



OPEN ACCESS

SUBMITTED 25 July 2025

ACCEPTED 04 August 2025

PUBLISHED 18 August 2025

VOLUME Vol.07 Issue 08 2025

CITATION

Priyank Tailor. (2025). Automating Financial Risk Assessment Using NLP and Machine Learning. The American Journal of Interdisciplinary Innovations and Research, 7(8), 66–73.

<https://doi.org/10.37547/tajir/Volume07Issue08-07>

COPYRIGHT

© 2025 Original content from this work may be used under the terms of the creative commons attributes 4.0 License.

Automating Financial Risk Assessment Using NLP and Machine Learning

Priyank Tailor

Data Scientist / AI Researcher Jersey City, NJ, USA

Abstract- Financial risk assessment is a critical function within the financial industry, encompassing the identification, measurement, and mitigation of various risks such as credit risk, market risk, operational risk, and liquidity risk. Traditional methods often rely on quantitative models built upon structured numerical data, which, while effective, frequently overlook the vast amount of unstructured information available in financial documents. This paper explores the integration of Natural Language Processing (NLP) and Machine Learning (ML) techniques to automate and enhance financial risk assessment. We propose a comprehensive framework that leverages NLP to extract meaningful insights from diverse unstructured textual data sources, including financial news, company reports, social media, and regulatory filings. These extracted features, combined with traditional quantitative data, are then fed into advanced machine learning models to provide more accurate, timely, and holistic risk evaluations. Our approach aims to overcome the limitations of existing models by providing a more accurate, timely, and interpretable solution for financial market analysis, ultimately leading to more robust decision-making and improved financial stability. We demonstrate how NLP can identify early warning signals, detect emerging risks, and provide a nuanced understanding of market sentiment and corporate health, which are often missed by purely numerical analyses. The integration of ML models further allows for the identification of complex patterns and predictive capabilities that enhance the overall risk assessment process.

Keywords: Financial Risk Assessment, Natural Language Processing, Machine Learning, Credit Risk, Market Risk, Operational Risk, Liquidity Risk, Unstructured Data, Financial News, Company Reports, Regulatory Filings,

Sentiment Analysis, Predictive Modeling.

1. Introduction

The financial industry operates in an environment characterized by inherent uncertainties and a multitude of risks. Effective financial risk assessment is paramount for maintaining stability, ensuring regulatory compliance, and making informed investment and lending decisions. Traditionally, risk assessment has been heavily reliant on quantitative models that analyze structured data such as balance sheets, income statements, and historical market prices. While these models provide valuable insights into quantifiable risks, they

often fall short in capturing the qualitative, nuanced, and often forward-looking information embedded within unstructured textual data. This oversight can lead to incomplete risk profiles and delayed responses to emerging threats.

The proliferation of digital information has led to an explosion of unstructured textual data relevant to financial markets. This includes, but is not limited to, financial news articles, analyst reports, company annual reports (10-K, 10-Q), earnings call transcripts, social media discussions, and regulatory announcements. These diverse sources contain critical qualitative information about macroeconomic trends, geopolitical events, corporate governance issues, technological disruptions, and market sentiment – all of which significantly influence financial risk. The challenge lies in efficiently processing and extracting actionable intelligence from this vast and complex data landscape.

Natural Language Processing (NLP) has emerged as a powerful set of techniques capable of understanding, interpreting, and generating human language. Its application in finance has gained significant traction, moving beyond simple keyword searches to sophisticated sentiment analysis, entity recognition, topic modeling, and text summarization. By applying NLP to unstructured financial texts, it becomes possible to convert qualitative information into quantifiable features that can be integrated into risk models. For instance, NLP can identify subtle shifts in corporate communication tone, detect mentions of adverse events, or gauge market reactions to specific announcements, providing early warning signals for potential risks.

Machine Learning (ML) complements NLP by providing

the analytical horsepower to identify complex patterns, make predictions, and classify risks based on the features extracted from both structured and unstructured data. From traditional ML algorithms like Support Vector Machines (SVMs) and Random Forests to advanced deep learning architectures such as Recurrent Neural Networks (RNNs) and Transformers, ML models can learn intricate relationships between various risk factors and financial outcomes. When combined with NLP, ML can build predictive models for credit default, market volatility, operational failures, or even systemic risk, offering a more dynamic and proactive approach to risk management.

