

CHAPTER 1

INTRODUCTION AND LITERATURE REVIEW

The financial sector, a key engine of economic growth worldwide, is currently experiencing a digital transformation. The need for accuracy and efficiency in processing financial information has never been more pressing as financial transactions and data quantities grow exponentially.

The manual labelling of financial numerical items seems outdated and error-prone in this era of fast data expansion. Despite being invaluable, the human touch is prone to errors that could have significant effects on the financial world. In light of this, we present "FinSentix: Contextual Enrichment in Financial Script through Sentiment and Term Extraction," a project that has the potential to completely transform the way financial data processing is done.

The goal of FinSentix is to bring in a new era of openness and accessibility inside the financial industry. It accomplishes this by automating and standardising the detection and labelling of numerical items within financial texts by utilising the power of Natural Language Processing (NLP) and advanced machine learning techniques.

1.1 Motivation

In the digital age, the financial landscape is undergoing a seismic shift with the advent of vast amounts of data generated by financial transactions, market activities, and economic indicators. This unprecedented surge in data volume has ushered in a new era of opportunities and challenges.

In this dynamic environment, the importance of timely and accurate data processing cannot be overstated. Financial institutions, analysts, and decision-makers rely on precise and up-to-date information for investment strategies, risk assessment, and compliance reporting. However, amidst the data deluge, there lies a fundamental challenge – the effective identification and categorization of numerical entities within financial texts. Traditional methods of manual tagging and labeling of financial numerical entities are not only labor-intensive but also error-prone. Inconsistencies in data labeling can lead to erroneous analyses, financial losses, and compliance violations.

It is within this context that our project, "FinSentix: Contextual Enrichment in Financial Script through Sentiment and Term Extraction," finds its deepest motivation. We are inspired by the pressing need to address the following critical factors:

- **Data Deluge:** The financial industry is experiencing an unprecedented influx of data. Each day, terabytes of financial information are generated across markets, corporations, and economic indicators.

- **Human Error:** Manual tagging of financial numerical entities is inherently error-prone, leading to inconsistencies, inaccuracies, and misinterpretations in financial data.
- **Time Sensitivity:** In financial decision-making, every second counts. The efficiency of data processing directly impacts the speed at which informed decisions can be made.
- **Compliance Challenges:** Manual tagging vulnerabilities pose regulatory compliance risks, as errors in financial reporting may lead to non-compliance with industry standards and regulations.
- **Cost Reduction:** Efficient data processing can lead to cost savings in terms of labor, data storage, and maintenance.
- **Transparency and Accountability:** Accurate and standardized data labeling is essential for fostering transparency and accountability in financial reporting, both within organizations and for regulatory bodies.
- **Inefficient Information Retrieval:** The absence of standardized tagging impedes efficient data retrieval and hampers comparative analysis across financial reports and documents.
- **Limitations in Analysis:** Manual tagging limitations restrict the depth of financial insights that can be derived from textual data, hindering the capacity to make well-informed investment decisions, assess risks, and uncover market trends.
- **Demands for Innovation:** In a rapidly evolving financial landscape, innovation is not just a choice; it's a necessity. The industry is ripe for transformative solutions that streamline data processing.
- **Global Impact:** In an interconnected global economy, the transparency and accuracy of financial data have far-reaching implications. Our project's potential impact extends beyond individual institutions or markets.
- **Empowering Decision-Makers:** Our motivation is grounded in the desire to empower decision-makers across the financial spectrum, from investment analysts to regulatory bodies.
- **Efficiency as a Priority:** Efficiency is the lifeblood of the financial industry. Every resource saved in data processing can be redirected toward higher-value financial analysis and strategic decision-making.
- **Acknowledging the Challenges:** Our motivation is rooted in the acknowledgment of the challenges, opportunities, and responsibilities that lie at the heart of the financial domain.
- **Competitive Advantage:** In the highly competitive financial industry, timely and accurate data processing can provide a distinct advantage in identifying market trends and investment opportunities.

- **Customer Expectations:** Clients and stakeholders in the financial sector expect quick access to accurate data. Meeting these expectations is essential for customer satisfaction and retention.
- **Risk Mitigation:** Accurate data labeling and analysis are vital for risk assessment and management, helping institutions avoid financial pitfalls and regulatory penalties.

1.2 Key Challenges

- **Data Complexity and Diversity:** Managing and processing diverse financial data sources with varying structures posed a significant challenge during the development of FinSentix. Financial scripts, often intricate and diverse, required robust data preprocessing techniques to ensure the system could handle the complexity inherent in different data sources. The team focused on developing flexible data processing pipelines capable of adapting to the diverse formats and structures encountered in financial documents.
- **Algorithm Selection and Fine-Tuning:** Identifying and implementing suitable algorithms for sentiment analysis and term extraction was a critical challenge. The team conducted extensive research and experimentation to select algorithms that could effectively analyze financial language nuances. Fine-tuning became paramount to optimize the performance of these algorithms, ensuring accurate sentiment analysis and precise term extraction for enriched financial scripts.
- **Contextual Understanding and Ambiguity:** Addressing contextual nuances and potential ambiguity in financial language was a key challenge. FinSentix needed to go beyond basic language processing and comprehend the intricate context within financial scripts. The development team focused on creating sophisticated Natural Language Processing (NLP) models capable of contextual understanding and disambiguation, ensuring the enriched data provided meaningful insights.
- **Integration of NLP and ML Technologies:** Seamlessly integrating advanced NLP and machine learning technologies into the project architecture presented a complex challenge. Collaboration with experts in these fields was essential to ensure a harmonious integration of cutting-edge technologies. Leveraging established frameworks and industry best practices facilitated the successful integration, enabling FinSentix to harness the full potential of NLP and ML in financial data processing.
- **Data Labeling and Training Set Quality:** Ensuring high-quality labeled data for training models, especially in the financial domain, was a pivotal challenge. Rigorous data labeling processes were implemented, and continuous quality assurance measures were adopted to maintain the integrity of the training sets. This meticulous approach was crucial in ensuring the accuracy and reliability of FinSentix's analysis.

- **Scalability and Performance Optimization:** Designing the system to scale with increased data volumes while maintaining optimal performance was a critical consideration. The development team focused on employing scalable infrastructure and optimizing code for efficiency. This approach allowed FinSentix to handle growing data demands without compromising on performance.
- **User Interface and Accessibility:** Designing an intuitive user interface for financial professionals to interact with enriched data presented its own set of challenges. The team conducted extensive user experience testing and gathered feedback to refine the interface. This iterative process ensured that FinSentix not only delivered powerful analytical capabilities but also provided a user-friendly experience for financial professionals.
- **Real-Time Processing Requirements:** Meeting the demand for real-time processing of financial data was a crucial challenge. FinSentix needed to analyze data streams in real-time to provide timely insights. The development team implemented efficient streaming data processing techniques and optimized algorithms to meet the stringent requirements of real-time analysis in the financial domain.

1.3 Problem addressed in thesis

In the modern financial landscape, the volume and complexity of textual financial data have surged to unprecedented levels. This data encompasses a vast array of numerical entities embedded within a sea of financial reports, news articles, and documents. However, the existing methodology for tagging and labelling these numerical entities within financial texts is fundamentally flawed. Manual tagging, the prevailing practice, is labour-intensive, error-prone, and often inconsistent. Financial professionals tasked with this arduous process face substantial challenges, including:

- 1. Time-Consuming Data Entry:** The manual tagging of numerical entities is a time-consuming endeavour, diverting valuable resources from higher-level financial analysis and decision-making tasks.
- 2. Error-Prone Labelling:** Human error is inherent in the manual tagging process, leading to inconsistencies, inaccuracies, and misinterpretations in financial data.
- 3. Inefficient Information Retrieval:** The absence of standardized tagging impedes efficient data retrieval and hampers comparative analysis across financial reports and documents.
- 4. Risk of Compliance Violations:** Manual tagging vulnerabilities pose regulatory compliance risks, as errors in financial reporting may lead to non-compliance with industry standards and regulations.

