

Fine-Tuning Transformer Models for Financial Sentiment Classification: A Deep Learning Approach

Iman Hamdan and Matt Hashemi

AAI-511 – Neural Networks and Deep Learning

University of San Diego

Dr. Esmacili

December 15, 2024

ABSTRACT

This project presents a comprehensive implementation of transformer-based deep learning models for financial sentiment classification, specifically focusing on fine-tuning FinBERT using SP500 sentiment data from 2015. Our methodology involves converting continuous sentiment scores into categorical labels using quantile-based thresholds, generating synthetic financial text representations, and applying domain-specific fine-tuning techniques. The project achieved 78.6% accuracy with FinBERT, significantly outperforming traditional machine learning baselines including Random Forest (71.4%) and Logistic Regression (64.3%). Results indicate that FinBERT substantially outperforms traditional approaches in capturing financial sentiment patterns, with practical implications for automated financial analysis systems. The work contributes to understanding transformer applications in finance while providing practical insights for implementing deep learning solutions in financial technology applications.

Table of Contents

| | |
|--|------------------------------|
| ABSTRACT | 1 |
| TABLE OF CONTENTS | Error! Bookmark not defined. |
| 1. Introduction..... | 3 |
| 2. Literature Review and Background | 4 |
| 2.1 Evolution of Sentiment Analysis in Finance | 4 |
| 2.2 Transformer Architecture and BERT | 5 |
| 2.3 Domain-Specific Language Models: FinBERT | 6 |
| 2.4 Financial Applications and Challenges | 6 |
| 3. Methodology | 7 |
| 3.1 Dataset Description | 7 |
| 3.2 Data Preprocessing and Label Creation..... | 7 |
| 3.3 Synthetic Text Generation | 8 |
| 3.4 FinBERT Implementation | 9 |
| 3.5 Baseline Models and Evaluation..... | 9 |
| 4. Results and Analysis | 10 |
| 4.1 Model Performance Comparison | 10 |
| 4.2 Feature Analysis and Model Insights..... | 10 |
| 4.3 Temporal and Correlation Analysis | 11 |
| 5. Discussion and Limitations | 12 |
| 5.1 Implications and Practical Applications | 12 |
| 5.2 Limitations and Constraints | 13 |
| 5.3 Future Research Directions | 13 |
| 6. Conclusion..... | 14 |
| 6.1 Key Achievements | 14 |
| 6.2 Practical Impact..... | 15 |
| 6.3 Course Relevance and Technical Contributions | 15 |
| 6.4 Final Reflections | 16 |
| References | 17 |

1. Introduction

The intersection of artificial intelligence and financial markets has become increasingly significant as financial information volume continues to grow exponentially. Traditional approaches to financial sentiment analysis have relied on rule-based systems and conventional machine learning techniques, which struggle to capture nuanced language patterns in financial communications (Liu & Zhang, 2019). The emergence of transformer-based architectures, particularly BERT and its domain-specific variants, has revolutionized natural language processing tasks across industries, including finance (Devlin et al., 2018).

Financial sentiment analysis presents unique challenges distinguished from general text classification. Financial language is characterized by specialized terminology, complex syntactic structures, and subtle semantic relationships that significantly impact market interpretations (Malo et al., 2014). These challenges motivated the development of domain-specific language models such as FinBERT, pre-trained on financial corpora to understand financial communication intricacies (Araci, 2019).

This project addresses how transformer-based deep learning models, specifically FinBERT, can be effectively applied to classify financial sentiment using real-world SP500 data. The research objectives include: (1) implementing a robust preprocessing pipeline converting continuous sentiment scores into categorical labels, (2) developing synthetic text generation techniques creating meaningful financial narratives from numerical data, (3) fine-tuning FinBERT models on domain-specific financial sentiment data, (4) comparing transformer-based approaches with

traditional machine learning methods, and (5) analyzing temporal patterns and correlations in financial sentiment data.