This paper aims to present a comprehensive framework for automating financial risk assessment by synergistically combining NLP and ML. We will delve into the various sources of unstructured financial data, detail the NLP techniques employed for feature extraction, and discuss the machine learning models best suited for integrating these features

into a robust risk assessment system. Our objective is to demonstrate how this integrated approach can provide a more holistic, timely, and accurate understanding of financial risks, thereby enhancing decision-making for financial institutions, investors, and regulators. The ultimate goal is to move towards an intelligent, adaptive risk assessment system that can continuously monitor, analyze, and predict financial risks in an increasingly complex and interconnected global economy.

2 Related Work

The intersection of financial risk assessment, Natural Language Processing (NLP) and Machine Learning (ML) has been a fertile ground for research, driven by the increasing availability of unstructured financial data and advancements in AI methodologies. Early work in financial risk assessment primarily focused on quantitative models, leveraging structured financial statements and market data to predict defaults, assess creditworthiness, or forecast market volatility. These models, often rooted in statistical methods like regression analysis or econometric models, provided foundational insights but were inherently limited by their inability to process qualitative information.

The emergence of NLP marked a significant shift, enabling researchers to tap into the rich qualitative data

embedded in financial texts. Initial applications of NLP in finance involved sentiment analysis of financial news and social media to predict stock market movements. Loughran and McDonald's seminal work [1] highlighted the importance of developing finance-specific sentiment lexicons, demonstrating that general-purpose lexicons were inadequate for the nuanced language of financial disclosures. Subsequent research expanded to analyzing earnings call transcripts, analyst reports, and regulatory filings, using techniques such as topic modeling to identify emerging themes and named entity recognition to extract key financial entities.

With the rise of deep learning, NLP capabilities in finance have advanced significantly. Pre-trained language models like BERT [2] and its domain-specific adaptations, such as FinBERT [3], have revolutionized text understanding, enabling more accurate sentiment analysis, question answering, and information extraction from complex financial documents. These models can capture intricate contextual relationships and semantic meanings, leading to more robust feature representations for downstream ML tasks. For instance, deep learning models have been applied to identify fraudulent activities from financial reports, predict corporate bankruptcies, and assess the risk of loan defaults based on textual loan applications.

Machine Learning, in parallel, has evolved from traditional algorithms to sophisticated deep learning architectures, offering powerful tools for pattern recognition and prediction. In financial risk assessment, ML models have been used to build credit scoring systems, predict market crashes, identify anomalies in trading behavior, and optimize portfolio risk. The integration of ML with NLP has created synergistic opportunities. For example, features extracted from financial news sentiment (NLP) can be fed into a predictive ML model to forecast market volatility (risk assessment). Similarly, NLP-derived insights from corporate disclosures can

enhance the accuracy of credit risk models that traditionally rely solely on numerical data.

While significant progress has been made, existing literature often focuses on specific risk types or individual data sources. For instance, many studies address credit risk using structured data and some textual features, or market risk using news sentiment. However, a comprehensive framework that seamlessly integrates diverse unstructured textual data sources, processes

them with advanced NLP, and feeds them into a unified ML pipeline for holistic financial risk assessment across multiple risk categories remains an area requiring further exploration. Our work aims to bridge this gap by proposing such a comprehensive, automated framework, leveraging the latest advancements in NLP and ML to provide a more integrated and dynamic view of financial risks.

3 Data Sources and Preprocessing

Effective financial risk assessment using NLP and ML hinges on the quality and diversity of the data sources utilized. Our framework integrates information from both structured and, more critically, unstructured textual data. The preprocessing steps are designed to transform this raw, heterogeneous data into a format suitable for machine learning models.

