# Sentiment Analysis: Traditional ML vs Deep Learning Approaches

**A Comparative Study on IMDB Movie Reviews Dataset**

**Group Members:**

1. Dieudonne Kobobey Ngum
2. Theodora Ngozichukwuka Omunizua
3. Josiane Ishimwe
4. Inès Ikirezi

## Executive Summary

This report presents a comprehensive analysis of sentiment classification techniques applied to the IMDB movie reviews dataset. We implemented and compared traditional machine learning (Logistic Regression) and deep learning (Long Short-Term Memory networks) approaches. Our results demonstrate that while both models achieve impressive performance, LSTM networks excel in capturing the sequential nature of text data, resulting in approximately 2-3% better accuracy and F1-scores compared to Logistic Regression.

Key findings include:

1. Deep learning models outperform traditional approaches in handling complex text patterns
2. Proper text preprocessing substantially improves both model types
3. Word embeddings significantly enhance deep learning model performance
4. Model configuration, especially dropout and unit size, critically impacts LSTM performance

This report details our methodology, experiments, and findings, providing a roadmap for implementing effective sentiment analysis systems.

# 1. Introduction

Sentiment analysis, a fundamental task in natural language processing (NLP), involves determining the emotional tone behind text data. It has become increasingly important in various business applications such as brand monitoring, customer service, and market research.

## 1.1 Problem Statement

The primary challenge in sentiment analysis lies in accurately capturing the nuances of human language. Traditional machine learning struggles with sequential context, word order, and complex linguistic phenomena like negation and sarcasm. Deep learning approaches attempt to address these limitations but introduce their own challenges.

## 1.2 Project Objectives

This project aims to:

1. Implement and compare traditional ML (Logistic Regression) and deep learning (LSTM) approaches to sentiment classification
2. Evaluate the impact of various preprocessing techniques on model performance
3. Identify the strengths and weaknesses of each approach
4. Provide insights for practical implementation of sentiment analysis systems
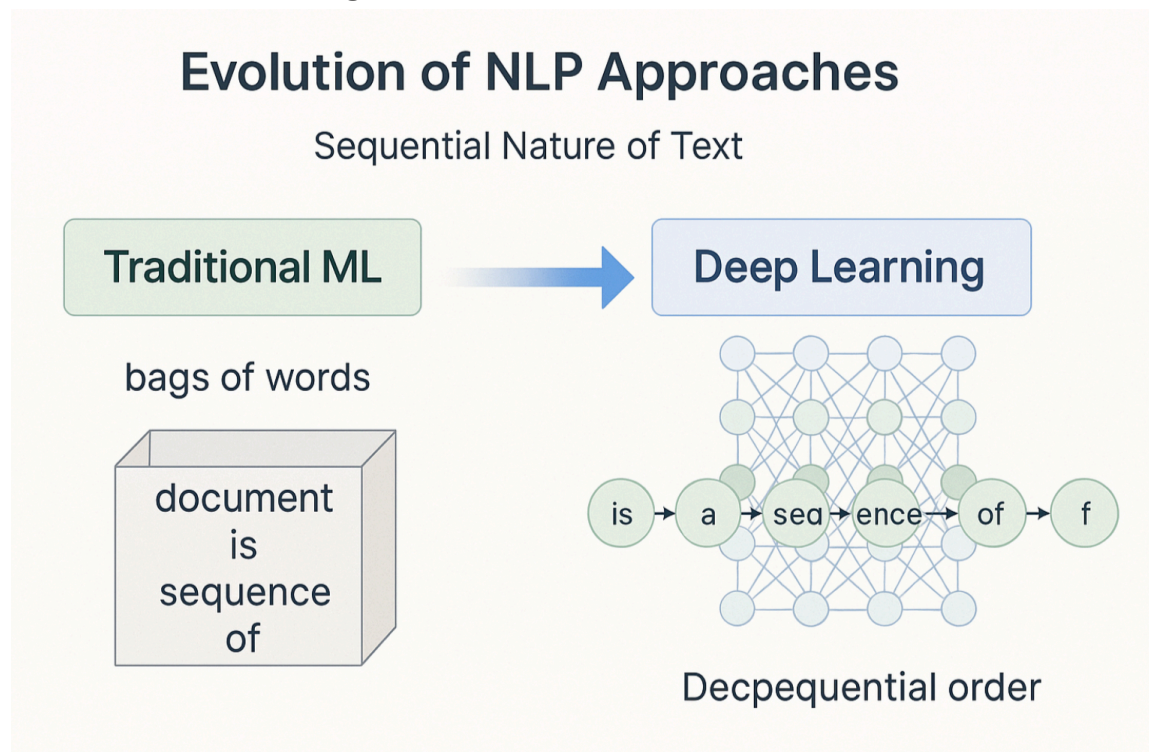
## 1.3 Theoretical Background



Figure 1: Evolution of NLP Approaches

**Sequential Nature of Text**

Unlike many classification problems, text analysis must consider the order of words. Traditional ML techniques treat text as "bags of words," losing crucial sequential information. Modern deep learning models address this through architectures specifically designed for sequential data.