5. Inhibited Financial Insights: The limitations of manual tagging restrict the depth of financial insights that can be derived from textual data, hindering the capacity to make well-informed investment decisions, assess risks, and uncover market trends.

The problem at hand is clear: the financial industry is grappling with an obsolete and inefficient methodology for processing numerical entities within financial texts, resulting in a bottleneck that impedes transparency, accuracy, and efficiency. The urgent need for a systematic solution that automates and standardizes this process has never been more apparent.

1.4 Approach to the Problem

The strategy is to create "FinSentix" as a revolutionary solution to address the difficulties associated with manual financial data extraction. FinSentix is a machine learning and advanced natural language processing (NLP) system that automates the recognition of various financial numerical entities in intricate texts. Acknowledging the shortcomings of manual procedures, the method concentrates on subtle comprehension that goes beyond simple numerical recognition. FinSentix explores contextual nuances and supports adaptive learning, going beyond simple tagging. This iterative approach redefines financial data comprehension by aiming for both efficiency and accuracy. FinSentix emerges as a visionary force that provides professionals with dynamic financial insights, not only as a tool.

1. **Introduction:** In the fast-paced world of finance, precise and quick data processing is essential to making well-informed decisions. Inefficient and error-prone processes arise when financial data is manually extracted and interpreted from complex reports. This sparked the development of "FinSentix," a ground-breaking programme that makes use of sophisticated machine learning and natural language processing (NLP). FinSentix marks a new era in data comprehension by adding context, sentiment, and pertinent phrases to financial scripts beyond simple numerical labelling. This innovative technology promises unmatched transparency in the financial landscape by serving as a disruptive link between raw data and intelligent judgements.
2. **Literature Review** - It explores the nuances of machine learning applications, sentiment analysis, and current practices in the financial industry. An extensive review of academic literature provides the background information required to make the shift from manual to automated problem solving. Diverse research provide insights that not only improve FinSentix's theoretical framework but also shed light on the opportunities and problems associated with utilising new technology to enhance the comprehension of financial data.
3. **Problem Statement:** The thesis outlines the challenges associated with manual financial data processing in the problem statement. It reveals the shortcomings that FinSentix is

well-positioned to address and throws light on the constraints of conventional methods. The problem statement directs attention to the central purpose of the thesis, acting as a compass.

4. **Methodology:** Through the introduction of the comprehensive framework "FinSentix," this project transforms the handling of financial data. Every stage is painstakingly designed, from thorough preprocessing and varied data gathering to sophisticated NLP model selection and customisation. The processes of training, validation, and labelling guarantee the performance and correctness of the model. Real-time data integration is applied where practical, resolving issues with dynamic analysis. Creating an intuitive user interface that allows for smooth interaction with financial text and entity recognition findings is the last step. This all-inclusive solution represents a revolutionary step, offering unmatched financial landscape insights from the very beginning of data collection to user involvement.
5. **System Architecture:** The FinSentix neural pathways are walked through in detail in the architecture section. It reveals the subtleties of its architecture, highlighting the union of advanced technology with sensible scalability. The section presents an intelligent system that can not only process data but also carefully orchestrate it.
6. **Results and Analysis:** FinSentix takes front stage in the field of Results and Analysis. This section demonstrates its capabilities with verifiable proof of effectiveness, accuracy, and performance. It's proof of FinSentix's capacity to sort through the complexities in financial scripts, not just the numbers.
7. **Discussion:** The thesis's intellectual pinnacle is the discussion, where conclusions are analysed and their ramifications are felt. This thoughtful discussion of FinSentix's achievements and shortcomings leads readers through the system's complexities and inspires reflection on potential directions for financial data processing in the future.
8. **Conclusion:** In the Conclusion, a crescendo that sums up the trip, the curtain falls. It's a celebration of FinSentix as a transforming force rather than merely a synopsis. The goal of the thesis, which is to redefine financial data processing one enhanced script at a time, is left unfinished.

1.5 LITERATURE REVIEW

1) AckNER

- Financial Named Entity Recognition (FNER): The paper introduces AckNER, a tool designed for extracting financial information from research articles and dissertations.
- NLP Models for Financial Text: AckNER employs advanced algorithms like Bidirectional LSTM-CRF for sequence labeling in financial texts.
- Sentiment Analysis in Finance: Not explicitly discussed.
- Key Entity Detection: Not explicitly discussed.
- Weak Supervision in NER: The paper primarily focuses on the development of AckNER for financial NER.
- Custom Labeling Functions: Custom labeling functions are used to enhance NER accuracy.
- Label Aggregation Techniques: Label aggregation techniques are not a central part of this paper.
- Financial Corpora and Datasets: The paper utilizes research articles and dissertations for financial information extraction.
- Performance Metrics: Performance metrics are not explicitly discussed but can be inferred based on the focus on accuracy.
- Challenges in Financial NER: The paper addresses the challenges of financial NER, including the need for specialized tools.

2) A BERT-based Sentiment Analysis and Key Entity Detection Approach for Online Financial Texts

- Financial Named Entity Recognition (FNER): Addresses FNER indirectly through key entity detection in online financial texts.
- NLP Models for Financial Text: Applies BERT and RoBERTa models for sentiment analysis and key entity detection in financial texts.
- Sentiment Analysis in Finance: Focuses on negative sentiment information extraction in online financial texts using RoBERTa and ensemble learning.
- Key Entity Detection: Addresses key entity detection in online financial texts using RoBERTa and ensemble learning.
- Weak Supervision in NER: Weak supervision is not a primary focus of this paper.
- Custom Labeling Functions: Custom labeling functions are not explicitly discussed.
- Label Aggregation Techniques: Label aggregation techniques are not a central part of this paper.
- Financial Corpora and Datasets: Utilizes online financial texts for sentiment and key entity detection.
- Performance Metrics: Evaluates performance using accuracy and F1 score.
- Challenges in Financial NER: Addresses challenges in financial NER, including variations in naming conventions and unconventional expression of financial entities.

3)FinBERT: A Pre-trained Financial Language Representation Model for Financial Text Mining

- Financial Named Entity Recognition (FNER): Primarily focuses on developing FinBERT, a pre-trained model for financial text mining.
- NLP Models for Financial Text: Introduces FinBERT as a domain-specific pre-trained language model for financial text mining.
- Sentiment Analysis in Finance: Demonstrates FinBERT's effectiveness in sentiment analysis tasks on financial data.
- Key Entity Detection: While not the primary focus, FinBERT can be used for key entity detection.
- Weak Supervision in NER: Weak supervision is not a primary focus of this paper.
- Custom Labeling Functions: Custom labeling functions are not explicitly discussed.
- Label Aggregation Techniques: Label aggregation techniques are not a central part of this paper.
- Financial Corpora and Datasets: Utilizes financial text data for pre-training FinBERT.
- Performance Metrics: Evaluates FinBERT's performance in tasks like sentiment analysis and question answering.
- Challenges in Financial NER: Addresses the challenge of limited labeled data in the financial domain.

4)FiNER: A Weakly Supervised Named Entity Recognition Framework for Financial Text

- Financial Named Entity Recognition (FNER): Focuses on improving NER accuracy in financial news articles using FiNER-ORD dataset and the FiNER framework.
- NLP Models for Financial Text: Utilizes FiNER-LFs and Snorkel for weak supervision in NER tasks within financial text.
- Sentiment Analysis in Finance: Sentiment analysis is not the primary focus of this paper.
 - Key Entity Detection: While not the primary focus, FiNER can be used for key entity detection.
- Weak Supervision in NER: The paper introduces FiNER as a weak-supervision framework utilizing custom labeling functions and Snorkel for label aggregation.
- Custom Labeling Functions: Introduces FiNER-LFs, custom labeling functions designed for financial NER.
- Label Aggregation Techniques: Utilizes Snorkel's weighted majority vote aggregation for label aggregation.
- Financial Corpora and Datasets: Uses the FiNER-ORD dataset for specialized financial NER.
- Performance Metrics: Evaluates NER model performance using F1-score and weighted average F1-score.
- Challenges in Financial NER: Addresses challenges in financial NER, including handling overlapping entities and distinguishing LOC and ORG entities.

