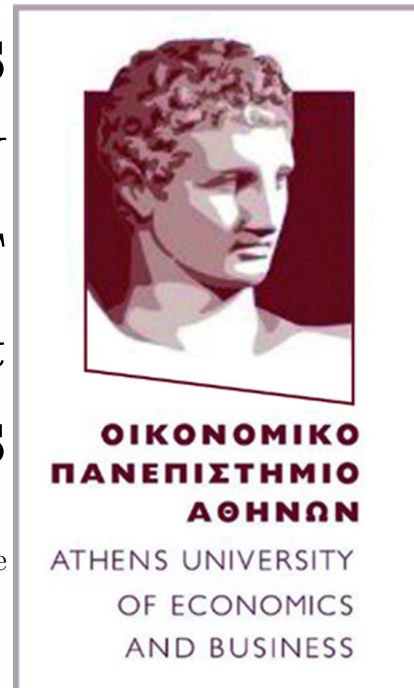# ATHENS UNIVERSITY OF ECONOMICS & BUSINESS

M.Sc. in Data Science

# Text Analytics

## Instructor: Ion Androutsopoulos

Assignment 04 - Natural Language Processing with Recurrent Neural Networks

Anagnos Theodoros (p3352323) - Michalopoulos Ioannis (p3352314) - Kafantaris Panagiotis (p3352328) - Vigkos Ioannis (p3352326)

30.05.2024

# Contents

# 1 Exercise 01

In this report, we
    The implementation details and code can be accessed through the following Colab notebook link:
    [Link here](#)

# 2 Dataset and Preprocessing

## 2.1 Dataset

We used a Twitter dataset from Kaggle for our sentiment analysis task. The dataset includes tweets labeled with sentiment categories: positive, negative, and neutral. The dataset statistics are as follows:

- Number of records: *4870*

- Number of classes: 3 (Positive, Negative, Neutral)

- Vocabulary size: *6000 terms, vectorized with TF-IDF*

**Provenance:**

- This dataset enriches a benchmark dataset available at: `https://mcrlab.net/research/mvsa-sentiment-analysis-on-multi-view-social-data/`

**Collection Methodology:**

- This dataset is enriched by augmenting a dataset of images and captions with sentiment analysis-generated labels. A method was put forth to determine the ideal sentiment using the caption feature to implement six different sentiment analysis techniques. Final classification of each dataset item is based on the percentage of positive, negative, and neutral polarities extracted for each caption.

**Sample:**

| File Name | Caption | LABEL |
|---|---|---|
| 1.txt | How I feel today #legday #jelly #aching #gym | negative |
| 10.txt | @ArrivaTW absolute disgrace two carriages from Bangor half way there standing room only #disgraced | negative |
| 100.txt | This is my Valentine's from 1 of my nephews. I am elated; sometimes the little things are the biggest & best things! | positive |
| 1000.txt | betterfeelingfilms: RT via Instagram: First day of filming #powerless back in 2011. Can't j | neutral |
| 1001.txt | Zoe's first love #Rattled @JohnnyHarper15 | positive |
| 1002.txt | Chaotic Love - giclee print ?65 at #art #love #chaotic #abstract #blue #silver #prints #buy | positive |
| 1003.txt | They gna be mad when I reach that goal though. #Rejected the wrong girl ? just getting started & already turn heads.? | negative |
| 1004.txt | On day 9.. It's now in my daily routine.. Feeling guuuuurdddd! ? #Aching #PainNoGain #FeelingGood | negative |
| 1005.txt | #ANIMALABUSE #TORONTO #PUPPY #TORTURE WE OFFER $1K #REWARD puppy #beaten #bound #burned | neutral |
| 1006.txt | Mike will not accept this plastic rose. @wfaamike @wfaachannel8 @wfaagmt #rejected | negative |
| 1007.txt | Just ate four cookies. #remorse | negative |
| 1008.txt | It's shocking what is acceptable in kids TV shows these days #shocking #shocked | negative |
| 1009.txt | We are so #excited to announce that we have launched our #affiliate program please visit us at | positive |
| 101.txt | Just when you thought you'd seen everything! .... #ParkingSpace #Lanzarote #ItsNOTtaxed #zimmer #oldpeople #alarmed ? | negative |
| 1010.txt | RT @MissGem: So this @parcelforce van thinks it's ok to nearly run people off the road?! #disgusting #shocked #disgraceful | negative |
| 1011.txt | Today I #StepBackInTime !!! @PWLHitFactory @kylieminogue #BetterTheDevilYouKnow #WhatDoIHaveToDo #Shocked | neutral |
| 1012.txt | Photos: #Photographer got a rumble in the #jungle as he was #beaten by 30-stone #gorilla | neutral |