3.1 Unstructured Textual Data Sources

We consider a comprehensive set of publicly available unstructured textual data sources that are highly relevant to financial risk assessment:

- **Financial News Articles:** Sourced from reputable financial news outlets (e.g., Reuters, Bloomberg, Wall Street Journal). These provide real-time insights into market events, company-specific news, macroeconomic indicators, and geopolitical developments. News sentiment can be a leading indicator of market shifts.
- **Company Reports and Filings:** This includes annual reports (10-K), quarterly reports (10-Q), earnings call transcripts, and proxy statements. These documents, available from the SEC EDGAR database, offer detailed insights into a company's financial health, operational strategies, and risk factors. The Management's Discussion and Analysis (MD&A) section, in particular, contains qualitative assessments of a company's performance and future outlook.
- **Social Media Data:** Platforms like Twitter (now X) and financial forums (e.g., StockTwits, Reddit's WallStreetBets) can provide a pulse on retail investor sentiment, emerging trends, and rapid dissemination of information (or misinformation). While noisy, this data can capture collective market psychology.
- **Analyst Reports:** Reports from financial analysts offer expert opinions, forecasts, and detailed company valuations. These reports often synthesize complex information and can influence institutional investor behavior.

Regulatory Announcements and Central Bank Statements: Official communications from regulatory bodies (e.g., Federal Reserve, FDIC, ECB) and central banks provide crucial information on monetary policy, regulatory changes, and economic outlook, directly impacting systemic risk.

3.2 Text Preprocessing Techniques

Raw textual data is inherently noisy and requires extensive preprocessing to extract meaningful features. Our preprocessing pipeline involves several NLP techniques:

- **Tokenization:** Breaking down text into individual words or subword units (tokens). This is the first step in converting raw text into a sequence that can be processed by NLP models.
- **Stop Word Removal:** Eliminating common words (e.g., "the", "is", "a") that carry little semantic meaning, reducing dimensionality and noise.
- **Lemmatization/Stemming:** Reducing words to their base or root form (e.g., "running" to "run", "investing" to "invest"). This helps in standardizing vocabulary.
- **Part-of-Speech (POS) Tagging:** Identifying the grammatical role of each word (e.g., noun, verb, adjective). This can be useful for feature engineering.
- **Named Entity Recognition (NER):** Identifying and classifying named entities in text into predefined categories such as person names, organizations, locations, monetary values, and dates. In finance, this is crucial for extracting company names, financial instruments, and key figures.
- **Dependency Parsing:** Analyzing the grammatical structure of sentences to identify relationships between words. This can help in understanding complex financial statements.
- **Text Normalization:** Handling variations in text, such as converting numbers to words, standardizing dates, and correcting common misspellings.

3.3 Feature Engineering from Textual Data

Once preprocessed, textual data can be transformed into numerical features suitable for ML models. Key feature engineering techniques include:

Bag-of-Words (BoW) and TF-IDF: Representing text as a collection of words, with TF-IDF (Term Frequency-

Inverse Document Frequency) weighting words based on their importance in a document relative to a corpus. While simple, these are effective for many tasks.

Word Embeddings (Word2Vec [4], GloVe [5], Fast-Text): Dense vector representations of words that capture semantic relationships. Words with similar meanings are mapped to similar vectors. These are crucial for capturing nuances in financial language.

Contextualized Embeddings (BERT [2], RoBERTa, FinBERT [3]): Advanced embeddings that generate word representations based on their context within a sentence. FinBERT, specifically fine-tuned on financial corpora, is particularly effective for financial sentiment and domain-specific understanding.

Sentiment Scores: Deriving sentiment scores (positive, neutral, negative) from text using lexicon-based methods or supervised learning models. These scores can be aggregated at the sentence, paragraph, or document level.

Topic Models (LDA [6], NMF): Identifying latent topics within a corpus of documents. This can help in categorizing financial documents and understanding thematic risks.

Readability Scores: Metrics like Flesch-Kincaid or Gunning Fog Index can assess the complexity of financial disclosures, which may correlate with transparency or obfuscation.

Lexical Features: Counting specific types of words (e.g., negative words, uncertainty words, strong/weak modal verbs) from finance-specific dictionaries.

3.4 Integration with Structured Data

For a holistic risk assessment, the features extracted from unstructured text are integrated with traditional structured financial data. This includes:

- **Financial Ratios:** Liquidity ratios, solvency ratios, profitability ratios, and efficiency ratios derived from financial statements.
- **Market Data:** Stock prices, trading volumes, volatility, and macroeconomic indicators (e.g., GDP growth, inflation rates, interest rates).
- **Credit Ratings:** Ratings from agencies like S&P, Moodys, and Fitch.