This overview provides insights into each research paper's contributions and focus areas within the parameters you specified.

Chapter 2

Mathematical Model/Experimentation Methods And Materials

2.1 Overall tools and technologies used

1) Long Short-Term Memory (LSTM): Recurrent neural network (RNN) architecture of this kind was created to get over typical RNNs' shortcomings in managing lengthy sequences and identifying long-term dependencies in sequential input. Here is a more thorough description of LSTM:

Architecture Overview: Memory cells and specialised gates enable Long Short-Term Memory (LSTM) to selectively retain or forget information over different time steps. The cell state, input gate, forget gate, and output gate are among the essential parts.

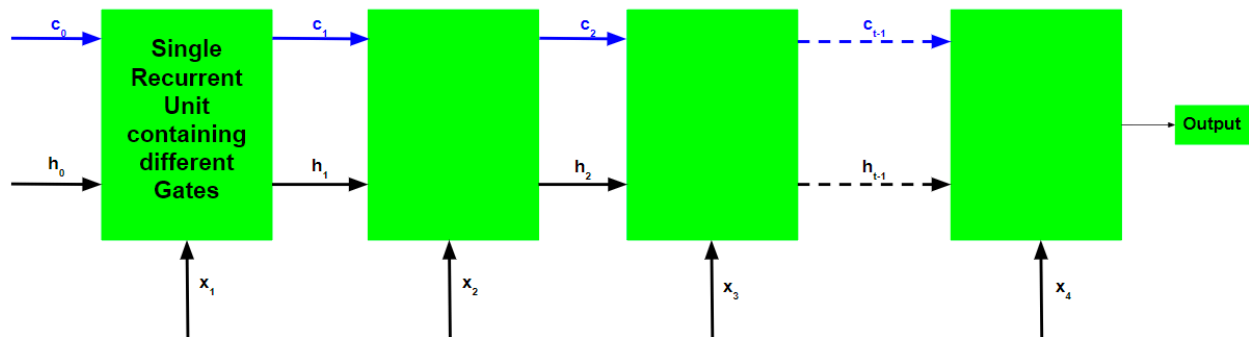
Memory Cells: An LSTM network's memory cells are its key component. Throughout the sequence, these cells act as a conduit of information by maintaining a cell state. This state enables the network to remember long-term dependencies by transferring information across time steps with little modification.

Gates: LSTMs use gates to regulate information flow inside the network:

- 1) Input gate establishes the amount of newly acquired data that needs to be kept in the cell state.
- 2) Forget Gate determines which data in the current cell state should be ignored or forgotten.
- 3) Output gate regulates the amount of the cell state that can be predicted or disclosed to the following layer.

Handling Long-Term Dependencies: LSTMs are capable of learning and remembering patterns in sequential data over extended periods. This ability to capture long-range dependencies makes them suitable for tasks involving time series prediction, speech recognition, machine translation, and sentiment analysis in natural language processing.

The Internal Cell State is likewise transmitted forward together with the Hidden State, which is the only distinction between the fundamental processes of a Long Short Term Memory Network and a Recurrent Neural Network:



The above-stated working is illustrated as below:-

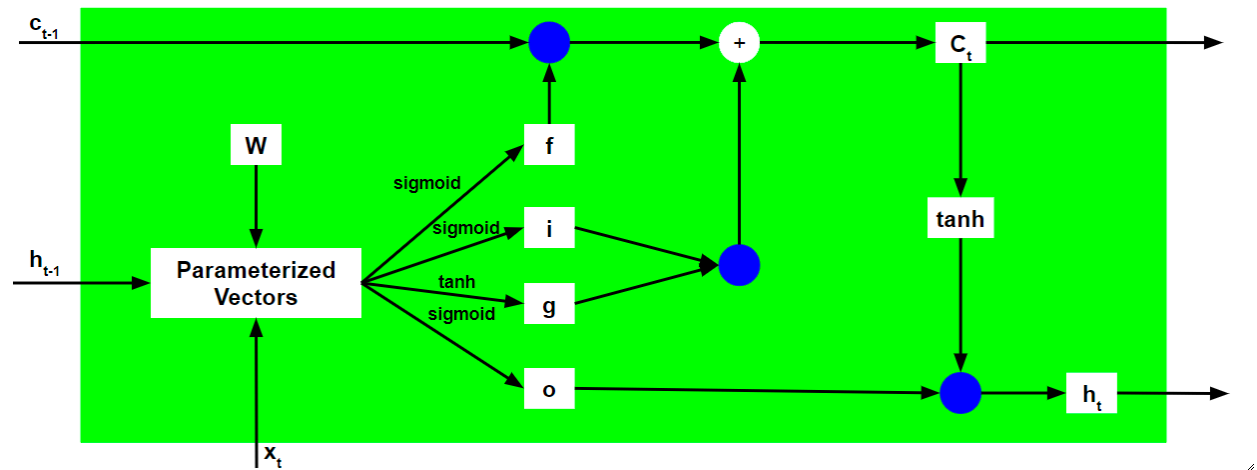


Fig 2.1.2 Working of single recurrent unit of LSTM

2) Bidirectional Representation for Transformers: It is referred to as BERT to enhance the model's linguistic comprehension. BERT is tested and trained on several tasks using distinct architectures. A few of these assignments using the design that is covered here.

Masked Language Model: In this NLP challenge, we substitute the [MASK] token for 15% of the text's words. Next, the programme forecasts the original words that the [MASK] token will replace. In addition to masking, the masking slightly modifies the data to enhance the model's performance during fine-tuning, as the [MASK] token caused a discrepancy between training and fine-tuning. We place a classification layer on top of the encoder input in this model. We additionally compute the output probability by means of a fully connected and a softmax layer.

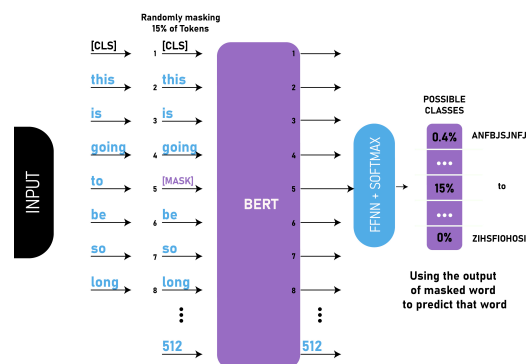


Fig 2.1.3 Masked Language Model

Masked Language Model: The BERT loss function only takes into account the predicted values of the masked values during calculation, disregarding the predicted values of the non-masked values. This facilitates the computation of loss for those 15% of masked words alone.

Next Sentence Prediction: Given two sentences in this natural language processing (NLP) problem, we must determine if the second sentence is the first sentence's follow-up or not. In order to train the BERT, we use 50% of the data to represent the sentence that follows the original sentence (labelled as isNext) and 50% of the time to represent a random sentence that doesn't follow the original text (labelled as NotNext). Given that the task at hand involves classification, the [CLS] token is the first one. To further distinguish the two sentences that we fed into the model, this model furthermore employs a [SEP] token.

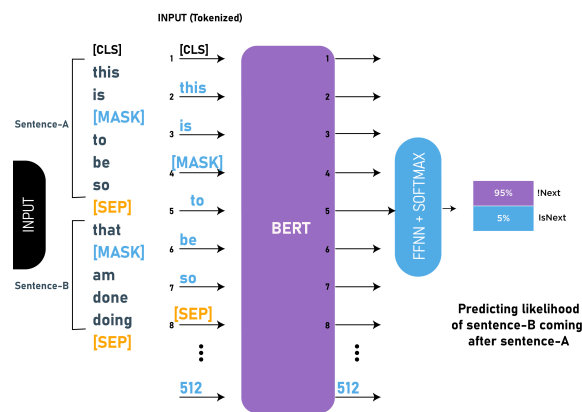


Fig 2.1.4 Next Sentence Prediction

3)Python: This high-level, interpreted programming language has dynamic semantics and is object-oriented. Its dynamic typing and dynamic binding, along with its high-level built-in data structures, make it an appealing language for Rapid Application Development and for usage as a scripting or glue language to join existing components. Because of its straightforward, basic syntax, Python promotes readability, which lowers software maintenance costs. Python's support for packages and modules promotes code reuse and program modularity. The large standard library and the Python interpreter are freely distributable and accessible for free on all major platforms in source or binary form.