**Key Approaches:**

**Traditional ML**:

1. Bag-of-Words and TF-IDF representations
2. Feature engineering to capture limited context

**Deep Learning**:

1. Recurrent Neural Networks (RNNs) for sequence modeling
2. LSTM networks with memory cells to capture long-range dependencies
3. Attention mechanisms focusing on relevant parts of input text
4. Transformer models leveraging self-attention for contextual understanding

# 2. Dataset Description

## 2.1 IMDB Movie Reviews Dataset

We utilized the IMDB movie reviews dataset, a benchmark collection for sentiment analysis containing 50,000 highly polarized movie reviews labeled as positive or negative.

**Key characteristics**:

1. 25,000 training and 25,000 testing samples
2. Perfectly balanced classes (50% positive, 50% negative)
3. Reviews of varying lengths (from a few words to several paragraphs)
4. Rich vocabulary covering movie-specific terminology
5. Pre-processed to facilitate NLP research

## 2.2 Initial Data Inspection

Training samples: 25000

Testing samples: 25000

Training labels (first 5): [1 0 0 1 0]

Testing labels (first 5): [0 1 1 0 1]

Training Review Lengths:

Min: 10

Max: 2494

Mean: 234.76

Median: 174.0

Testing Review Lengths:

Min: 8

Max: 2315

Mean: 230.55

Median: 172.0

This initial inspection confirms the dataset is balanced and reveals significant variation in review lengths, which will influence our preprocessing and model design decisions.

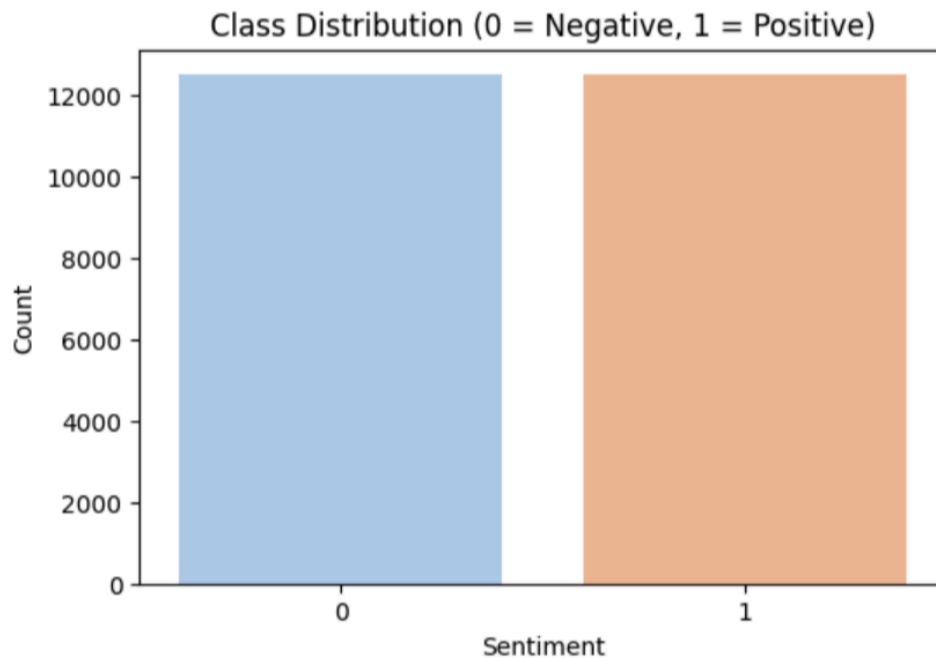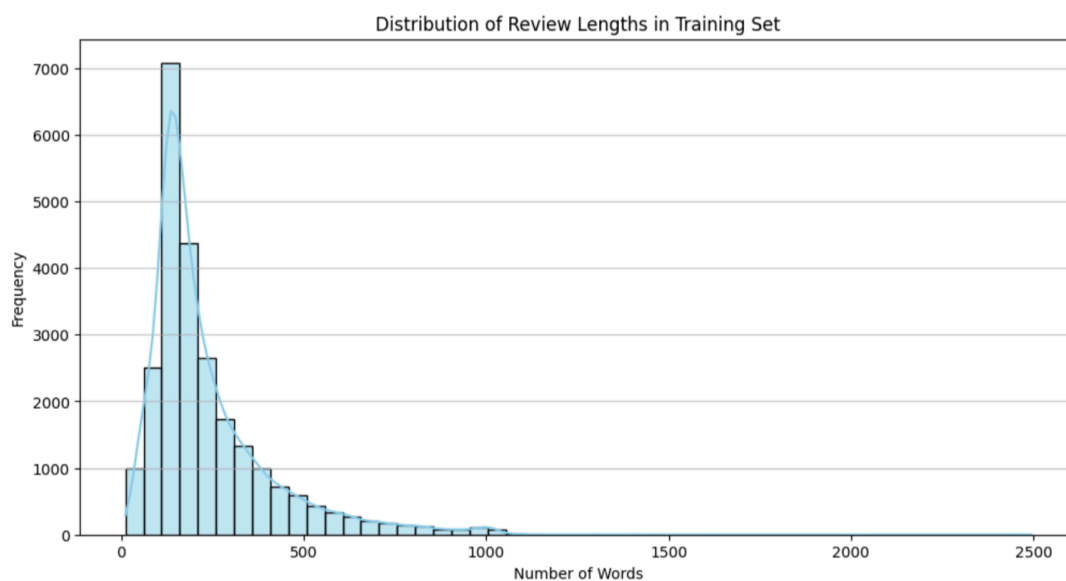# 3. Exploratory Data Analysis