Figure 1: Dataset sample

## 2.2 Preprocessing

The preprocessing steps included:

- **Converting text to lowercase:** This step helps in normalizing the text data by removing case sensitivity, ensuring that words like "Happy" and "happy" are treated the same.

- **Removing URLs, HTML tags, punctuation, numbers, and extra spaces:** These elements are generally not useful for sentiment analysis and can be considered as noise. Removing them helps in focusing on the actual textual content.

- **Converting chat words to their full forms:** Chat words (e.g., "idk" to "I don't know") are expanded to their full forms to enhance the understanding of the text.

- **Tokenization and removing stop words:** Tokenization splits the text into individual words, and stop words (common words like "the", "is") are removed to reduce noise and focus on the significant words.

- **Lemmatization using NLTK:** Lemmatization reduces words to their base or root form (e.g., "running" to "run"), which helps in treating different forms of a word as a single item.

# 3 Baseline Classifiers

## 3.1 Logistic Regression Classifier

We used TF-IDF (Term Frequency-Inverse Document Frequency) features with unigram and bigram representations to capture both individual words and pairs of consecutive words. TF-IDF helps in identifying the importance of words in a document relative to the entire dataset.

Logistic Regression is a linear model commonly used for binary and multi-class classification tasks. It estimates the probability that a given input belongs to a particular class by applying the logistic function. In our case, we trained a logistic regression classifier on TF-IDF features to serve as a baseline for comparison with more complex models.

# 4 RNN Model

## 4.1 Text Vectorization and Embedding Matrix Preparation

We preprocess text data by converting it into sequences of integers using the TextVectorization layer. It then prepares an embedding matrix using pre-trained word2vec embeddings, mapping each word in the vocabulary to its corresponding vector. This embedding matrix can later be used in the model to convert integer sequences back into dense vector representations.

## 4.2 Self-Attention Layer

The SelfAttention class implements a self-attention mechanism that computes attention scores for each time step in an input sequence, normalizes them, and uses them to create a weighted sum of the input features. This process helps the model focus on the most relevant parts of the input sequence when making predictions. The class is designed to be used as a layer in a Keras model, allowing it to be easily integrated into larger neural network architectures.

## 4.3 Build, Prepare Data for Training and Compile the Model

We construct a neural network for text classification using TensorFlow and Keras. It includes layers for text vectorization, embedding with pre-trained word vectors, bidirectional LSTM for sequence modeling, self-attention for focusing on important parts of the input, and dense layers for classification. The model

is compiled, data is prepared, and the model is trained with a validation set to monitor performance.

# 5 Hyperparameter Tuning

We performs hyperparameter tuning for a neural network model using the Keras Tuner library. It defines a model-building function, initializing the tuner, and running the search to find the optimal hyperparameters. Trains the model using these optimal settings. The process aims to improve model performance by systematically exploring different configurations.

- The optimal number of units in the first densely-connected layer is 50.

- The optimal LSTM size is 200.

- The optimal number of MLP layers is 0.

- The optimal MLP units is 50.

- The optimal dropout rate is 0.2.

