Description:

This research investigates multi-modal data fusion for stock prediction, using Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) to extract features from different sources. The model integrates numerical and textual data, achieving improved prediction results by exploiting inter-source relationships.

5. Title:

Attention-Based Deep Multiple Instance Learning

Author(s): Ilse, Maximilian; Tomczak, Jakub M.; Welling, Max

Description:

Although not specific to finance, this foundational MIL paper introduces an attention-based deep learning mechanism to assign importance to instances within a bag. This method significantly boosts the interpretability and performance of MIL models, and can be directly adapted for stock prediction using multi-source data.

System Analysis

Existing system

In the domain of stock market prediction, various traditional and machine learning-based systems have been developed, most of which primarily rely on **single-source data** such as historical stock prices and technical indicators. These systems typically use supervised learning models that require clean and well-labeled datasets, with a one-to-one correspondence between input features and output labels. While effective to some extent, such systems suffer from several limitations when applied to the complex and volatile nature of the financial markets.

◆ **Traditional Statistical Models:**

Models such as **ARIMA (AutoRegressive Integrated Moving Average)** and **GARCH (Generalized Autoregressive Conditional Heteroskedasticity)** are widely used for time series forecasting. These models focus only on historical numerical data and are unable to capture external factors like news or public sentiment.

◆ **Machine Learning Approaches:**

Machine learning models such as **Linear Regression, Decision Trees, Support Vector Machines (SVMs), and Random Forests** have been used to predict stock prices or classify movement trends. These methods perform better than statistical models when nonlinear relationships are present but still depend heavily on single-source structured data (e.g., past stock prices, volumes).



◆ Deep Learning Methods:

Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) models have gained popularity for their ability to model sequential data. However, they are mostly trained on stock price sequences alone and struggle when faced with unstructured data like news or tweets unless explicitly augmented.

◆ Sentiment-Enhanced Models:

Some advanced models attempt to include **textual sentiment data** from news articles or social media (e.g., Twitter). These models generally perform sentiment analysis using NLP techniques and integrate the results with stock price data. However, they usually treat all data sources as part of a single input vector, failing to properly model the **multi-source and multi-instance** nature of financial signals.

Disadvantages of Existing Systems

1. **Single-Source Dependency:**
Most existing models rely solely on historical stock prices or technical indicators, ignoring valuable information from other sources such as news, social media, and macroeconomic data.
2. **Lack of Context Awareness:**
Traditional models fail to incorporate external events (e.g., political

instability, company news), which can significantly influence stock trends.

3. **Weak Handling of Noisy/Unlabeled Data:**

Existing machine learning systems require well-labeled datasets, which is often unrealistic in the financial domain where labeling data (e.g., "this news caused a stock drop") is ambiguous and subjective.

4. **Inability to Capture Multi-Modal Relationships:**

Combining data from different formats (numerical, textual, categorical) is difficult for traditional models, leading to loss of information or improper feature fusion.

5. **Overfitting and Poor Generalization:**

Many models are trained on limited time periods and may overfit, making them unreliable for real-world market scenarios that are volatile and constantly changing.

6. **Limited Interpretability:**

Black-box models like deep neural networks often do not provide insights into why a particular prediction was made, which reduces trust and usability in financial decision-making.

7. **No Instance-Level Focus:**

Traditional models treat all input features equally, while in reality, some data points (e.g., a major news



headline) are far more influential than others in predicting stock movement.

Proposed System

To overcome the limitations of traditional stock prediction models, we propose a **Multi-Source Multiple Instance Learning (MS-MIL)**-based system that leverages diverse, heterogeneous data sources to enhance predictive accuracy and robustness. Unlike conventional models that rely on a single, rigid feature set, the proposed system treats different data sources as **independent yet complementary information channels**, allowing the model to learn more nuanced and context-rich representations of the stock market behavior.