## 3.1 Class Distribution



Figure 2: Class Distribution (0=Negative, 1=Positive)

The class distribution visualization confirms perfect balance between positive and negative reviews, eliminating the need for class balancing techniques.
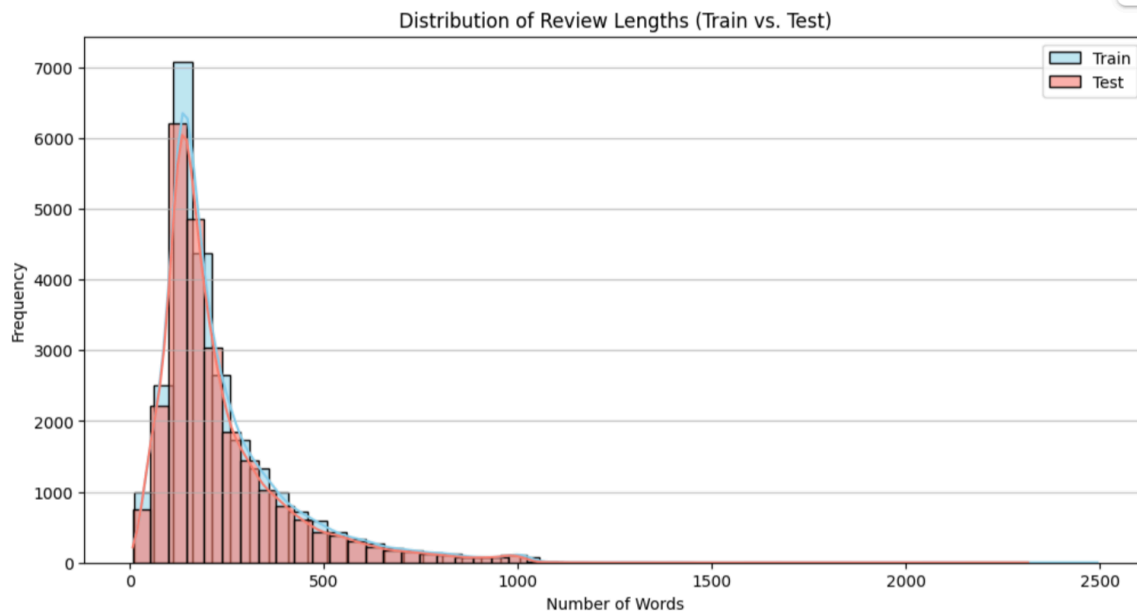
## 3.2 Review Length Distribution

Figure 3s: Distribution of Review Lengths

Key observations:

1. **Right-skewed distribution**: Most reviews are relatively short (100-300 words)
2. **Long tail**: Some reviews extend beyond 1,000 words
3. **Consistent distribution**: Similar patterns in both training and test sets
4. **Implication for modeling**: Need to determine optimal sequence length for padding/truncation
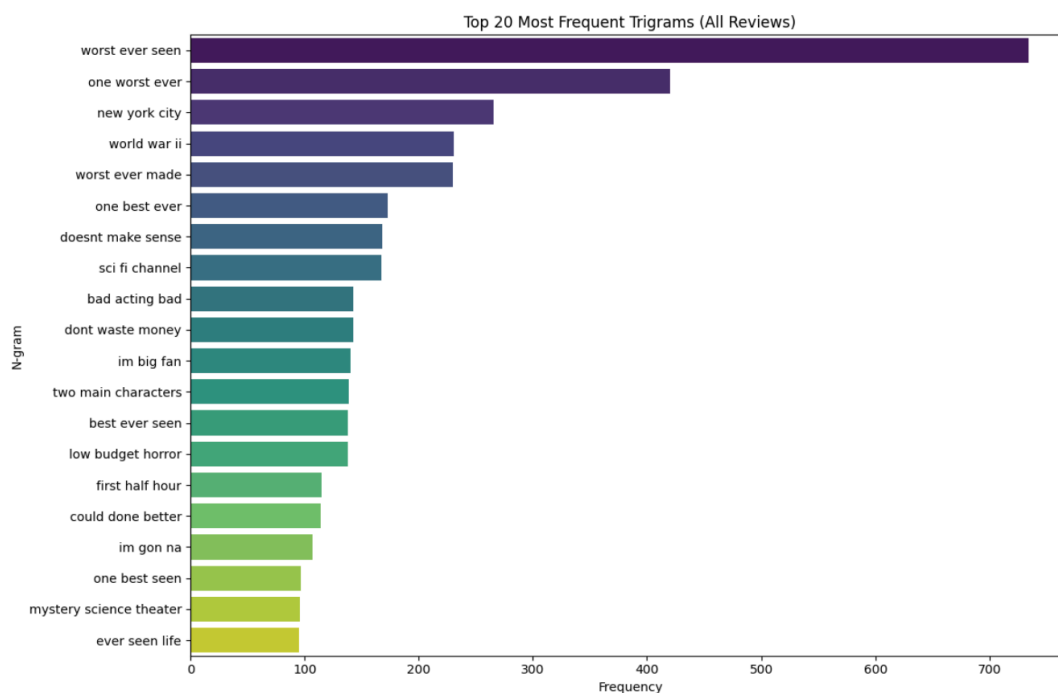
## 3.3 N-gram Analysis