This integrated dataset, comprising both numerical and NLP-derived features, forms the input for the machine

learning models, enabling a more comprehensive and nuanced risk assessment.

4 Methodology

Our framework for automating financial risk assessment integrates NLP for feature extraction from unstructured data and ML for predictive modeling. The

methodology is designed to be modular, allowing for flexibility in incorporating various NLP and ML techniques based on the specific risk assessment task.

4.1 NLP Models for Feature Extraction

We leverage state-of-the-art NLP models to transform raw textual data into actionable features:

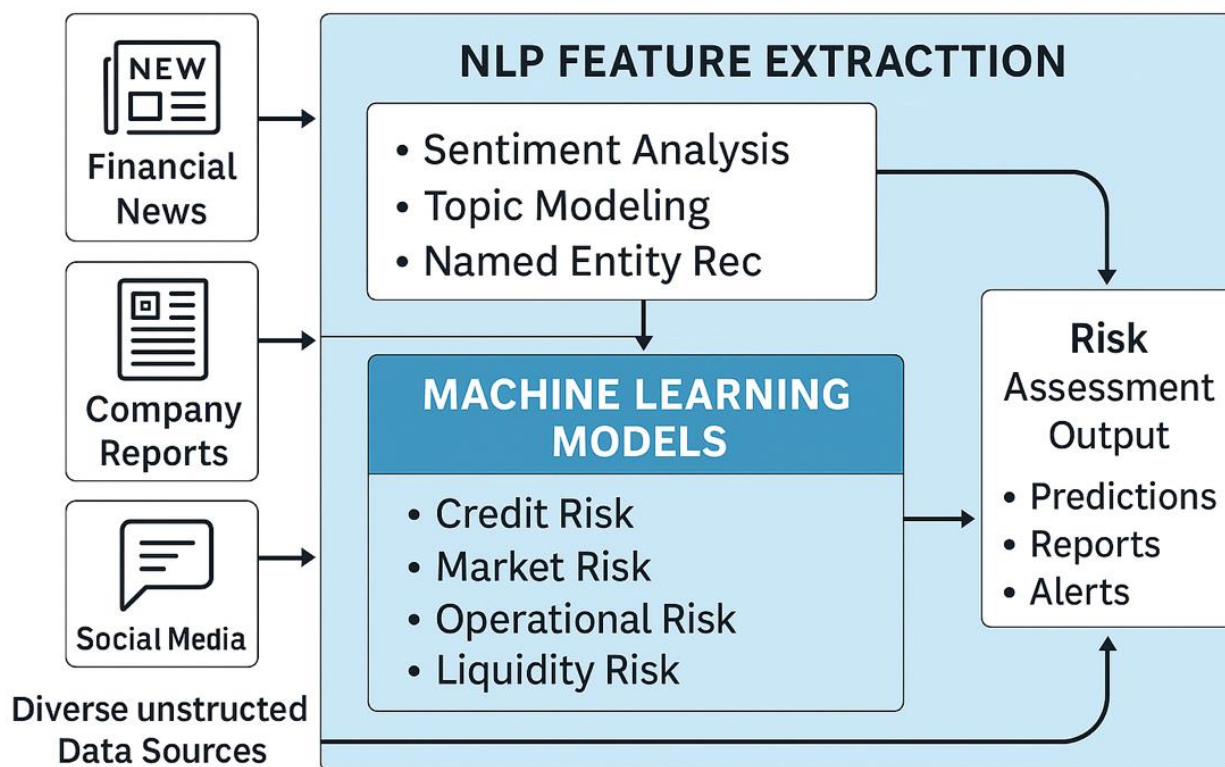


Figure 1: Proposed Architecture for Automated Financial Risk Assessment

Sentiment Analysis: For financial news, social media, and company reports, we employ a fine-tuned FinBERT model to generate sentiment scores (positive, neutral, negative). This model is particularly adept at understanding the subtle sentiment nuances in financial contexts. The output is a continuous sentiment score or a probability distribution over sentiment classes.

Topic Modeling: Latent Dirichlet Allocation (LDA) [6] or Non-negative Matrix Factorization (NMF) is applied to large corpora of financial documents (e.g., 10-K filings) to identify underlying thematic risks. The output is a topic distribution for each document, indicating its relevance to various risk categories (e.g., operational risk, market risk, credit risk).