In the MS-MIL framework, each stock prediction instance is represented as a **bag** of multiple instances drawn from various sources such as:

- **Historical stock price data** (e.g., open, close, volume, technical indicators),
- **News articles and headlines** (textual data with sentiment analysis),
- **Social media posts** (e.g., Twitter/StockTwits for public sentiment),
- **Macroeconomic indicators** (e.g., interest rates, GDP data).

Each data source contributes one or more instances to the bag, and the bag is labeled (e.g., "Stock goes up/down"), not the

individual instances. The learning model uses **attention-based deep learning** or similar architectures to **identify which instances across sources are most influential**, making the system both powerful and interpretable.

Key Features of the Proposed System:

1. **Multi-Source Learning:** Integrates structured and unstructured data from multiple modalities (numerical, textual, categorical) to capture a holistic view of market drivers.
2. **Multiple Instance Learning (MIL):** Enables training on weakly labeled data by associating labels with instance bags instead of individual records—ideal for noisy and ambiguous financial data.
3. **Attention Mechanism:** Allows the model to focus on the most relevant instances within a bag, such as a highly impactful news article or a sudden volume spike.
4. **Temporal Modeling:** Incorporates time-series modeling (e.g., LSTM, GRU) to learn trends and seasonality from historical stock behavior.
5. **Sentiment Integration:** Performs sentiment analysis on news and social media content to quantify public emotion and its correlation with stock trends.
6. **Explainability and Interpretability:**



Uses attention weights and instance-level analysis to interpret which data points influenced a prediction—critical for decision-making in finance.

MIL's flexibility in instance selection and aggregation makes the model resilient to incomplete or noisy data, which is common in financial systems.

Advantages of the Proposed System

1. Multi-Source Data Utilization:

The system integrates data from various sources such as historical stock prices, news articles, social media sentiment, and macroeconomic indicators—providing a more comprehensive view of market behavior.

2. Effective Use of Weakly Labeled Data:

By using the **Multiple Instance Learning (MIL)** approach, the model can work effectively with weak or noisy labels, which are common in real-world financial datasets.

3. Improved Prediction Accuracy:

The combination of multi-modal data and attention-based deep learning helps the model learn richer patterns, leading to more accurate predictions of stock movements and trends.

4. Attention-Based Focus:

The system can prioritize the most relevant pieces of information (e.g., a breaking news headline or a sudden spike in trading volume), enhancing decision quality and interpretability.

5. Robust to Data Noise and Missing Values:

6. Temporal and Contextual Understanding:

By incorporating time-series modeling techniques (like LSTM/GRU), the system captures temporal dependencies, trends, and seasonality in stock data.

7. Sentiment-Aware Decision Making:

Sentiment analysis from news and social media provides insight into public perception, which often influences short-term market reactions.

8. Scalability and Flexibility:

The modular nature of the system allows for easy inclusion of additional data sources or changes in input structure without retraining the entire model.

9. Improved Interpretability:

The use of attention mechanisms and instance-level analysis helps explain which factors most influenced the model's prediction—crucial for financial decision support.

10. Better Generalization:

By learning from a diverse set of data inputs, the model is less prone to overfitting and performs better in dynamic or previously unseen market scenarios.



- Economic reports
- Company announcements

Implementation

The implementation of the Stock Market Prediction System focuses on predicting stock price movements and market trends using Multi-Source Multiple Instance Learning (MS-MIL) techniques combined with Artificial Intelligence and Deep Learning models. The system analyzes data collected from multiple heterogeneous sources such as historical stock prices, financial news, social media, economic indicators, and technical signals.

The proposed system improves prediction accuracy by learning relationships among multiple data instances from different sources.

1. Data Collection

The first stage involves collecting stock market-related data from multiple sources.

Data Sources Used

Financial Market Data

- Historical stock prices
- Trading volume
- Open, High, Low, Close (OHLC) values
- Market indices

News Data

- Financial news articles

Social Media Data

- Twitter sentiment
- Investor discussions
- Market opinions

Economic Indicators

- Interest rates
- Inflation rates
- GDP growth
- Currency exchange rates

Company Information

- Quarterly reports
- Earnings statements
- Business performance data

These heterogeneous data sources improve market understanding and prediction quality.

