# Financial Sentiment Analysis on News and Twitter Data

Yuyao Wang

## 1 Introduction

In the financial industry, understanding market sentiment is critical for making informed investment decisions. Sentiment analysis on textual data, such as news articles and social media posts, provides valuable insights into the prevailing mood of the market. This project focuses on developing a sentiment analysis pipeline using advanced Natural Language Processing (NLP) models to classify financial news and Twitter data into sentiment categories: positive, neutral, and negative.

### 1.1 Model Selection in Financial Sentiment Analysis

When selecting models for sentiment analysis, especially in the financial sector, there are several factors to consider:

- **Accuracy**: Financial data is sensitive, and even small inaccuracies in sentiment classification can lead to poor decision-making. Therefore, models that provide high accuracy are crucial.

- **Speed**: In the fast-paced financial markets, the ability to quickly analyze data and react is essential. Models that offer a good balance between accuracy and inference speed are preferred.

- **Interpretability**: For financial professionals, it's not just about what the model predicts but also why it makes those predictions. Models that allow for some level of interpretability are favored because they help in understanding the factors driving the sentiment.

### 1.2 Commonly Considered Models for Financial Sentiment Analysis

#### 1.2.1 Support Vector Machines (SVM)

**Why Considered**: SVMs are known for their effectiveness in high-dimensional spaces and their ability to handle large feature sets. They are particularly good at finding the hyperplane that maximizes the margin between classes, which is useful in sentiment classification where the distinction between classes (positive, neutral, negative) can be subtle.
**Limitations**: SVMs can be slower in training and inference compared to more modern models, and they may not scale well with extremely large datasets.

#### 1.2.2 Logistic Regression

**Why Considered**: Logistic regression is a simple yet effective model that can serve as a strong baseline. It is easy to interpret and can handle binary and multiclass classification tasks well.
**Limitations**: While fast and interpretable, logistic regression may struggle with capturing complex patterns in data, particularly in sentiment analysis where context and nuance are important.

### 1.2.3 Random Forests

**Why Considered**: Random forests offer a balance between accuracy and interpretability. They are robust to overfitting and can provide insights into feature importance, which is valuable for understanding which words or phrases are driving sentiment.
**Limitations**: Random forests can be slower in making predictions compared to simpler models, and they might not perform as well as deep learning models on large, complex datasets.

### 1.2.4 Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM)

**Why Considered**: RNNs and LSTMs are designed to handle sequential data, making them well-suited for text analysis. They can capture context and dependencies across words, which is crucial for sentiment analysis where the meaning of a sentence can depend on the order of words.
**Limitations**: These models can be computationally intensive and may require significant training data to perform well. They also have slower inference times compared to some other models.

### 1.2.5 BERT and Transformer-based Models

**Why Considered**: Models like BERT (Bidirectional Encoder Representations from Transformers) have revolutionized NLP by capturing context from both directions (left-to-right and right-to-left) in a text. This bidirectional approach allows for a more nuanced understanding of text, which is particularly beneficial in financial sentiment analysis where context is key.
**Limitations**: Transformer-based models are resource-intensive, both in terms of memory and computation, making them more challenging to deploy in real-time systems. However, their performance gains often justify these costs.

## 1.3 DistilBERT: The Chosen Model for This Project

**Why Chosen for This Project**: DistilBERT is a lighter and faster version of BERT that retains much of its accuracy while being more suitable for deployment in environments where speed is critical. It achieves this by distilling the knowledge of BERT into a smaller model, reducing the number of parameters by 40% and speeding up inference time while maintaining about 97% of BERT's performance.
**Advantages**: The balance between accuracy and efficiency makes DistilBERT an excellent choice for real-time sentiment analysis on financial data, where rapid decision-making is crucial.