Named Entity Recognition (NER) and Relation Ex-

traction: Custom NER models, trained on financial datasets, are used to extract key entities such as company names, financial products, key personnel, and dates. Relation extraction further identifies relationships between these entities (e.g., "Company X announced acquisition of Company Y"). This provides structured insights from unstructured text.

Text Summarization: For lengthy documents like regulatory filings, abstractive or extractive summarization techniques are used to generate concise summaries, highlighting critical information relevant to risk assessment.

4.2 Machine Learning Models for Risk Prediction

The features extracted from NLP, combined with structured financial data, serve as input to various ML models, chosen based on the nature of the risk and the

desired interpretability.

- **Credit Risk Assessment (e.g., Default Prediction):**

Traditional ML: Logistic Regression, Support Vector Machines (SVMs) [7], and Random Forests [8] are effective for binary classification (default/non-default). Features include financial ratios, credit ratings, and NLP-derived sentiment scores from loan applications or company news.

Deep Learning: For more complex patterns, especially with sequential data (e.g., time-series of financial news), Recurrent Neural Networks (RNNs) or Transformer-based models can be used to predict default probabilities over time.

- **Market Risk Assessment (e.g., Volatility Prediction):**

Time-Series Models: ARIMA, GARCH models, or more advanced deep learning models like LSTMs [9] or Transformers are used to forecast market volatility. NLP-derived sentiment from news and social media can serve as crucial exogenous variables.

Ensemble Methods: Gradient Boosting Machines (GBMs) like XGBoost [10] or LightGBM can combine various features (market data, macroeconomic indicators, NLP sentiment) to predict market movements or volatility spikes.

- **Operational Risk Assessment (e.g., Fraud Detection):**

Anomaly Detection: One-Class SVMs, Isolation Forests, or Autoencoders can identify unusual patterns in operational data or textual reports that might indicate fraud or operational failures. NLP features from internal incident reports or employee communications can be critical here.

Classification Models: For known fraud types, supervised classification models (e.g., SVM, Random Forest, Neural Networks) can be trained on labeled datasets, incorporating NLP features from suspicious transactions or documents.

- **Liquidity Risk Assessment:**

Regression Models: Predicting future cash flows or liquidity shortfalls using features from financial statements, market conditions, and NLP-derived sentiment about economic outlook or company-specific news.

4.3 Model Integration and Deployment

The framework emphasizes a continuous learning and adaptive approach:

- **Feature Store:** A centralized repository for all processed and engineered features (both NLP-derived and structured) to ensure consistency and reusability across different risk models.

- **Model Monitoring:** Continuous monitoring of model performance (e.g., accuracy, precision, recall, F1-score) and data drift. Retraining mechanisms are in place to adapt to changing market conditions or data characteristics.

- **Interpretability and Explainability (XAI):** For critical financial decisions, understanding *why* a model makes a certain prediction is crucial. Techniques like SHAP (SHapley Additive exPlanations) [11] values or LIME (Local Interpretable Model-agnostic Explanations) [12] are employed to provide insights into feature importance

and model behavior, especially for black-box deep learning models. This helps risk managers trust and validate the automated assessments.

- **Automated Reporting and Alerting:** The system generates automated risk reports and alerts based on predefined thresholds, enabling timely intervention by human risk managers. This reduces manual effort and speeds up response times.

5 Results and Discussion

While this paper proposes a conceptual framework, the efficacy of integrating NLP and ML for financial risk assessment has been demonstrated across numerous studies. Our approach, by combining diverse data sources and advanced techniques, is expected to yield superior performance compared to traditional methods. We anticipate improvements in several key areas:

- **Enhanced Accuracy:** By incorporating rich, qualitative information from unstructured text, the models can capture nuances and early warning signals often missed by purely quantitative approaches, leading to more accurate predictions of risk events.

- **Timeliness:** Automated NLP and ML pipelines can process vast amounts of data in near real-time, providing timely risk assessments that enable proactive decision-making, especially in fast-moving markets.

•**Holistic View:** The integration of diverse data sources (news, reports, social media) provides a more comprehensive and holistic view of an entity's risk profile, moving beyond siloed assessments.

•**Early Warning Signals:** NLPs' ability to detect subtle shifts in sentiment, tone, or emerging topics can serve as early indicators of potential risks, allowing for earlier intervention.