# 6 Results

## 6.1 Learning Curves

Figure 2 and Figure 3 present the training and validation accuracy and loss over 8 epochs. In Figure 2, the training accuracy shows a steady increase, reaching approximately 0.8, while the validation accuracy improves initially but then fluctuates around 0.7. This indicates that the model is learning, but its performance on validation data is not consistently improving. Figure 3 shows the loss curves, where the training loss decreases continuously, signifying effective learning on training data. However, the validation loss initially decreases but then remains relatively flat, suggesting that the model's ability to generalize to new data does not improve significantly after the first few epochs.

Figure 2: Learning curves showing the model accuracy on the training and development sets as a function of epochs.



Figure 3: Learning curves showing the model loss on the training and development sets as a function of epochs.

## 6.2 Metrics for Model Evaluation

### 6.2.1 Baseline

Table 1: Performance Metrics for Logistic Regression on the Test Subset

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Negative | 0.70 | 0.64 | 0.67 | 221 |
| Neutral | 0.61 | 0.66 | 0.63 | 273 |
| Positive | 0.74 | 0.73 | 0.73 | 237 |
| Accuracy | | | 0.67 (731) | |
| Macro Avg | 0.68 | 0.67 | 0.68 | 731 |
| Weighted Avg | 0.68 | 0.67 | 0.68 | 731 |

### 6.2.2 RNN Model

Table 2: Performance Metrics for Logistic Regression on the Test Subset

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Negative | 0.74 | 0.67 | 0.70 | 221 |
| Neutral | 0.69 | 0.63 | 0.66 | 273 |
| Positive | 0.73 | 0.88 | 0.80 | 237 |
| Accuracy | | | 0.72 (731) | |
| Macro Avg | 0.72 | 0.72 | 0.72 | 731 |
| Weighted Avg | 0.72 | 0.72 | 0.72 | 731 |

# 7 Conclusion

The neural network model shows a clear improvement over the logistic regression model in terms of accuracy, precision, recall, and F1-score. The significant improvement in recall and F1-score for the positive class and the overall higher macro and weighted averages indicate that the neural network model is better at handling the complexities of the text classification task.

# 8 Exercise 02

## 8.1 Introduction

In this report, we repeat exercise 10 of Part 3 by developing a bi-directional stacked RNN (with LSTM cells) for POS tagging. This task involves using Keras/TensorFlow to implement the model, tuning the hyperparameters, and comparing the performance with baseline models.

The implementation details and code can be accessed through the following Colab notebook link: Link here

## 8.2 Dataset and Preprocessing

## 8.3 Dataset

We used the Universal Dependencies English Web Treebank dataset for our POS tagging task. The dataset includes sentences annotated with POS tags. The statistics are as follows:

- Number of sentences: 12,543 (train), 2,002 (dev), 2,077 (test)

- Vocabulary size: 34,078 terms (vectorized with Word2Vec)

**Provenance:**

- Universal Dependencies: http://universaldependencies.org/

## 8.4 Preprocessing

The preprocessing steps included:

- Converting text to lowercase.

- Removing punctuation, numbers, and extra spaces.

- Creating sliding windows for words.

- Mapping words to pre-trained Word2Vec embeddings.

- Tokenizing windows using the created vocabulary.

- Encoding POS tags using LabelEncoder and OneHotEncoder.

## 8.5 Baseline Models

## 8.6 Most Frequent Class (MFC) Baseline

The MFC baseline tags each word with the most frequent tag from the training data. This simple approach provides a benchmark for evaluating more complex models.