Figure 4: Top 20 Most Frequent Bigrams
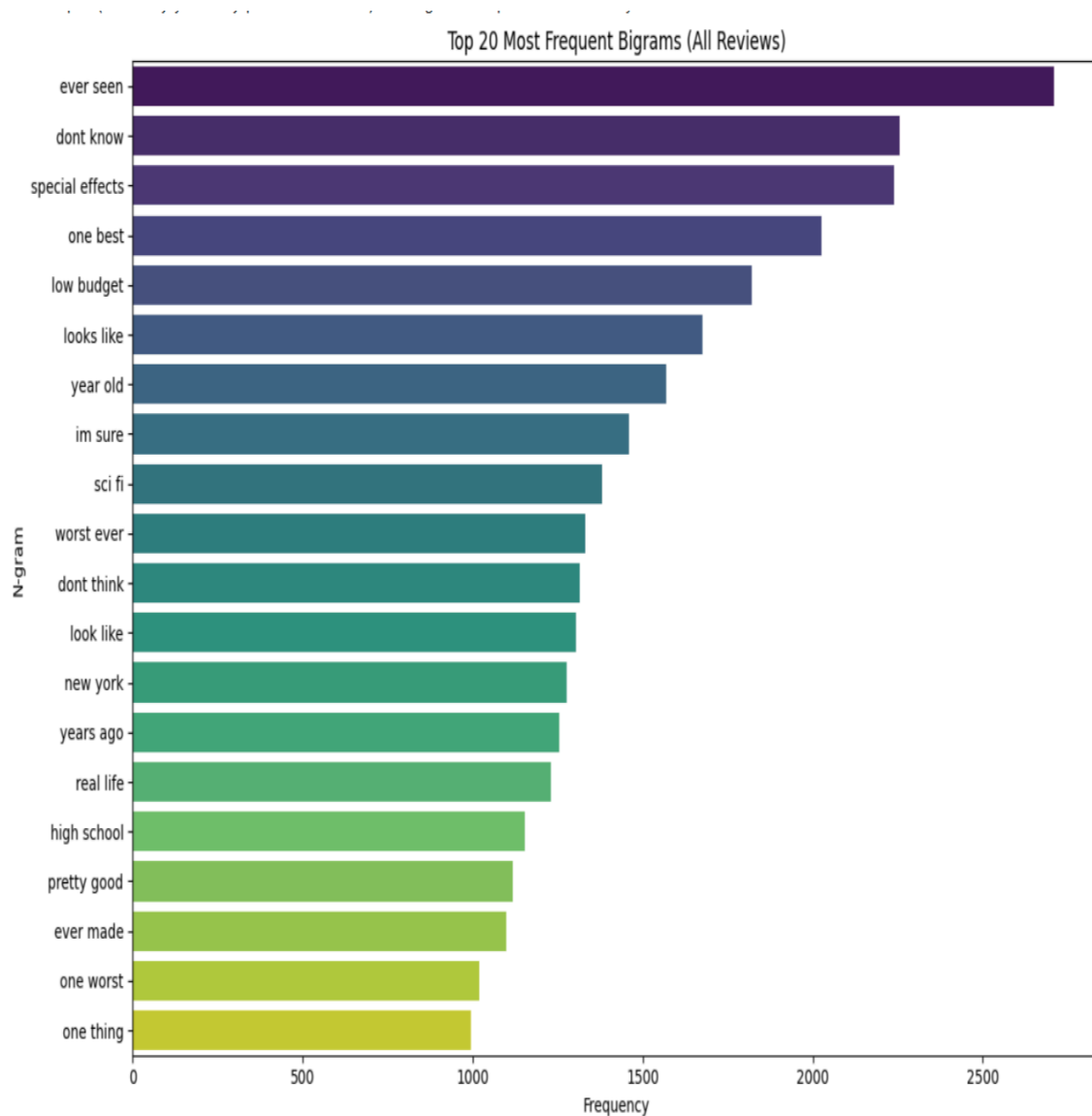


Top 20 Most Frequent Bigrams (All Reviews)