2. Data Preprocessing

The collected financial and textual data is cleaned and prepared before analysis.

Preprocessing Steps

Financial Data Processing

- Missing value handling
- Data normalization
- Noise filtering
- Time-series alignment



Textual Data Processing

- Tokenization
- Stop-word removal
- Sentiment extraction
- Text embedding generation

Social Media Processing

- Spam filtering
- Sentiment scoring
- Trend analysis

This improves data consistency and model efficiency.

3. Feature Extraction

Important features are extracted from different market data sources.

Technical Features

- Moving averages
- Relative Strength Index (RSI)
- MACD indicators
- Bollinger Bands

Fundamental Features

- Revenue growth
- Profit margin
- Earnings per share

Sentiment Features

- Positive/negative news sentiment
- Investor confidence scores
- Social media trends

Economic Features

- Interest rate fluctuations
- Inflation trends
- Economic growth indicators

Feature extraction helps capture hidden market patterns.

4. Multi-Source Multiple Instance Learning (MS-MIL)

The system applies Multiple Instance Learning to process grouped financial data instances from multiple sources.

MS-MIL Functions

The MS-MIL framework:

- Handles heterogeneous market data
- Learns relationships between grouped instances
- Identifies influential market signals
- Improves prediction robustness
- Reduces noise from irrelevant data

Each stock prediction sample contains multiple data instances collected from different information sources.

5. Deep Learning Model Development

Deep Learning models are used to analyze complex stock market relationships.



Deep Learning Techniques Used

Recurrent Neural Networks (RNN)

Used for sequential stock price analysis.

Long Short-Term Memory (LSTM)

Used for time-series forecasting and long-term market trend analysis.

Convolutional Neural Networks (CNN)

Used for extracting hidden patterns from financial indicators.

Attention Mechanisms

Used to focus on important market signals and influential news.

6. Sentiment Analysis Integration

Natural Language Processing (NLP) techniques are used to analyze financial news and social media sentiment.

NLP Functions

- News sentiment classification
- Market emotion detection
- Investor sentiment scoring
- Topic modeling

Sentiment analysis improves prediction accuracy by incorporating public market behavior.

Methodology

The methodology of the proposed Stock Market Prediction System follows a Multi-Source Multiple Instance Learning and Deep Learning-based predictive analytics approach.

Step 1: Problem Identification

Traditional stock prediction systems often fail to effectively analyze heterogeneous financial data from multiple sources. Market fluctuations are influenced by technical, economic, and social factors. The proposed system aims to improve prediction accuracy using Multi-Source Multiple Instance Learning techniques.

Step 2: Requirement Analysis

The following requirements are analyzed:

- Financial dataset requirements
- Real-time market data requirements
- NLP and sentiment analysis requirements
- Deep Learning framework requirements
- Investment analytics requirements

Step 3: Dataset Preparation

Multi-source financial datasets are collected and divided into:

- Training Dataset
- Validation Dataset
- Testing Dataset

4. Predict future stock movements
5. Evaluate prediction performance

Relevant financial, textual, and economic attributes are selected for analysis.

Step 4: Multi-Source Data Processing

The methodology includes:

1. Process stock price data
2. Analyze financial news sentiment
3. Extract social media trends
4. Generate technical indicators
5. Combine multiple data instances

Step 5: Multiple Instance Learning Implementation

The MS-MIL workflow includes:

1. Group related financial data instances
2. Analyze inter-source relationships
3. Identify influential market signals
4. Learn prediction patterns from grouped instances

This improves stock market forecasting performance.

Step 6: Deep Learning Model Training

The Deep Learning workflow includes:

1. Input multi-source features
2. Train LSTM/CNN prediction models
3. Apply attention mechanisms

Technologies Used

- Python
- Deep Learning
- Machine Learning
- TensorFlow / PyTorch
- LSTM / CNN Models
- Natural Language Processing (NLP)
- Pandas & NumPy
- Scikit-learn
- Flask / Django
- MySQL / MongoDB

Result :



This image shows the User Login Page of a Stock Market Prediction System based on Multi-Source Multiple Instance Learning, where users enter their credentials to access the application.