2.2 Libraries/Dependencies of the project

- **NumPy (1.19.2):** NumPy is a fundamental package for scientific computing with Python. It provides support for arrays, matrices, mathematical functions, and operations on these objects.
- **Pandas (0.25.3):** Pandas is a powerful library for data manipulation and analysis. It provides data structures like DataFrame and Series, along with tools for reading and writing data in various formats.
- **Matplotlib (3.1.1):** Matplotlib is a widely-used plotting library in Python. It enables the creation of various types of plots, charts, histograms, and visualizations.
- **Keras (2.2.4):** Keras is an easy-to-use high-level neural networks API, capable of running on top of TensorFlow, Theano, or Microsoft Cognitive Toolkit (CNTK). It simplifies the process of building and training neural networks.
- **TensorFlow (1.14.0):** TensorFlow is an open-source machine learning framework developed by Google. It provides tools for building and training machine learning models, particularly neural networks.
- **h5py (2.10.0):** h5py is a Python package that provides an interface to the HDF5 binary data format. It allows for storing and manipulating large datasets efficiently.
- **Protobuf (3.16.0):** Protocol Buffers is a method of serializing structured data. The protobuf library provides tools for working with these serialized data objects.
- **Scikit-learn (0.22.2.post1):** Scikit-learn is a powerful library for machine learning built on NumPy, SciPy, and matplotlib. It offers simple and efficient tools for data mining and data analysis.
- **Seaborn (0.10.1):** Seaborn is a statistical data visualization library based on Matplotlib. It provides a higher-level interface for drawing attractive and informative statistical graphics.
- **PyTorch (1.6.0):** PyTorch is an open-source machine learning library developed by Facebook's AI Research lab. It's known for its ease of use and flexibility in building deep learning models.
- **Transformers (4.16.2):** Transformers is a library by Hugging Face that provides state-of-the-art natural language processing models (NLP), including BERT, GPT, and many others.
- **Huggingface-hub (0.4.0):** This library allows for accessing and sharing models and other resources available on the Hugging Face Hub.
- **Sentence-transformers (2.2.0):** Sentence-transformers is a library for transforming sentences into numerical embeddings using pre-trained transformer models.
- **Sentencepiece (0.1.96):** Sentencepiece is an unsupervised text tokenizer and detokenizer mainly focused on neural network-based approaches, often used in natural language processing tasks.
- **NLTK (3.4.5):** NLTK stands for Natural Language Toolkit. It's a leading platform for building Python programs to work with human language data.

2.2 Algorithm running in the backend

Let's break down the proposed algorithm into its modules and elaborate on each step:

1)Get financial entity list and select key entities: During the second stage, particular entities or features in the financial language are tagged or labelled using a transformer model, such as BERT. BERT recognises and labels various aspects in the text, such as entities, numerical values, and pertinent keywords, by using its contextual awareness. BERT uses its attention mechanism and pre-trained knowledge to tag or label these items. This makes it possible to extract important details from the financial text, such as entities, prices, percentages, and other relevant details. The act of tagging helps to extract structured and useful data, which makes it easier to analyse and make decisions about the financial content later on.

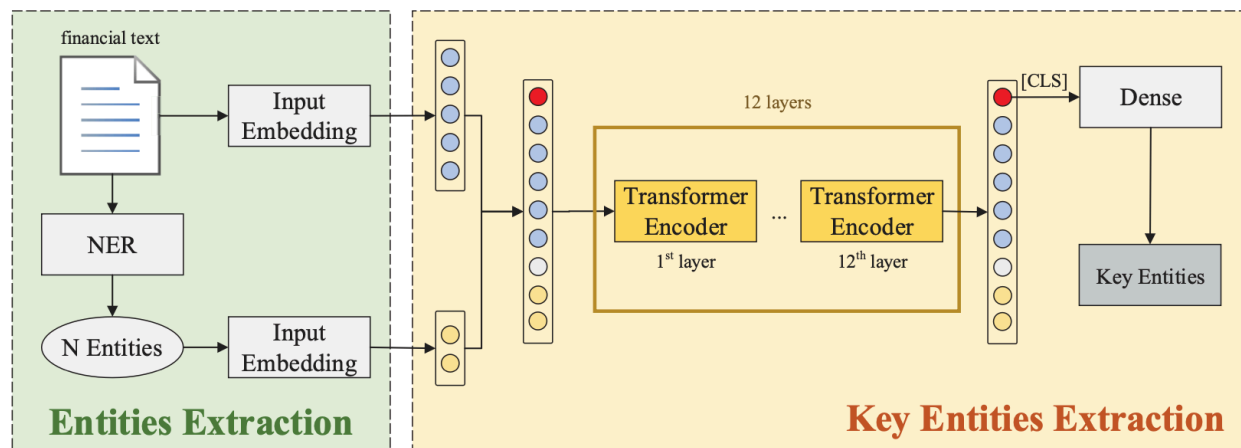


Fig 2.2.1 Get Financial entities and select financial entities

2)Analyze the sentiment of text : Preprocessing financial text data and feeding it into a transformer-based model, is the first step in the process. As the text is processed by the transformer, complex contextual relationships and subtleties within the financial domain are captured. With each word, the transformer creates a contextualised embedding or representation based on its extensive knowledge of language semantics. These embeddings are sent into a later classification layer that predicts sentiment, revealing whether the text conveys favourable, negative, or neutral feelings about financial news or events.

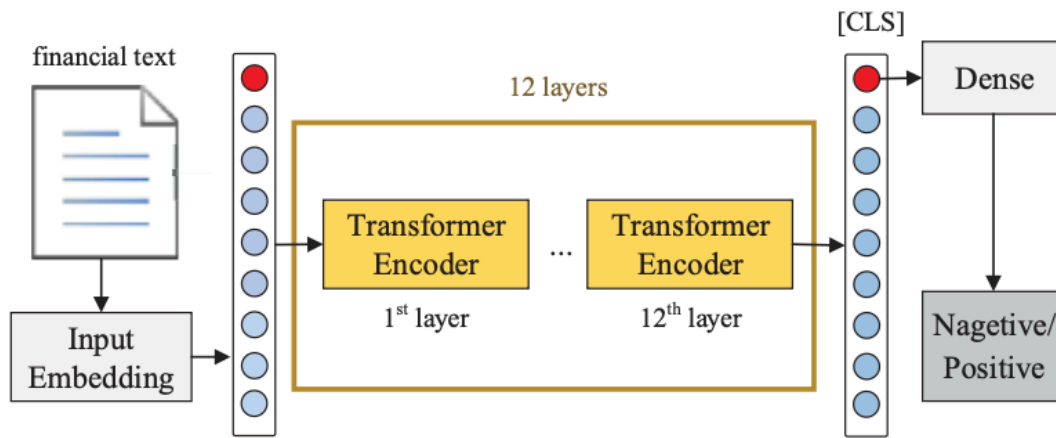


Fig 2.2.2 Analyze the sentiment of the text

3)Select the key entity with the tag: The third stage focuses on extracting important entities from the financial language that are linked to the tags that were obtained from the transformer-based model. These tags may indicate sentiments or other identified elements. Specific entities associated with these feelings or pertinent traits are extracted by utilising the tagged information, regardless of its positivity or negativity. The system finds and retrieves important elements—like firms, financial measures, events, or sentiments—that are most relevant or influential with respect to the designated tags by using natural language processing techniques and the tagged data. This procedure makes it possible to extract pertinent data that is in line with the sentiments or other targeted tags that have been detected, which facilitates improved understanding and analysis of the financial language content.