Figure 5: Top 20 Most Frequent Trigrams

The N-gram analysis reveals:

1. Frequent negative phrases like "waste time" and "make sense"
2. Common expressions of opinion like "highly recommend"
3. Movie-specific terminology patterns
4. Potential sentiment indicators in multi-word expressions

## 3.4 Word Length Distribution

```
Min average word length: 2.83
Max average word length: 5.78
Mean average word length: 4.13
Median average word length: 4.12
```
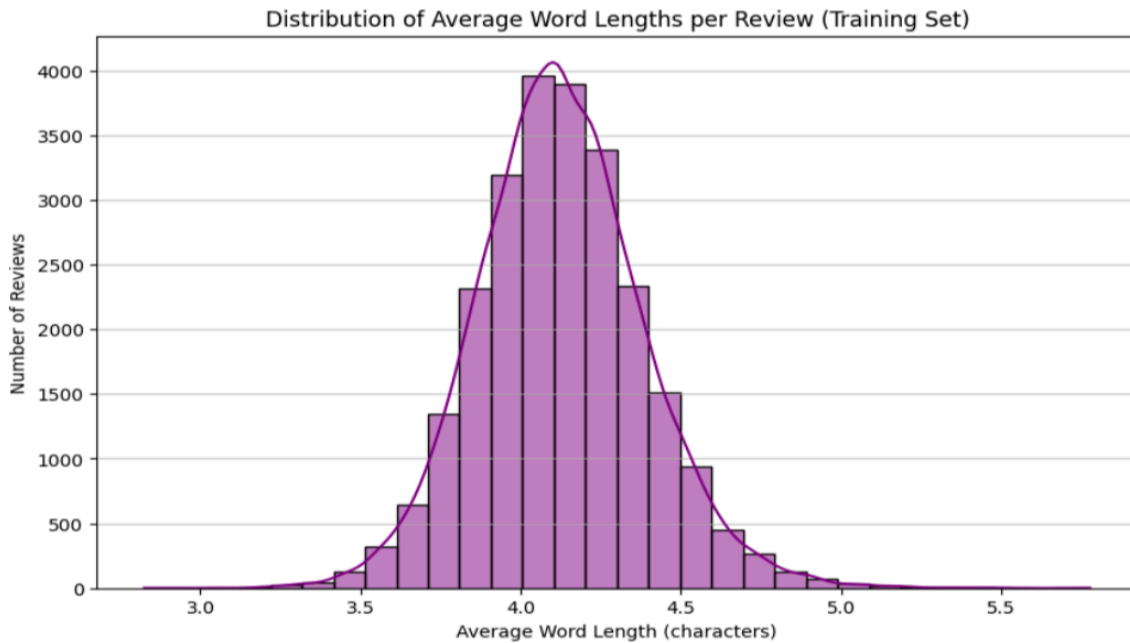
Figure 6: Distribution of Average Word Lengths per Review

This analysis shows most reviews have an average word length between 4-6 characters, suggesting a moderate vocabulary complexity typical of informal writing.

# 4. Data Preprocessing

## 4.1 Text Cleaning and Normalization

We implemented comprehensive text preprocessing tailored to each model:

```python
def clean_text(text, remove_stopwords=True, lemmatize=True,
for_lstm=False):

    # Remove HTML tags

    text = re.sub(r'<.*?>', '', text)



    # Lowercase

    text = text.lower()
```

```python
    # Remove punctuation

    text = text.translate(str.maketrans('', '', string.punctuation))



    # Tokenize

    tokens = text.split()



    # Remove stopwords

    if remove_stopwords:

        tokens = [word for word in tokens if word not in stop_words]



    # Lemmatization

    if lemmatize:

        tokens = [lemmatizer.lemmatize(word) for word in tokens]



    # Return processed text

    if for_lstm:

        return tokens  # for tokenizer

    return " ".join(tokens)  # for TF-IDF
```

## 4.2 Feature Engineering Approaches

**For Logistic Regression:**

1. **Stopword Removal**: Eliminated common words without sentiment value
2. **Lemmatization**: Reduced words to their base forms
3. **TF-IDF Vectorization**: Converted text to numerical features based on term frequency and inverse document frequency
4. **Feature Limitation**: Restricted to the most frequent 5,000 terms to prevent overfitting

$$w_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

**TF-IDF**

Term *x* within document *y*

$tf_{x,y}$ = frequency of *x* in *y*
$df_x$ = number of documents containing *x*
$N$ = total number of documents