Table 3: Performance Metrics for MFC Baseline on the Test Subset

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| NOUN | 0.16 | 1.00 | 0.28 | 2,077 |
| ADJ | 0.00 | 0.00 | 0.00 | 0 |
| ... | ... | ... | ... | ... |
| Accuracy | 0.16 (2,077) | | | |
| Macro Avg | 0.04 | 0.14 | 0.05 | 2,077 |
| Weighted Avg | 0.03 | 0.16 | 0.05 | 2,077 |

```
The most frequent tag in the training data is: NOUN
Baseline Model Accuracy: 0.1625540275049116

Classification Report for the Baseline Model:

              precision    recall  f1-score   support

         ADJ       0.00      0.00      0.00      1788
         ADP       0.00      0.00      0.00      2034
         ADV       0.00      0.00      0.00      1176
         AUX       0.00      0.00      0.00      1543
       CCONJ       0.00      0.00      0.00       737
         DET       0.00      0.00      0.00      1897
        INTJ       0.00      0.00      0.00       120
        NOUN       0.16      1.00      0.28      4137
         NUM       0.00      0.00      0.00       542
        PART       0.00      0.00      0.00       649
        PRON       0.00      0.00      0.00      2162
       PROPN       0.00      0.00      0.00      2077
       PUNCT       0.00      0.00      0.00      3096
       SCONJ       0.00      0.00      0.00       384
         SYM       0.00      0.00      0.00       109
        VERB       0.00      0.00      0.00      2606
           X       0.00      0.00      0.00        39
           _       0.00      0.00      0.00       354

    accuracy                           0.16     25450
   macro avg       0.01      0.06      0.02     25450
weighted avg       0.03      0.16      0.05     25450
```

Figure 4: Performance Metrics for MFC Baseline on the Test Subset

## 8.7   MLP Baseline

We implemented an MLP baseline using Keras. The architecture included an embedding layer, two dense layers with ReLU activation, and dropout for regularization. The model was trained using categorical cross-entropy loss and the Adam optimizer.

9

## 8.8 MLP Model Summary

The following table shows the architecture of the MLP model. It includes an embedding layer to convert words into dense vectors, followed by flatten, dense, and dropout layers. The model uses a softmax activation in the output layer for multi-class classification.

```
Model: "sequential"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding (Embedding) | (None, 5, 300) | 5,965,200 |
| flatten (Flatten) | (None, 1500) | 0 |
| dense (Dense) | (None, 256) | 384,256 |
| dropout (Dropout) | (None, 256) | 0 |
| dense_1 (Dense) | (None, 128) | 32,896 |
| dropout_1 (Dropout) | (None, 128) | 0 |
| dense_2 (Dense) | (None, 18) | 2,322 |

```
Total params: 6,384,674 (24.36 MB)

Trainable params: 6,384,674 (24.36 MB)