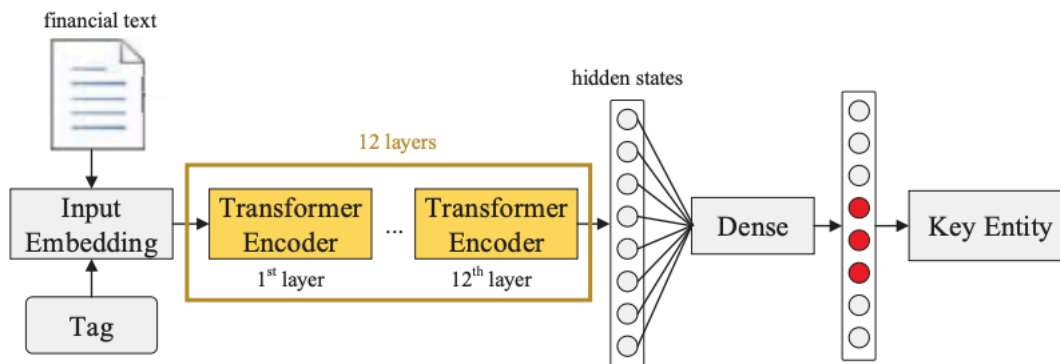


Fig 2.2.3 Select the key entity with the tag

Chapter 3

Results And Discussions

3.1 Technical Workflow

1. Data Collection and Processing:

Data Collection: Financial news datasets are collected from the Hugging Face website, presumably containing text-based news articles or snippets related to finance.

Data Preprocessing: The collected dataset undergoes several preprocessing steps:

1. Removal of stop words: Commonly occurring words (e.g., "and," "the") that do not contribute much to the meaning are removed.
2. Special symbols removal: Punctuation marks, symbols, or characters irrelevant to sentiment analysis are eliminated.
3. Stemming and Lemmatization: Words are reduced to their root form to normalize the text.
4. Business Data Cleaning: Specific cleaning tasks related to the financial domain are performed to refine the dataset further.

2. NLP Model Selection:

- **Numeric Value Identification:** An NLP regular expression algorithm is used to extract and identify numeric values present within the financial news text. These could be currency amounts, percentages, or any numerical data related to finance.
- **BERT Embeddings:** The preprocessed text data is passed through the BERT (Bidirectional Encoder Representations from Transformers) algorithm. BERT generates numeric vectors for each word in the news text, replacing each word with its average frequency in the context of the sentence.

3. Labelling: The financial dataset from Hugging Face's Financial PhraseBank is utilized, as it already contains annotated sentiments (positive, negative, neutral) associated with the financial phrases or sentences. This labeled dataset aids in supervised learning for sentiment analysis.

4. Training and Validation: The BERT-processed dataset is divided into an 80-20 split, where 80% of the data is used for training the model, and the remaining 20% is held out for testing the trained model's performance.

5. LSTM Training: An LSTM (Long Short-Term Memory) neural network is employed using the 80% training data to learn patterns and relationships within the financial news text and its corresponding sentiments.

6. Model Evaluation and Real-Time Prediction:

- **Model Evaluation:** The trained LSTM model is tested on the unseen 20% test data to calculate prediction accuracy. This step validates the model's performance in predicting sentiments from the financial news test data.

- **Real-Time Prediction:** A model is built to predict prices, percentages, and news sentiments from new test data. This model is deployed to predict these details from the uploaded financial news test data in real-time.

This proposed algorithm encompasses data preprocessing, NLP model selection, labeling with sentiment annotations, data splitting for training and testing, LSTM model training, evaluation, and real-time prediction to analyze financial news sentiment and potentially other numerical information present in the news articles. The ultimate goal is to create a predictive model capable of understanding sentiments in financial news and possibly other related numerical information.

3.2 Application Features Demo

1) To run project double click on run.bat file to get below screen

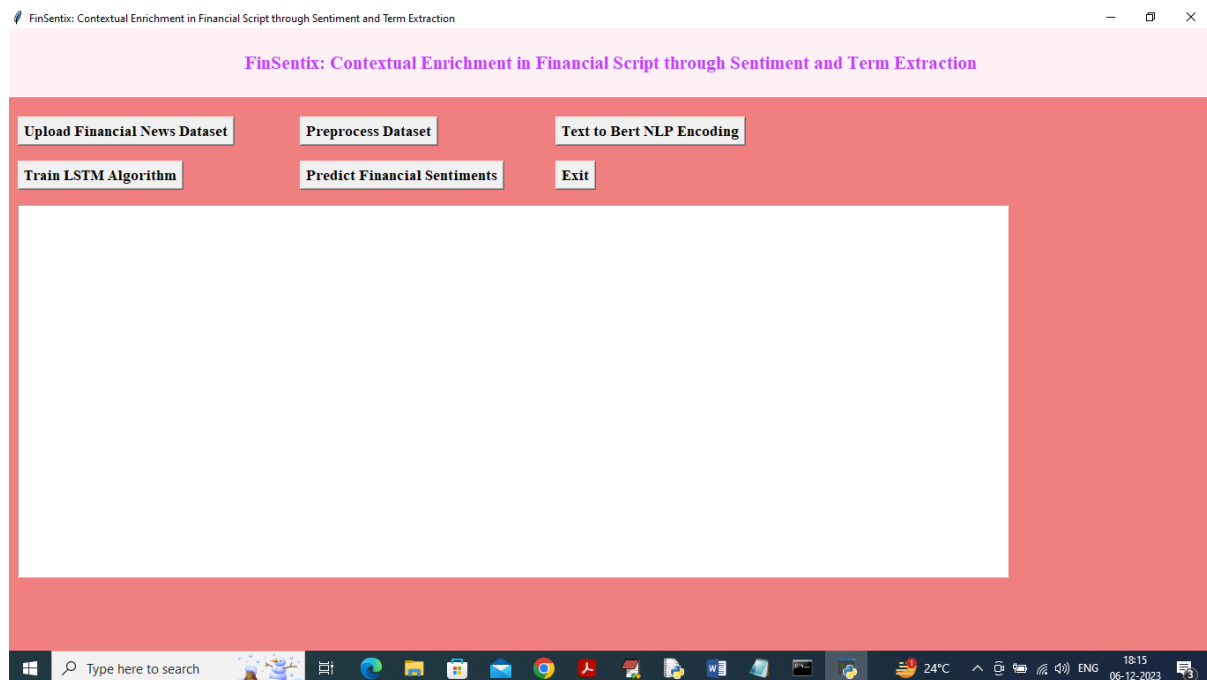


Fig 3.2.1 First Page

2) In above screen click on 'Upload Financial News Dataset' button to upload dataset and get below page

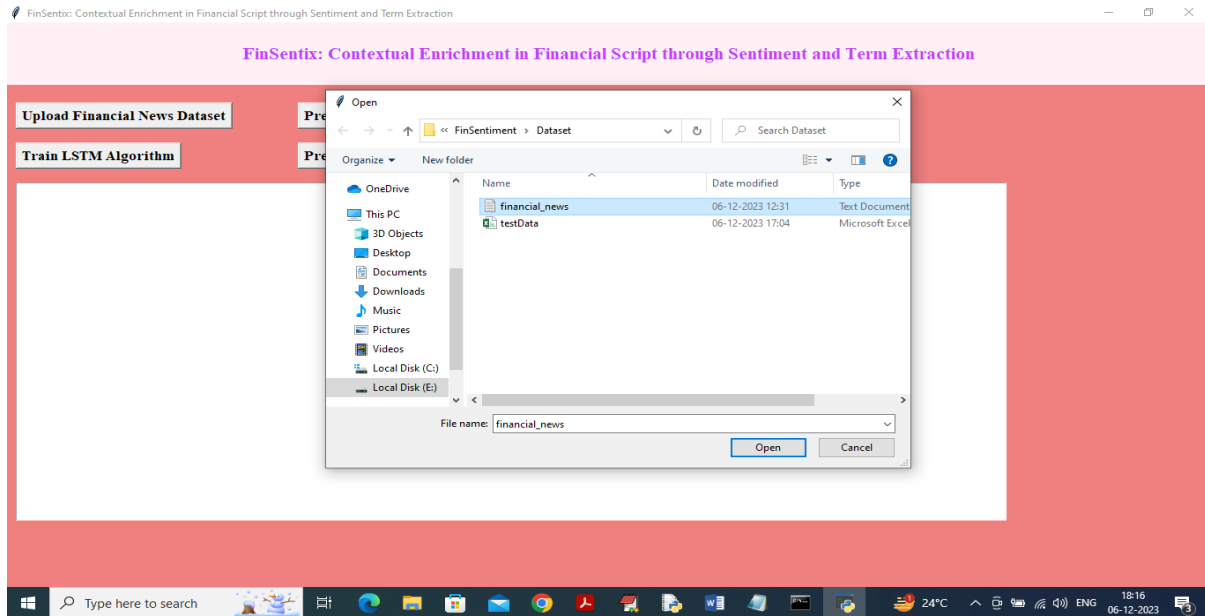


Fig 3.2.2 Upload Financial Dataset

3) In above screen selecting and uploading financial news dataset and then click on 'Open' button to load dataset and get below page