This screen shows the Event, Sentiment Analysis, and Vector Generation stage of the Stock Market Prediction system, where news and market data are converted into numerical feature vectors for machine learning.



This image shows the output of a Stock Market Prediction System, where test data is analyzed using a machine learning model to predict future stock trends.

The prediction results indicate whether the stock price is expected to Rise (green) or Decline (red), helping investors make informed decisions.

Conclusion

This work presents an advanced **Adaptive Multi-Modal Speech Transformer Decoder** designed to effectively integrate multiple input modalities—namely audio, visual, and textual data—for enhanced speech recognition and understanding. The proposed system addresses the key limitations of existing approaches, such as poor robustness in noisy environments, computational inefficiency, and rigid fusion strategies.

Through dynamic **cross-modal attention**, **modality reliability estimation**, and **multi-level fusion**, the system demonstrates superior performance in challenging scenarios where single-modality models typically struggle. The adaptive nature of the fusion mechanism ensures the system remains resilient even in the presence of degraded or missing modalities, offering a more reliable and intelligent decoding solution.

Furthermore, the implementation leverages state-of-the-art transformer architectures with optimized resource usage, making it suitable for real-time applications in domains such as virtual assistants, automated transcription services, and surveillance systems.

Overall, the study highlights the importance of context-aware and reliability-driven multi-modal integration, providing clear



evidence that under the right conditions and design, multiple modalities **significantly improve** the accuracy and robustness of speech decoding systems.

References

1. H. Inaguma, S. Dalmia, T. Hori, S. Watanabe, "Multimodal Transformer with Missing Modality Imagination for Simultaneous Translation," *arXiv preprint arXiv:2004.12655*, 2020.
2. Y. Lee, C. Kim, H. Kim, "Audio-Visual Transformer for Robust Speech Recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2795–2809, 2021.
3. A. M. R. Dabre, K. Sudoh, S. Kurohashi, "A Survey of Multimodal Machine Translation: Challenges and Future Directions," *ACM Computing Surveys (CSUR)*, vol. 55, no. 2, pp. 1–38, 2023.
4. S. Afouras, J. S. Chung, A. Zisserman, "Deep Audio-Visual Speech Recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 1, pp. 215–229, 2022.
5. Y. Akbari et al., "VATT: Transformers for Multimodal Self-Supervised Learning from Raw Video, Audio and Text," *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
6. D. Chen, H. Hu, X. Wang, and L. Zhang, "Multimodal Transformer with Multi-View Visual Representation for Video Captioning," *CVPR*, 2021.
7. H. Li, Y. Tao, L. Wang, "Multi-Modal Speech Recognition Using Visual Information for Noise Robustness," *ICASSP*, 2019.
8. A. Vaswani et al., "Attention is All You Need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
9. J. Huang, W. Xie, M. Wu, "Lip Reading with Cross-Attention Transformer," *ICASSP*, 2022.
10. S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, "Hybrid CTC/Attention Architecture for End-to-End Speech Recognition," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1240–1253, 2017.



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

the Department of Computer Applications at QIS College of Engineering & Technology, Ongole an Autonomous college in Prakasam dist. SHE completed his undergraduate degree in BCA (Computers) from ANU. With a keen interest in research and practical learning, she is actively involved in academic projects and technical activities related to his field.

AUTHORS PROFILE



Mr. K. Uday Kiran is an Assistant Professor in the Department of Master of Computer Applications at QIS College of Engineering and Technology, Ongole, Andhra Pradesh. He earned his Master of Computer Applications (MCA) from Bapatla Engineering College, Bapatla. His research interests include Machine Learning, Programming Languages. He is committed to advancing research and fostering innovation while mentoring students to excel in both academic and professional pursuits.



Ms. N. Lalitha Adithya yadav is a postgraduate student pursuing a MCA in