The significance extends beyond academic exploration, as financial institutions increasingly rely on automated sentiment analysis systems for risk management, algorithmic trading, and investment decision-making. The project demonstrates practical applications of neural networks and deep learning techniques in real-world financial scenarios, aligning with AAI-511 course objectives while contributing to understanding how pre-trained language models can be adapted for specialized domains.

2. Literature Review and Background

2.1 Evolution of Sentiment Analysis in Finance

Financial sentiment analysis has evolved significantly over two decades, transitioning from simple keyword-based approaches to sophisticated deep learning architectures. Early systems relied on lexicon-based methods using predefined dictionaries of financial terms (Loughran & McDonald, 2011). While computationally efficient, these approaches suffered from limited contextual understanding and inability to capture financial language nuances.

Machine learning techniques marked significant advancement, with Support Vector Machines, Naive Bayes classifiers, and ensemble methods demonstrating improved performance over rule-based systems (Pang et al., 2002). However, traditional approaches faced limitations in handling complex linguistic structures and domain-specific terminology. The feature engineering process

required extensive domain expertise and manual effort to extract meaningful representations from raw text data.

Deep learning architectures, particularly RNNs and LSTM networks, addressed many traditional limitations by automatically learning hierarchical feature representations from raw text (Kim, 2014). These models demonstrated superior performance in capturing sequential dependencies and contextual relationships within financial texts, though computational requirements and training complexity presented practical challenges for large-scale applications.

2.2 Transformer Architecture and BERT

The Transformer architecture revolutionized natural language processing by introducing the attention mechanism as a replacement for recurrent architectures (Vaswani et al., 2017). Self-attention enables models to capture long-range dependencies and parallel processing capabilities, significantly improving performance and computational efficiency.

BERT built upon Transformer architecture by introducing bidirectional training and masked language modeling objectives (Devlin et al., 2018). Unlike previous models processing text unidirectionally, BERT considers entire word context by examining both directions simultaneously, enabling deeper understanding of language context and semantic relationships.

The pre-training and fine-tuning paradigm introduced by BERT became the standard approach for transfer learning in natural language processing. Models are first pre-trained on large-scale unlabeled text corpora using self-supervised learning objectives, then fine-tuned on specific downstream tasks with limited labeled data, leveraging pre-training knowledge to achieve superior specialized task performance (Rogers et al., 2020).

2.3 Domain-Specific Language Models: FinBERT

BERT's success motivated researchers to develop domain-specific variants for specialized applications. FinBERT represents one of the most successful BERT adaptations for financial applications, pre-trained on large financial text corpora including news articles, earnings reports, and analyst communications (Araci, 2019).

Domain-specific pre-training involves continued training of base BERT models on financial corpora, allowing adaptation to unique financial language characteristics. This process results in improved performance on financial NLP tasks, including sentiment analysis, named entity recognition, and document classification (Yang et al., 2020). Studies demonstrate FinBERT consistently outperforms general-purpose BERT models on financial sentiment analysis tasks, highlighting domain-specific adaptation importance.

2.4 Financial Applications and Challenges

Financial sentiment analysis has found numerous practical applications, from algorithmic trading to risk management and investment research. Algorithmic trading systems incorporate sentiment signals as decision-making features, using sentiment scores to identify market trends and predict price movements (Zhang & Skiena, 2010). Risk management applications utilize sentiment analysis to monitor market sentiment and identify potential systematic risk sources (Tetlock, 2007).

Despite significant advances, challenges remain. Financial language's dynamic nature presents ongoing adaptation challenges, requiring regular model updates and retraining to maintain performance (Xing et al., 2018). Data quality and labeling consistency represent additional

challenges, as sentiment interpretation subjectivity can lead to training data inconsistencies (Mohammad, 2016). Regulatory considerations also impact deployment, requiring model interpretability and explainability in regulated environments (Rudin, 2019).

3. Methodology

3.1 Dataset Description

The primary dataset consists of SP500 sentiment data from 2015, containing 202 daily observations spanning January 1 to October 9, 2015. The dataset includes four primary variables: date timestamps, daily average news sentiment scores, daily average Twitter sentiment scores, and SP500 opening prices. News sentiment scores range from -0.267 to 0.144, Twitter sentiment scores from -0.003 to 0.000, and SP500 opening prices from \$45.98 to \$47.45.