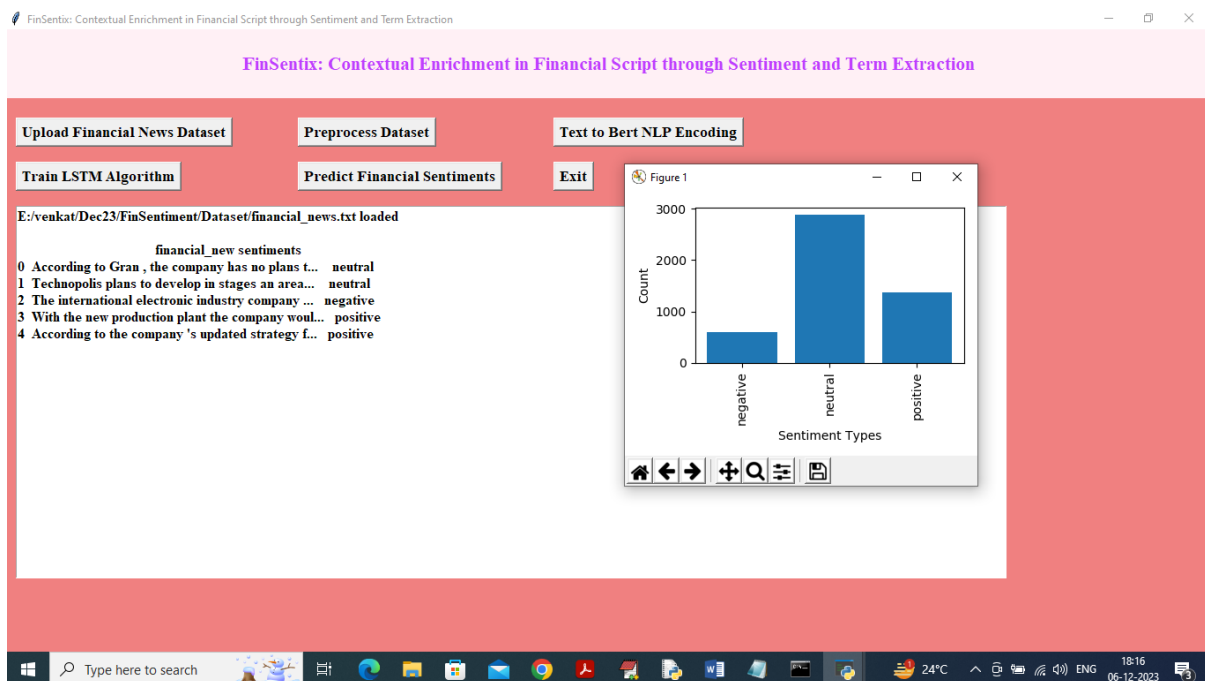


Fig 3.2.3 Load dataset and result for sentiments

4) In above screen financial news data loaded and displaying few records from them and in graph x-axis represents dataset sentiments type and y-axis represents count of those records and now close above graph and then click on 'Preprocess Dataset' button to clean NEWS test and get below output

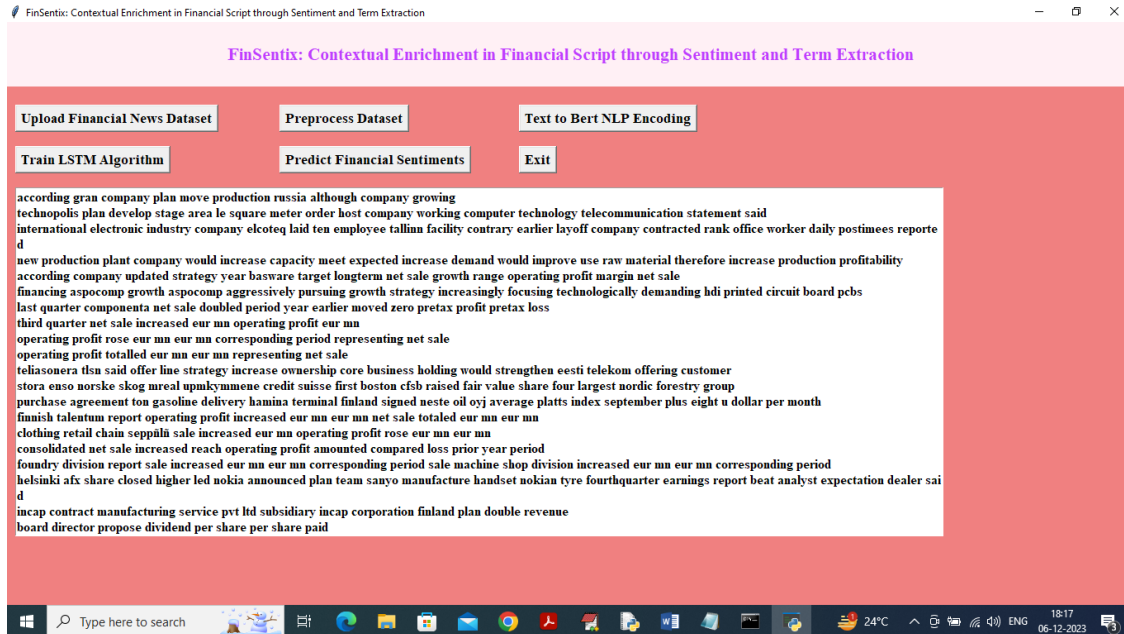


Fig 3.2.4 Load Data for preprocessing

5) In above screen each line represents one financial news and from that line we removed special symbols, stop words. Stemming and lemmatization applied and now click on 'Text to BERT NLP Encoding' button to convert text news to BERT vector and get below output

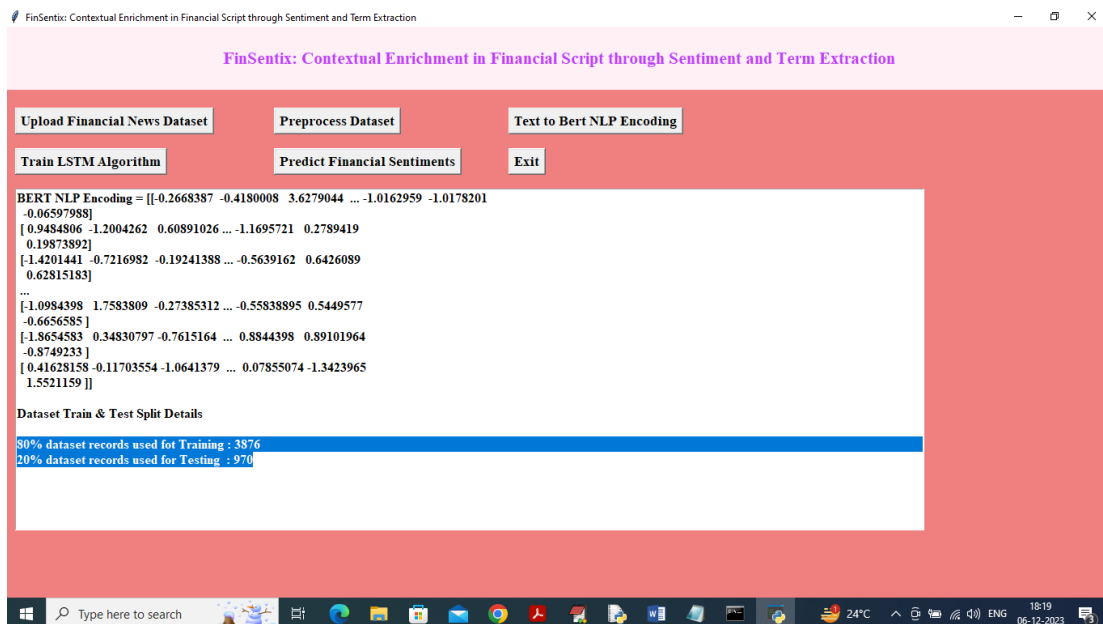


Fig 3.2.5 Vectors obtained

6) In above screen all text data converted to BERT encoding vector and then can see train and test size in blue colour and now click on 'Train LSTM Algorithm' button to train LSTM and get below output