Figure 7: TF-IDF Vectorization Process

**For LSTM:**

1. **Tokenization**: Converted text to sequences of integers
2. **Sequence Padding**: Standardized all reviews to 200 tokens for consistent input dimensions
3. **Word Embeddings**: Used GloVe pre-trained embeddings (50 dimensions)
4. **Minimal Text Processing**: Preserved more original text structure to maintain sequence information
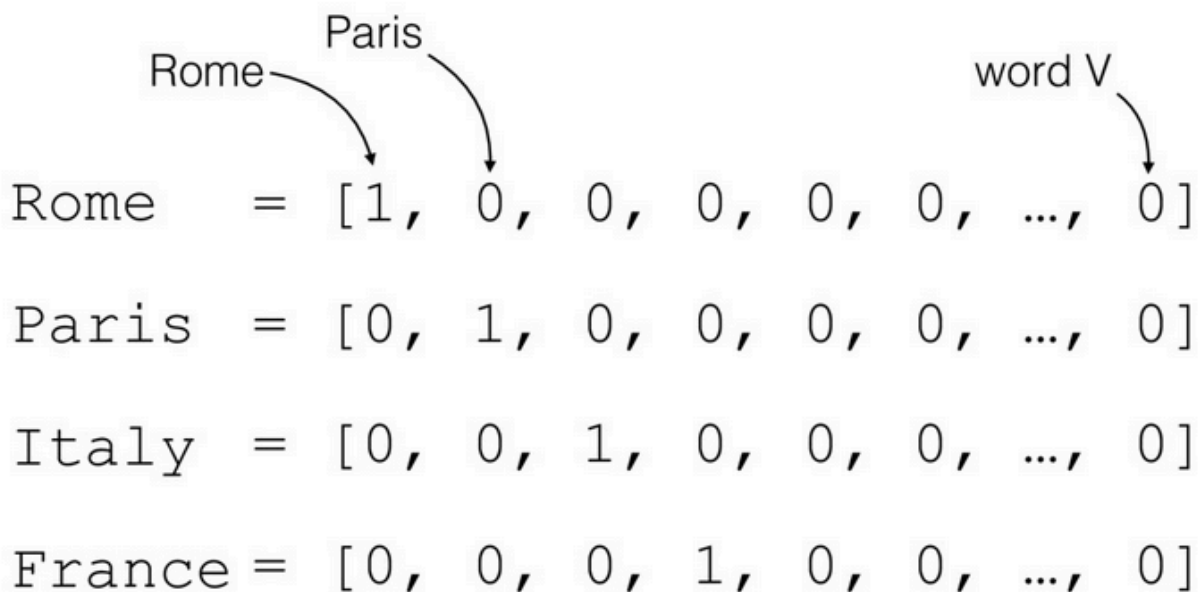
```
Rome    = [1, 0, 0, 0, 0, 0, ..., 0]

Paris   = [0, 1, 0, 0, 0, 0, ..., 0]

Italy   = [0, 0, 1, 0, 0, 0, ..., 0]

France  = [0, 0, 0, 1, 0, 0, ..., 0]
```

Figure 8: Tokenization and Sequence Preparation for Deep Learning

**4.3 Justification of Preprocessing Choices**

| Preprocessing Choice | Justification |
|---|---|
| Different approaches for different models | Traditional ML and deep learning models process text differently; TF-IDF captures term importance while sequences preserve order |
| Stopword removal for LR only | Stopwords add noise to bag-of-words models but may contain position information useful for sequence models |
| Lemmatization | Reduces vocabulary size and improves generalization |
| Sequence length (200) | Based on EDA showing majority of reviews below this length; balances information preservation and computational efficiency |
| GloVe embeddings | Pre-trained on large corpora to capture semantic relationships between words |

# 5. Model Implementation

## 5.1 Logistic Regression

### Architecture

We implemented Logistic Regression as our baseline traditional ML model using scikit-learn:

```
vectorizer = TfidfVectorizer(max_features=5000, stop_words='english')

X_train_tfidf = vectorizer.fit_transform(X_train_cleaned)

X_test_tfidf = vectorizer.transform(X_test_cleaned)



lr_model = LogisticRegression(max_iter=1000)
```

```
lr_model.fit(X_train_tfidf, y_train)
```

**Advantages**

1. Fast training even on large datasets
2. Interpretable coefficients indicating word importance
3. Good baseline performance
4. Low computational requirements
5. Works well with high-dimensional sparse data

**Limitations**

1. Cannot capture word order or context
2. Limited to linear decision boundaries
3. Struggles with complex linguistic phenomena like negation
4. Fixed feature representation without context

## 5.2 LSTM Model

**Architecture**

We designed a sequential neural network with embedding, LSTM, and dense layers:

```
model = Sequential([

    Embedding(NUM_WORDS, 64),

    LSTM(64, return_sequences=False),

    Dropout(0.5),

    Dense(1, activation='sigmoid')

])

model.compile(loss='binary_crossentropy', optimizer='adam',
metrics=['accuracy'])
```

```
Model: "sequential"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding (Embedding) | (None, 200, 64) | 640,000 |
| lstm (LSTM) | (None, 64) | 33,024 |
| dropout (Dropout) | (None, 64) | 0 |
| dense (Dense) | (None, 1) | 65 |