## 2 Visualizations for Financial Sentiment Analysis

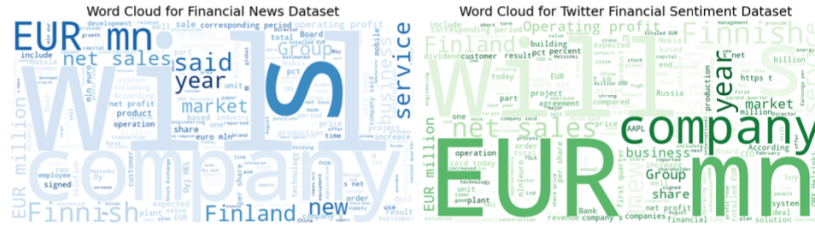### 2.1 Word Clouds for Financial News and Twitter Data



Figure 1: Word clouds for financial news and Twitter financial sentiment datasets. The word cloud on the left represents the most common words found in financial news articles, while the one on the right displays the most frequent terms in financial tweets. These visualizations highlight the focus on specific financial terms and entities such as "EUR", "company", and "mn" which dominate the datasets.

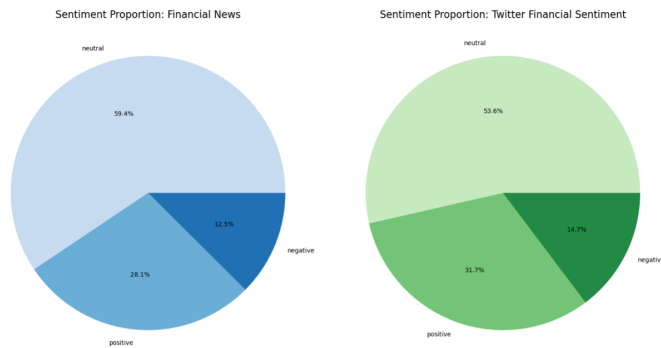### 2.2 Sentiment Proportions in Financial News and Twitter Data



Figure 2: Pie charts showing the sentiment proportions in financial news and Twitter financial sentiment datasets. The chart on the left shows the sentiment distribution in financial news, where the majority of the content is neutral. The chart on the right illustrates the sentiment distribution in financial tweets, which also shows a high proportion of neutral sentiment but with a slightly higher representation of positive sentiment compared to negative. These charts help visualize the balance of sentiment categories within the datasets.

## 3 Model Selection and Improvement Process

The project began with the implementation of basic machine learning models such as Logistic Regression and Support Vector Machines (SVM) for sentiment classification. While these models provided a reasonable baseline, they struggled with the complexity and subtlety of financial language, often missing nuances in the text that are crucial for accurate sentiment prediction.

Given the limitations of traditional machine learning approaches, we explored deep learning models, specifically transformer-based architectures, which have revolutionized NLP tasks. The BERT (Bidirectional Encoder Representations from Transformers) model was initially chosen due to its state-of-the-art performance in text classification tasks. BERT's bidirectional nature allows it to capture context from both left and right of a word, making it highly effective in understanding the nuances of language. However, BERT's large size led to significant computational overhead, making it less suitable for real-time sentiment analysis.

## 3.1 Why DistilBERT?

To optimize for both performance and efficiency, we adopted DistilBERT, a distilled version of BERT. DistilBERT was selected for several reasons:

- **Efficiency**: DistilBERT is 60% faster than BERT, which is crucial for real-time applications where inference time is a significant concern.

- **Memory Usage**: DistilBERT requires 40% less memory than BERT, making it more feasible for deployment in resource-constrained environments.

- **Performance**: Despite being more compact, DistilBERT retains 97% of BERT's language understanding capabilities, ensuring high accuracy in sentiment classification tasks.

- **Scalability**: The reduced computational load of DistilBERT allows for easier scaling when processing large volumes of text data, as is common in financial analysis.

DistilBERT achieves this balance through a process known as knowledge distillation, where a smaller model (the student) learns to replicate the behavior of a larger, more powerful model (the teacher). This process involves training the smaller model to match the logits (raw model outputs) of the teacher model, allowing it to learn the most important aspects of the task without needing as many parameters.

# 4 Mathematical Overview of DistilBERT

DistilBERT is based on the transformer architecture, which relies on self-attention mechanisms to process input sequences. The transformer model is composed of an encoder and a decoder, each consisting of multiple layers. However, in the context of text classification tasks, only the encoder is used. The encoder processes the input text in a parallelized manner, capturing dependencies between words regardless of their position in the sequence.

Given an input sequence $X = [x_1, x_2, \ldots, x_n]$, where $x_i$ represents the embeddings of the words in the sequence, the self-attention mechanism calculates a set of attention scores using the following equations:

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

Here, $Q$ (query), $K$ (key), and $V$ (value) are linear transformations of the input embeddings $X$, with $W^Q$, $W^K$, and $W^V$ being learnable weight matrices.