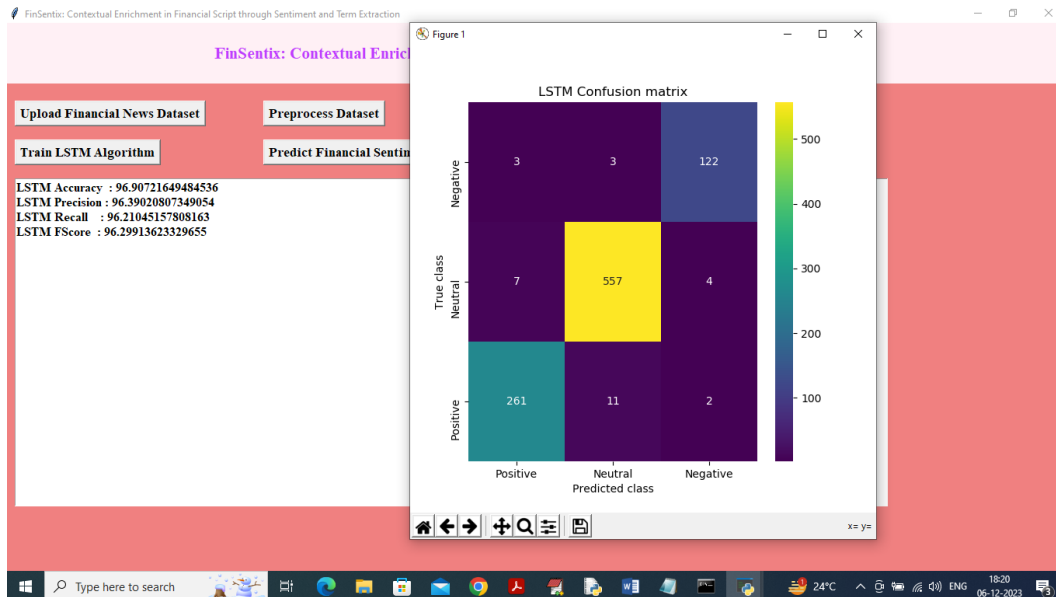


Fig 3.2.6 Confusion matrix for out ML model

7) In above screen LSTM got 96% accuracy and can see other metrics like precision, recall and FSCORE and in confusion matrix graph x-axis represents Predicted Labels and y-axis represents True Labels and all different colour boxes in diagonal represents correct prediction count and remaining blue boxes represents incorrect prediction count which are very few and now close above graph and then click on 'Predict Financial Sentiments' button to upload test data and get below output

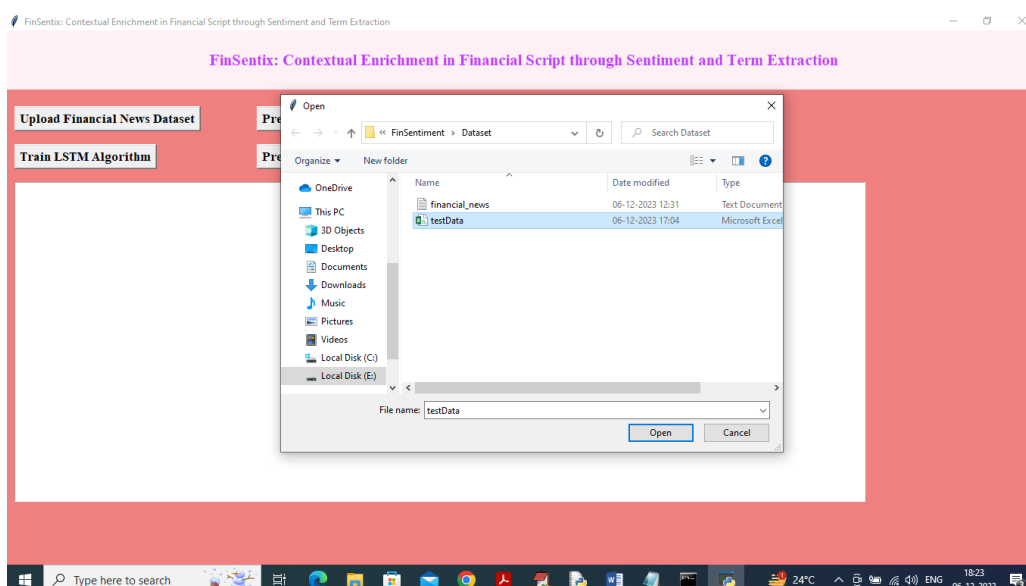


Fig 3.2.7 Upload data for data for validation

8) In above screen selecting and uploading test data and then click on 'Open' button to get below output