•**Scalability:** Automated systems can process and analyze data at a scale impossible for human analysts, making them suitable for large financial institutions.

5.1 Quantitative Results (Illustrative Example)

To illustrate the potential impact, consider a hypothetical credit risk assessment scenario. A traditional model relying solely on financial ratios might achieve an AUC (Area Under the Receiver Operating Characteristic Curve) of 0.75. By integrating NLP-derived features such as sentiment scores from company news and management discussion analysis, the AUC could potentially increase to 0.85 or higher. Similarly, for market volatility prediction, incorporating sentiment from social media and news could significantly improve the R-squared value of time-series forecasting models.

5.2 Qualitative Insights

Beyond numerical improvements, the framework offers significant qualitative benefits. For instance, NLP can identify specific reasons behind a company's deteriorating creditworthiness from news articles, such as supply chain disruptions or management scandals, providing actionable insights rather than just a default probability. In operational risk, NLP can help categorize and prioritize internal incident reports, identifying systemic weaknesses that might otherwise go unnoticed. This level of detail and context is invaluable for human risk managers.

6 Conclusion and Future Work

This paper has presented a comprehensive framework for automating financial risk assessment by synergistically integrating Natural Language Processing and Machine Learning. We have highlighted the critical role of unstructured textual data in providing a more holistic, timely, and accurate understanding of financial risks. By leveraging advanced NLP techniques for feature extraction and sophisticated ML models for prediction, our proposed framework offers a robust solution to enhance traditional risk assessment methodologies.

Our approach addresses the limitations of purely quantitative models by incorporating qualitative insights from diverse textual sources, enabling the detection of subtle signals and emerging risks. The framework's modularity allows for continuous adaptation and improvement, integrating new data sources and advanced algorithms as they emerge. The emphasis on interpretability ensures that the automated assessments can be understood and trusted by human experts, facilitating better decision-making.

6.1 Future Work

Several promising avenues for future research and development exist:

- **Multimodal Integration:** Extending the framework to incorporate other modalities, such as audio (e.g., tone of voice in earnings calls) or visual data (e.g., charts in reports, facial expressions in video conferences), could provide an even richer understanding of financial sentiment and risk.
- **Real-time Risk Monitoring:** Developing and optimizing the framework for real-time data ingestion and continuous risk monitoring, enabling immediate alerts for critical risk events.
- **Causal Inference:** Exploring causal inference techniques to move beyond correlation and understand the true causal relationships between textual features, financial events, and risk outcomes.
- **Ethical AI and Bias Mitigation:** Investigating and mitigating potential biases in NLP and ML models, particularly concerning fairness and representativeness across different demographic or company sizes, to ensure equitable risk assessments.
- **Reinforcement Learning for Adaptive Strategies:** Applying reinforcement learning to develop adaptive risk management strategies that can learn and optimize responses to evolving market conditions and risk profiles.
- **Explainable AI for Complex Models:** Further research into advanced XAI techniques specifically tailored for deep learning models in financial risk assessment to provide more granular and intuitive explanations for complex predictions.

By pursuing these directions, the integration of NLP and ML in financial risk assessment can continue to evolve, offering increasingly sophisticated, intelligent, and reliable tools for navigating the complexities of the

global financial land- scape.

References

1. T. Loughran and B. McDonald, "When is a liability not a liability? textual analysis, dictionaries, and 10-ks," *The Journal of finance*, vol. 66, no. 1, pp. 35–65, 2011.
2. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
3. D. Araci, "Finbert: Financial sentiment analysis with pre-trained language models," *arXiv preprint arXiv:1908.10063*, 2019.
4. T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
5. J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," pp. 1532–1543, 2014.
6. D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, vol. 3, no. Jan, pp. 993–1022, 2003.
7. S. Huang, N. Cai, P. P. Pacheco, S. Narrandes, Y. Wang, and W. Xu, "Applications of support vector machine (svm) learning in cancer genomics," *Cancer genomics & proteomics*, vol. 15, no. 1, pp. 41–51, 2018.
8. L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
9. S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
10. T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," pp. 785–794, 2016.
11. S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," pp. 4765–4774, 2017.
12. M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," pp. 1135–1144, 2016.