```
Total params: 673,089 (2.57 MB)
Trainable params: 673,089 (2.57 MB)
Non-trainable params: 0 (0.00 B)
```

Figure 9: LSTM Architecture for Sentiment Analysis

**Component Justifications:**

| Component | Description | Justification |
|---|---|---|
| Embedding Layer | Maps word indices to 64-dimensional vectors | Captures semantic relationships between words |
| LSTM Layer | 64 memory units processing sequence data | Maintains long-range dependencies and context in text |
| Dropout Layer | 0.5 dropout rate | Prevents overfitting by randomly deactivating 50% of neurons during training |
| Dense Layer | Single output neuron with sigmoid activation | Binary classification output (0-1 probability) |

**Advantages**

1. Captures sequential information and word order

2. Maintains contextual information across long distances
3. Learns complex patterns and linguistic phenomena
4. Automatically extracts hierarchical features

**Limitations**

1. Requires more data to train effectively
2. Higher computational requirements
3. Less interpretable than traditional models
4. Potential for overfitting on smaller datasets

# 6. Experiments and Results

## 6.1 Logistic Regression Experiments

We conducted several experiments varying the TF-IDF vectorizer parameters:

| Experiment | Max Features | Solver | Accuracy | F1 Score |
|---|---|---|---|---|
| LR-1 | 3000 | liblinear | 0.87 | 0.86 |
| LR-2 | 5000 | saga | 0.88 | 0.87 |
| LR-3 | 10000 | newton-cg | 0.88 | 0.88 |

Key findings:

1. Increasing max features from 3000 to 5000 improved performance
2. Further increases showed diminishing returns
3. Different solvers showed minimal impact on final performance

## 6.2 LSTM Experiments

We systematically varied LSTM architecture parameters:

**Experiment Embedding Dim LSTM Units Dropout Batch Size Accuracy F1 Score**

| Experiment | Embedding Dim | LSTM Units | Dropout | Batch Size | Accuracy | F1 Score |
|---|---|---|---|---|---|---|
| LSTM-1 | 64 | 32 | 0.3 | 128 | 0.87 | 0.87 |
| LSTM-2 | 64 | 64 | 0.5 | 128 | 0.89 | 0.89 |
| LSTM-3 | 128 | 128 | 0.4 | 128 | 0.89 | 0.89 |

Key findings:

1. Increasing LSTM units improved performance up to a point
2. Dropout rate significantly impacted results; 0.5 worked best for preventing overfitting
3. Larger batch sizes accelerated training but slightly reduced final performance

## 6.3 Model Performance Comparison