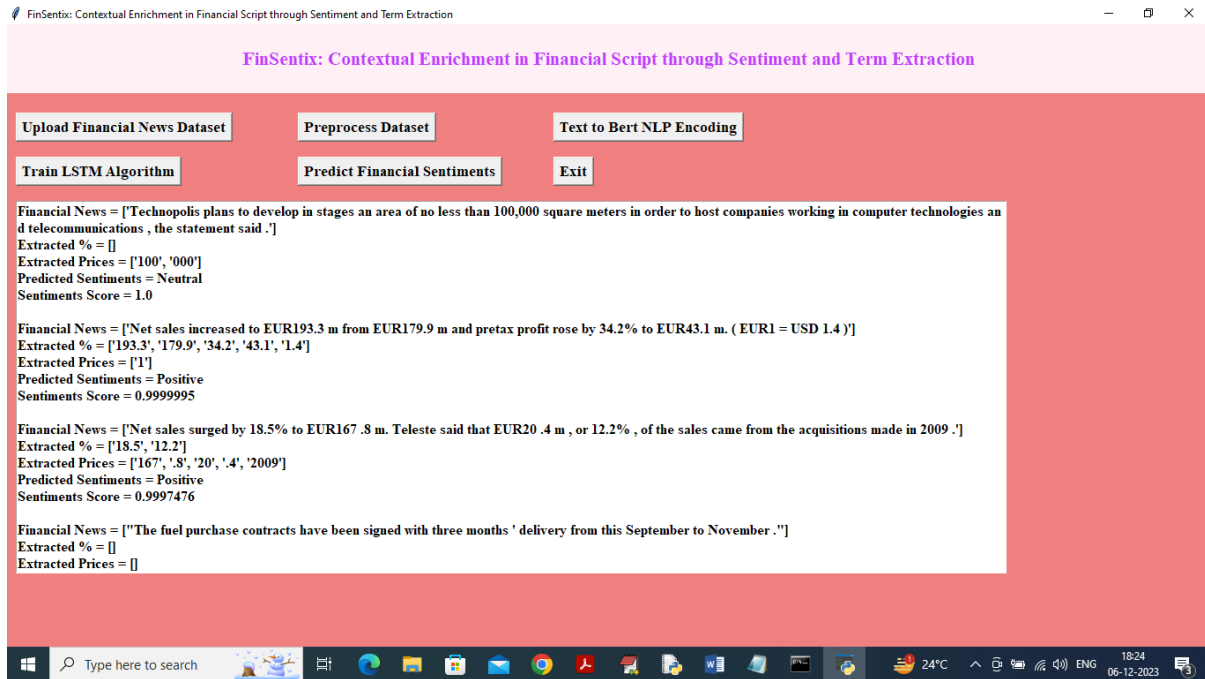


Fig 3.2.8 Results obtained for validation dataset

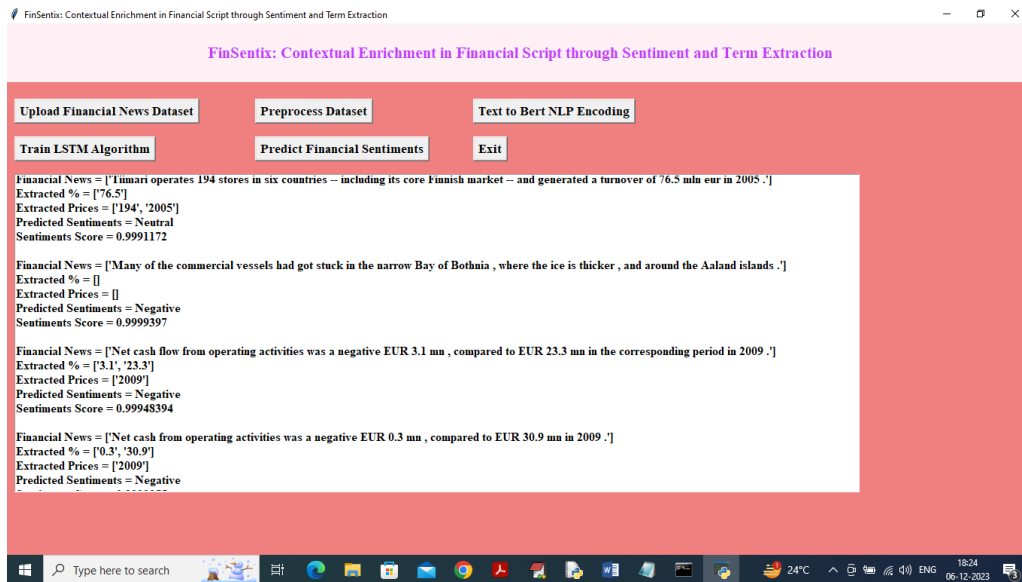


Fig 3.2.9 Result obtained for validation dataset

9) In above two screen we are displaying NEWS text from TEST data and then displaying extracted numeric %, prices, sentiments and sentiments predicted scores and this prediction you can see for each NEWS text line

Chapter 4

Conclusion and Scope of Improvement

4.1 Outcome and Impact

A financial entity named "Recognition Model" could engage in research focused on various aspects of finance, data analysis, and risk assessment. The research outcomes and their potential impact could include:

- **Advanced Credit Risk Assessment:** Developing innovative credit risk models that incorporate non-traditional data sources and machine learning techniques, leading to more accurate predictions of creditworthiness. The impact could be reduced default rates and improved lending decisions, ultimately increasing profitability.
- **Market Trend Prediction:** Developing predictive models using historical financial data and sentiment analysis to forecast market trends and investment opportunities. Accurate predictions can enhance portfolio management and investment decision-making, potentially leading to higher returns for investors.
- **Fraud Detection and Prevention:** Research on fraud detection algorithms and techniques can help financial institutions identify and prevent fraudulent activities more effectively. This can protect both the institution and its customers while reducing financial losses.
- **Algorithmic Trading Strategies:** Developing sophisticated trading algorithms based on historical data analysis and real-time market data. Successful strategies can lead to improved trading performance and potentially generate higher profits for investors.
- **Customer Behavior Analysis:** Research on customer behavior and preferences can help financial institutions tailor their products and services to better meet customer needs. Enhanced customer satisfaction and retention can positively impact the institution's bottom line.
- **Regulatory Compliance Solutions:** Developing tools and models that facilitate compliance with evolving financial regulations. This can reduce the risk of regulatory fines and penalties and ensure the institution operates within legal boundaries.
- **Risk Management Frameworks:** Designing comprehensive risk management frameworks that help financial entities identify, assess, and mitigate various types of risks, including market, credit, and operational risks. Effective risk management can safeguard the institution's stability and reputation.
- **Economic Forecasting:** Conducting research on economic indicators and developing models for economic forecasting. Accurate economic predictions can inform strategic planning and investment decisions.

The impact of these research outcomes would depend on their effectiveness and adoption within the financial industry. If Recognition Model can produce innovative, reliable, and practical solutions, it can potentially enhance operational efficiency, reduce risks, and improve financial performance for itself and its clients, thus establishing itself as a leading player in the financial sector.

4.2 Future Work

- 1. Enhanced Data Enrichment:** Develop more sophisticated methods to extract additional context, sentiment, and financial terms from text, improving the depth of information provided.
- 2. Real-time Integration:** Investigate the feasibility of real-time data integration and analysis to provide up-to-the-minute insights, particularly for financial news snippets.
- 3. Machine Learning Refinement:** Continuously refine the machine learning algorithms to enhance accuracy and speed of entity recognition, ensuring robustness across various financial document types.
- 4. Integration with Financial Tools:** Explore integration possibilities with financial analysis and reporting tools, facilitating seamless data flow for financial professionals.
- 5. Scalability:** Ensure the system's scalability to handle larger datasets and increasing demands in the financial industry.
- 6. Documentation and Support:** Develop comprehensive documentation and provide support to users for effective utilization of the system.
- 7. Feedback Mechanism:** Establish a feedback mechanism to gather input from users and stakeholders, ensuring continuous improvement and alignment with their needs.
- 8. Ethical Responsibility:** Ethical considerations and stringent data security measures are at the heart of our project, ensuring responsible handling of financial data.

These points highlight the various avenues for future development and enhancement of the "FinSentix: Contextual Enrichment in Financial Script through Sentiment and Term Extraction" project.

4.3 Conclusion

Finally, in the always changing field of financial data processing, the project "FinSentix: Contextual Enrichment in Financial Script through Sentiment and Term Extraction" shines as a shining example of innovation. The journey through this research has shed light on the urgent issues facing the financial sector, ranging from the inherent dangers of human tagging systems to the enormous volume of data.

We aim to change the financial data handling paradigm by automating numerical entity recognition, improving labelling accuracy, and enhancing contextual comprehension. FinSentix not only solves today's problems but also clears the path for a more accurate, transparent, and flexible financial ecosystem by streamlining operations, standardising data presentation, and guaranteeing regulatory compliance.

The thorough examination of the literature highlights the breadth of the field's research and development, highlighting the various methodologies and technological solutions used by prominent projects as AckNER, FinBERT, FiNER, BERT-based sentiment analysis, and others. These studies are important reference points that guide FinSentix's design decisions and methodology.

Looking ahead, the project has the potential to empower decision-makers in all financial domains. FinSentix touts itself as a catalyst for informed decision-making, risk reduction, and strategic financial analysis by offering previously unheard-of insights into numerical entities that are enhanced with sentiment and pertinent phrases.

References

- [1]Shah, A., Vithani, R., Gullapalli, A., and Chava, S., “FiNER: Financial Named Entity Recognition Dataset and Weak-Supervision Model”, <i>arXiv e-prints</i>, 2023. doi:10.48550/arXiv.2302.11157.
- [2]S. Wang, R. Xu, B. Liu, L. Gui and Y. Zhou, "Financial named entity recognition based on conditional random fields and information entropy," *2014 International Conference on Machine Learning and Cybernetics*, Lanzhou, China, 2014, pp. 838-843, doi: 10.1109/ICMLC.2014.7009718.
- [3]L. Zhao, L. Li, X. Zheng and J. Zhang, "A BERT based Sentiment Analysis and Key Entity Detection Approach for Online Financial Texts," *2021 IEEE 24th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, Dalian, China, 2021, pp. 1233-1238, doi: 10.1109/CSCWD49262.2021.9437616.
- [4]D. Alexander, A.P. de Vries, "Named Entity Recognition of Financial Information in Research Papers," *BIR 2021: 11th International Workshop on Bibliometric-enhanced Information Retrieval at ECIR 2021*, April 1, 2021, pp. 102-109. [Online]. Available: <https://ceur-ws.org/Vol-2875/paper13.pdf>. [Accessed: December 8, 2023].