Data quality assessment reveals complete coverage with no missing values, ensuring analysis reliability. Correlation analysis shows moderate positive correlation (0.287) between news sentiment and SP500 prices, minimal correlation (-0.037) between Twitter sentiment and prices, and low correlation (0.080) between news and Twitter sentiment, indicating largely independent information sources.

3.2 Data Preprocessing and Label Creation

The preprocessing pipeline converts continuous sentiment scores into categorical labels suitable for neural network classification. We implemented a quantile-based thresholding approach ensuring balanced class distributions while preserving sentiment intensity ordering. The 33rd and

67th percentiles serve as threshold values: negative sentiment below -0.048, neutral sentiment between -0.048 and 0.025, and positive sentiment above 0.025.

This approach resulted in balanced distribution: 67 negative samples (33.2%), 68 neutral samples (33.7%), and 67 positive samples (33.2%). The quantile-based method ensures sufficient observations for effective model training while creating meaningful sentiment distinctions, avoiding issues with fixed thresholds that might create imbalanced class distributions.

3.3 Synthetic Text Generation

A critical methodology component involves generating synthetic financial text representations from numerical sentiment data to enable FinBERT processing. This addresses the challenge of applying text-based transformer models to datasets containing sentiment scores rather than raw text documents.

The text generation algorithm constructs financial news-style descriptions integrating temporal information, market data, and sentiment characterizations. Each generated text follows the structure: "On [date], Microsoft (MSFT) opened at \$[price], reflecting [sentiment description] market sentiment in financial news and [sentiment description] social media sentiment. Financial analysts noted sentiment patterns with news outlets showing [score] sentiment levels while social media platforms indicated [score] sentiment scores."

The process incorporates domain-specific financial language and terminology, enhancing relevance for FinBERT processing. Conditional logic adds contextual information based on sentiment strength: positive sentiment above 0.1 includes optimistic language, while negative sentiment below -0.1 includes concern-related phrases.

3.4 FinBERT Implementation

The implementation utilizes the pre-trained ProsusAI/finbert model from Hugging Face Transformers library, based on BERT-base architecture with 12 transformer layers, 768 hidden dimensions, 12 attention heads, and approximately 110 million parameters. The model is specifically pre-trained on financial texts, enabling domain-specific language pattern capture.

Fine-tuning configuration employs transformer best practices: learning rate $2e-5$, batch size 8 for training and 16 for evaluation, maximum sequence length 128 tokens, 3 epochs with early stopping, AdamW optimizer with linear learning rate schedule, and 100 warmup steps. A custom PyTorch dataset class handles tokenization, padding, truncation, and label encoding for efficient batch processing.

3.5 Baseline Models and Evaluation

Traditional machine learning baselines include Logistic Regression and Random Forest classifiers operating on engineered features from numerical sentiment data. Features include raw sentiment scores, lagged values, moving averages, and interaction terms, totaling eight engineered features capturing temporal patterns and relationships.

The evaluation framework employs temporal validation respecting chronological data ordering: training set comprises first 80% of observations, validation set 10%, and test set final 20%. This approach prevents data leakage and provides realistic performance assessment. Performance metrics include accuracy, precision, recall, and F1-scores, with weighted averaging ensuring fair evaluation across classes.

4. Results and Analysis

4.1 Model Performance Comparison

Comprehensive evaluation reveals significant performance differences across implemented methods. FinBERT achieved highest performance with 78.6% accuracy and 78.9% F1-score, substantially outperforming traditional baselines. Random Forest achieved 71.4% accuracy with 72.1% F1-score, while Logistic Regression achieved 64.3% accuracy with 64.8% F1-score.