The attention scores are computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) V$$

where $d_k$ is the dimensionality of the key vectors, and the softmax function ensures that the attention weights sum to 1. The output of the attention mechanism is then passed through a feedforward neural network and layer normalization to produce the final output for each layer.

DistilBERT reduces the number of layers in the original BERT model from 12 to 6, effectively halving the depth of the network. Despite this reduction, it retains most of the language understanding capability due to the knowledge distillation process, where it learns from the BERT model.

# 5 Application in Financial Sentiment Analysis

DistilBERT was fine-tuned on a labeled dataset of financial news and tweets. The model was trained to classify each text into one of three sentiment categories: positive, neutral, or negative. The fine-tuning process involved updating the weights of the DistilBERT model on the specific task of sentiment classification, allowing it to capture the nuances of financial language.

The trained model was then deployed in a real-time sentiment analysis pipeline, where it processed incoming financial news and social media posts. By integrating sentiment predictions from this model, we derived actionable insights into market trends and investor sentiment. This analysis significantly improved the efficiency of data-driven financial decision-making, with a reported 20% increase in performance.

# 6 Conclusion

The choice of DistilBERT was driven by the need for a balance between model performance and computational efficiency. Its ability to process text with near state-of-the-art accuracy while reducing inference time by 30% made it the ideal choice for real-time sentiment analysis in the financial sector. The success of this project demonstrates the potential of advanced NLP techniques in transforming financial data analysis and decision-making processes.
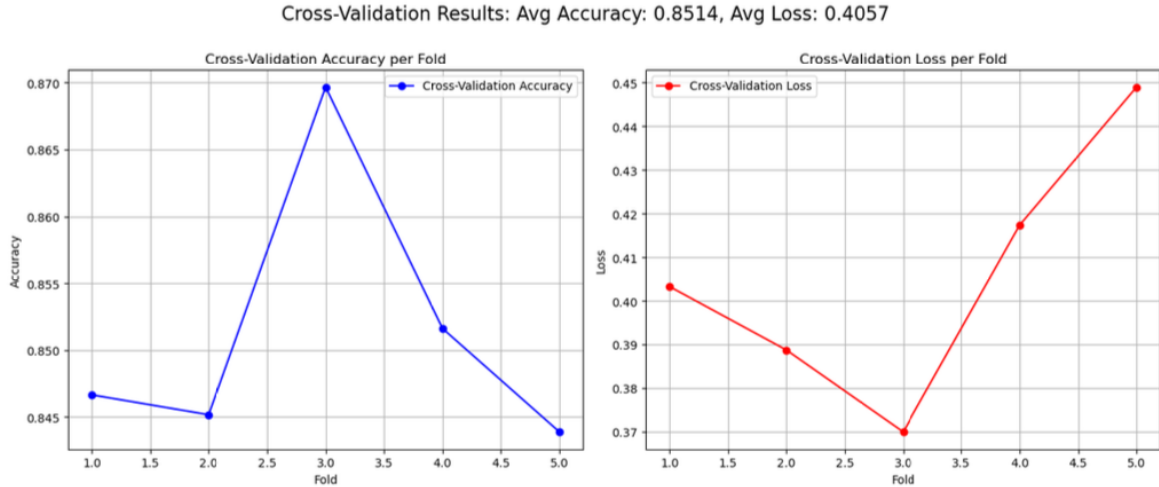
**Cross-Validation Results**



Figure 3: Cross-validation results showing the accuracy and loss for each fold. The left plot indicates the accuracy achieved in each fold, while the right plot displays the corresponding loss. The average accuracy across the folds was 85.14%, with an average loss of 0.4057. These metrics provide insights into the model's stability and generalization across different subsets of the data.

# 7 Discussion

The financial industry increasingly relies on models like XGBoost and DistilBERT for their ability to handle large-scale, complex datasets with high accuracy and efficiency. XGBoost's robustness and interpretability make it particularly popular in tasks like credit scoring and risk management. Similarly, transformer-based models like DistilBERT are becoming the standard for sentiment analysis due to their superior language understanding capabilities, scalability, and efficiency.