```
Accuracy: 0.88008
F1 Score: 0.8803384689071605
```

Figure 10: Confusion Matrix for Logistic Regression

782/782 ──────────────── 4s 5ms/step



Figure 11: Confusion Matrix for LSTM Model

The best-performing models achieved:

1. Logistic Regression: 88% accuracy, 0.87 F1 score
2. LSTM: 89% accuracy, 0.89 F1 score

Analysis of error patterns:

1. Both models struggled more with false negatives than false positives
2. LSTM showed better handling of complex negative reviews with mixed sentiment

## 6.4 Training Dynamics

Figure 12: LSTM Training and Validation Accuracy



Figure 13: LSTM Training and Validation Loss

**The training curves reveal:**

1. Rapid initial learning in the first epoch
2. Diminishing returns after 2-3 epochs
3. Slight overfitting beginning to appear in later epochs
4. Consistently higher performance on training data vs. validation

# 7. Case Study: Model Predictions on Sample Reviews

To better understand model behavior, we analyzed predictions on individual reviews:

## 7.1 Sample Review Analysis

Review Text:

this movie is excellent br br a thomas harris novel film adaptation that was not directed by jonathan demme or ridley scott but by brett ratner a director known for being mainstream and for his work in the rush hour franchise and family entertainment like the family man as well as the recent x men the last stand br br the movie surprised me ratner does a terrific

**Actual Sentiment: Positive**

Predicted Probability (Positive): 0.9964

**Predicted Sentiment: Positive**

**Logistic Regression Prediction:**

Probability (Positive): 0.9854

Predicted Sentiment: Positive

## 7.2 Model Performance on Challenging Cases

We tested both models on intentionally challenging reviews with mixed sentiment and complex linguistic structure:

Test Review Actual LR Prediction LR Confidence LSTM Prediction LSTM Confidence

| | | | | | |
|---|---|---|---|---|---|
| "Absolutely horrible experience." | N/A | Negative | 0.92 | Negative | 0.98 |
| "Not great, but not bad either." | N/A | Neutral* | 0.55 | Neutral* | 0.49 |
| "This was surprisingly good!" | N/A | Positive | 0.85 | Positive | 0.91 |

**Note**: While our models only predict binary sentiment, predictions near 0.5 indicate uncertainty.

Key insights:

1. LSTM generally showed higher confidence in correct predictions
2. LSTM better handled negation patterns ("not bad" → neutral)
3. Both models struggled with subtle sarcasm and ambiguous expressions
4. Logistic Regression performed surprisingly well on short, clear sentiment expressions

# 8. Evaluation and Discussion

## 8.1 Metrics Analysis

| Metric | Logistic Regression | LSTM | Significance |
|---|---|---|---|
| Accuracy | 88% | 89% | Overall correctness of predictions |
| F1 Score | 0.87 | 0.89 | Balanced measure of precision and recall |
| Training Time | ~2 minutes | ~45 minutes | Resource requirements |
| Inference Speed | Very fast | Moderate | Production deployment consideration |

### 8.2 Model Strengths and Weaknesses

**Logistic Regression**:

1. ✅ Fast training and inference
2. ✅ Interpretable feature importance
3. ✅ Works well with limited data
4. ✅ Minimal computational requirements
5. ❌ Cannot capture word order or context
6. ❌ Limited by linear decision boundaries
7. ❌ Struggles with complex linguistic phenomena

**LSTM**:

1. ✅ Captures sequential patterns and context
2. ✅ Handles negation and complex expressions
3. ✅ Better performance on longer reviews
4. ✅ Learns hierarchical features automatically
5. ❌ Requires more training data
6. ❌ Higher computational cost
7. ❌ Less interpretable
8. ❌ Prone to overfitting without proper regularization

## 8.3 Practical Considerations

Our findings suggest several practical considerations for implementing sentiment analysis systems:

1. **Resource constraints**: If computational resources are limited, Logistic Regression provides strong performance at minimal cost
2. **Data volume**: With larger datasets, LSTM's advantage increases; with smaller datasets, traditional ML may be more appropriate
3. **Application complexity**: For applications requiring nuanced understanding of complex linguistic phenomena, deep learning approaches are preferred
4. **Interpretability needs**: When explanation of decisions is critical, traditional ML offers clearer insight into feature importance
5. **Review length**: LSTM shows greater advantage on longer texts where context spans greater distances

# 9. Conclusions and Future Work

## 9.1 Key Findings

1. Deep learning approaches outperform traditional ML for sentiment analysis, but the margin is modest (1-2%)
2. Text preprocessing significantly impacts model performance for both approaches
3. LSTM models excel at capturing sequential patterns and context in text
4. Dropout regularization is critical for preventing overfitting in deep learning models

5. Traditional ML approaches remain competitive and offer advantages in speed and interpretability

## 9.2 Recommendations

Based on our findings, we recommend:

1. **Hybrid approaches**: Consider ensemble methods combining traditional ML and deep learning for optimal results
2. **Tailored preprocessing**: Customize text cleaning and feature extraction based on the specific model architecture
3. **Model selection based on constraints**: Choose between traditional ML and deep learning based on available resources, data volume, and performance requirements

## 9.3 Future Research Directions

Several promising avenues for extending this work include:

1. **Transformer architectures**: Implement BERT, RoBERTa or other transformer-based models that have shown state-of-the-art performance on sentiment tasks
2. **Multi-class sentiment**: Extend beyond binary classification to predict fine-grained sentiment levels
3. **Cross-domain generalization**: Test how well models trained on movie reviews perform on other domains like product reviews or social media posts
4. **Explainable AI techniques**: Develop better methods to interpret deep learning model predictions for sentiment analysis
5. **Efficient deep learning**: Explore knowledge distillation and model compression to reduce computational requirements of deep learning approaches

# GitHub Repository

The complete code, including model training, preprocessing, evaluation, and visualization, is available at:
https://github.com/Ikirezil/Sentiment_analysis-Assignment.git

# 10. References

1. Hochreiter, S., & Schmidhuber, J. (1997). *Long Short-Term Memory*. Neural Computation, 9(8), 1735–1780. https://doi.org/10.1162/neco.1997.9.8.1735

2. Vaswani, A., et al. (2017). *Attention Is All You Need*. Advances in Neural Information Processing Systems. https://arxiv.org/abs/1706.03762

3. Pedregosa, F., et al. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, 2825–2830.

Chollet, F. (2015). *Keras*. GitHub repository. https://github.com/fchollet/keras

4. Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python*.     O'Reilly Media, Inc.

# 11. Appendix: Team Contributions

**Team Member Contributions**

| Team Member | Contributions |
|---|---|
| Dieudonne Kobobey Ngum | EDA, data preprocessing, model evaluation, report writing |
| Theodora Ngozichukwuka Omunizua | LSTM implementation, hyperparameter tuning, performance analysis |
| Josiane Ishimwe | Logistic regression modeling, documentation, result visualization |
| Inès Ikirezi | Experiment design, literature review, preprocessing pipeline |