Results demonstrate clear performance hierarchy, with FinBERT providing 7.2 percentage point improvement over Random Forest and 14.3 percentage point improvement over Logistic Regression. Performance improvements represent substantial gains over 33.3% random baseline, with FinBERT achieving 135% improvement over random guessing. Consistent performance across accuracy and F1-score metrics indicates genuine classification performance gains rather than imbalanced class predictions.

| Model | Accuracy | F1-Score | Improvement vs Random |
|---------------------|----------|----------|-----------------------|
| FinBERT | 78.6% | 78.9% | +135% |
| Random Forest | 71.4% | 72.1% | +114% |
| Logistic Regression | 64.3% | 64.8% | +93% |
| Random Baseline | 33.3% | 33.3% | Baseline |

4.2 Feature Analysis and Model Insights

Feature importance analysis for Random Forest reveals news sentiment scores as most predictive features (41.2% importance), followed by 3-day moving average of news sentiment (19.6%), opening prices (14.5%), and news-Twitter sentiment interaction (9.8%). These results validate feature engineering approach and provide insights into information source relative importance.

Confusion matrix analysis shows FinBERT demonstrates strong performance across all sentiment classes with balanced precision and recall. Most misclassifications occur at boundaries between sentiment categories, particularly for observations with sentiment scores near quantile thresholds. This pattern is expected given continuous nature of underlying sentiment scores and threshold boundary arbitrariness.

Training progress analysis reveals FinBERT exhibited consistent improvement across epochs with stable convergence patterns. Training loss decreased steadily without overfitting indicators, achieving optimal performance within allocated 3 epochs. Computational efficiency analysis shows approximately 5 minutes per epoch training time and 50 milliseconds per document inference speed, making the approach practical for real-time applications.

4.3 Temporal and Correlation Analysis

The moderate positive correlation (0.287) between news sentiment and SP500 prices suggests news sentiment contains meaningful market condition information, supporting sentiment analysis validity for financial applications. Minimal correlation (-0.037) between Twitter sentiment and prices suggests social media sentiment was less predictive during this 2015 period, possibly reflecting earlier stage social media adoption in financial markets.

Temporal distribution of errors reveals no significant patterns, suggesting consistent model performance across different time periods within the dataset. This consistency indicates successful learning of generalizable patterns rather than overfitting to specific temporal conditions or market events during training period.

5. Discussion and Limitations

5.1 Implications and Practical Applications

FinBERT's superior performance demonstrates significant value of domain-specific language models for specialized applications. The 78.6% accuracy represents substantial improvement over traditional approaches, validating transformer-based architecture effectiveness for financial text analysis. This performance level suggests FinBERT-based systems could provide reliable sentiment classification for practical financial applications.

The synthetic text generation approach represents important methodological contribution, enabling application of powerful text-based models to traditionally numerical datasets. This technique opens possibilities for applying state-of-the-art NLP models to financial time series data and other numerical datasets in specialized domains.

Practical applications include automated financial news monitoring, investment decision support systems, risk assessment tools, algorithmic trading signal generation, and financial report analysis automation. The 78.6% accuracy level, while not perfect, represents significant improvement over manual analysis or rule-based systems, making it suitable for integration with existing financial technology infrastructure.

5.2 Limitations and Constraints

The primary limitation is restricted temporal scope and sample size, with only 202 observations from 10-month period in 2015. This may not capture full range of market conditions and sentiment patterns across different economic cycles. The Microsoft-specific focus limits generalizability to other companies and market sectors.

The synthetic text generation approach, while innovative, may not fully capture complexity and nuance of authentic financial communications. Structured templates create consistent but potentially artificial linguistic patterns that may not reflect real financial writing variability and sophistication.

The quantile-based labeling approach introduces arbitrary threshold effects that may not align with natural sentiment boundaries. Observations near threshold boundaries may be misclassified due to inherent sentiment categorization subjectivity rather than true model limitations.

Computational requirements of FinBERT fine-tuning and inference may limit scalability for real-time applications processing large text volumes. Limited interpretability of transformer-based models presents challenges for regulatory compliance and risk management applications in financial institutions.

5.3 Future Research Directions

Several promising directions emerge from this research. Multi-company and multi-sector analysis would provide insights into approach generalizability across different market segments.

Incorporating longer time series spanning multiple economic cycles would enable model stability assessment across varying market conditions.

Real-time deployment studies would provide valuable insights into practical implementation challenges in production environments. Enhanced feature engineering approaches could combine traditional machine learning strengths with transformer-based methods. Hybrid models incorporating both numerical features and text-based representations could potentially achieve superior performance while maintaining interpretability.

Alternative labeling strategies beyond quantile-based thresholds could provide more meaningful sentiment categories. Market impact-based labeling, where sentiment categories are defined based on subsequent price movements, could create more economically relevant classifications.

6. Conclusion

This research successfully demonstrates FinBERT application for financial sentiment classification using real-world SP500 data. Through comprehensive analysis and rigorous methodology, the study achieved important outcomes contributing to both academic understanding and practical applications in financial technology.

6.1 Key Achievements

The primary achievement is demonstrating FinBERT significantly outperforms traditional machine learning approaches for financial sentiment classification. With 78.6% accuracy, FinBERT achieved 7.2 percentage point improvement over Random Forest and 14.3 percentage

point improvement over Logistic Regression, representing 135% improvement over random baseline classification.

The innovative synthetic text generation methodology represents significant methodological contribution, successfully enabling text-based transformer model application to traditionally numerical datasets. The quantile-based labeling approach provides robust method for converting continuous sentiment scores into balanced categorical labels suitable for neural network training.

6.2 Practical Impact

The research demonstrates clear practical value for financial technology applications.

Performance levels achieved by FinBERT are sufficient for deployment in automated sentiment monitoring systems, investment decision support tools, and risk management applications. The correlation analysis revealing moderate positive relationship between news sentiment and SP500 prices validates sentiment analysis use for financial applications.

6.3 Course Relevance and Technical Contributions

This project effectively demonstrates core AAI-511 Neural Networks and Deep Learning concepts through practical transformer architecture implementation for domain-specific applications. The successful FinBERT fine-tuning showcases transfer learning principles, where pre-trained models are adapted for specialized tasks with limited labeled data.

The implementation demonstrates practical neural network training aspects, including hyperparameter selection, optimization strategies, and overfitting prevention. The use of modern

deep learning frameworks provides experience with current industry-standard tools and practices.

6.4 Final Reflections

This project successfully bridges the gap between academic deep learning research and practical financial applications, demonstrating how state-of-the-art neural network techniques can solve real-world business problems. The comprehensive methodology, rigorous evaluation, and practical considerations make this work valuable contribution to both academic understanding of transformer applications in finance and practical AI system deployment in financial technology.

The superior performance achieved through FinBERT fine-tuning validates domain-specific language model importance and transfer learning in specialized applications. The innovative approach to converting numerical sentiment data into text representations opens new possibilities for applying text-based AI models to traditionally numerical financial datasets.

As financial markets continue generating vast amounts of unstructured data, the techniques and frameworks developed in this project provide essential tools for extracting actionable insights and supporting automated decision-making systems. The work represents significant step forward in neural networks and deep learning application to financial sentiment analysis, with clear pathways for future research and practical implementation.

References

- Araci, D. (2019). FinBERT: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Loughran, T., & McDonald, B. (2011). When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. *The Journal of Finance*, 66(1), 35-65.
- Malo, P., Sinha, A., Korhonen, P., Wallenius, J., & Takala, P. (2014). Good debt or bad debt: Detecting semantic orientations in economic texts. *Journal of the Association for Information Science and Technology*, 65(4), 782-796.
- Mohammad, S. M. (2016). Sentiment analysis: Detecting valence, emotions, and other affectual states from text. *Emotion Measurement*, 201-237.
- Pang, B., Lee, L., & Vaithyanathan, S. (2002). Thumbs up? Sentiment classification using machine learning techniques. *Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing*, 79-86.
- Tetlock, P. C. (2007). Giving content to investor sentiment: The role of media in the stock market. *The Journal of Finance*, 62(3), 1139-1168.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.