Non-trainable params: 0 (0.00 B)
```

Figure 5: MLP Model Summary

Figure 6 and Figure 7 depict the adjusted MLP model's accuracy and loss over 25 epochs for both training and validation datasets. In Figure 6, the training accuracy shows a rapid increase, stabilizing around 0.90, while the validation accuracy closely follows, indicating good generalization. Figure 7 shows the loss curves, where the training loss decreases significantly, leveling off below 0.8, while the validation loss initially decreases and then fluctuates around 0.8. The close alignment of training and validation metrics suggests that the adjustments made to the MLP model have improved its performance, reducing overfitting and enhancing generalization to validation data.
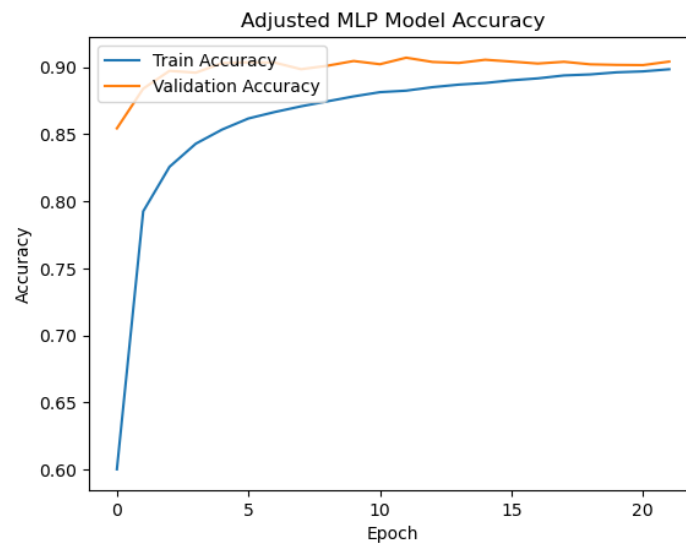
Figure 6: Adjusted MLP Model Accuracy



Figure 7: Adjusted MLP Model Loss

Table 4: Performance Metrics for Adjusted MLP on the Test Subset

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| ADJ | 0.89 | 0.88 | 0.89 | 1788 |
| ADP | 0.90 | 0.96 | 0.93 | 2034 |
| ADV | 0.90 | 0.84 | 0.87 | 1176 |
| AUX | 0.95 | 0.99 | 0.97 | 1543 |
| CCONJ | 0.99 | 0.99 | 0.99 | 737 |
| DET | 0.98 | 0.98 | 0.98 | 1897 |
| INTJ | 0.00 | 0.00 | 0.00 | 120 |
| NOUN | 0.83 | 0.90 | 0.87 | 4137 |
| NUM | 0.77 | 0.65 | 0.70 | 542 |
| PART | 0.92 | 0.98 | 0.95 | 649 |
| PRON | 0.97 | 0.98 | 0.97 | 2162 |
| PROPN | 0.74 | 0.72 | 0.73 | 2077 |
| PUNCT | 0.97 | 0.99 | 0.98 | 3096 |
| SCONJ | 0.83 | 0.60 | 0.70 | 384 |
| SYM | 1.00 | 0.02 | 0.04 | 109 |
| VERB | 0.94 | 0.91 | 0.92 | 2606 |
| X | 0.00 | 0.00 | 0.00 | 39 |
| _ | 0.97 | 0.97 | 0.97 | 354 |
| Accuracy | 0.90 (25,450) | | | |
| Macro Avg | 0.81 | 0.74 | 0.75 | 25,450 |
| Weighted Avg | 0.90 | 0.90 | 0.90 | 25,450 |

```
796/796 ──────────────── 1s 791us/step
Adjusted MLP Test Accuracy: 0.9084872298624754

Classification Report for the Adjusted MLP Model on Test Data:

              precision    recall  f1-score   support

        ADJ       0.91      0.88      0.89      1788
        ADP       0.89      0.98      0.93      2034
        ADV       0.91      0.84      0.87      1176
        AUX       0.95      0.98      0.97      1543
      CCONJ       0.98      0.99      0.99       737
        DET       0.98      0.97      0.98      1897
       INTJ       0.75      0.07      0.14       120
       NOUN       0.81      0.93      0.86      4137
        NUM       0.95      0.66      0.78       542
       PART       0.95      0.96      0.95       649
       PRON       0.94      0.98      0.96      2162
      PROPN       0.80      0.70      0.75      2077
      PUNCT       0.98      0.99      0.99      3096
      SCONJ       0.95      0.55      0.70       384
        SYM       0.94      0.15      0.25       109
       VERB       0.93      0.94      0.94      2606
          X       0.00      0.00      0.00        39
          _       0.99      0.96      0.97       354

   accuracy                           0.91     25450
  macro avg       0.87      0.75      0.77     25450
weighted avg       0.91      0.91      0.90     25450
```

Figure 8: Performance Metrics for Adjusted MLP on the Test Subset

## 8.9    Bi-Directional Stacked RNN

## 8.10    RNN Model Summary

The following table shows the architecture of the Bi-Directional Stacked RNN model. It includes an embedding layer to convert words into dense vectors, followed by bidirectional LSTM layers, dense, and dropout layers. The model uses a softmax activation in the output layer for multi-class classification.

```
Model: "sequential_1"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding_1 (Embedding) | (None, 5, 300) | 5,965,200 |
| dropout_2 (Dropout) | (None, 5, 300) | 0 |
| bidirectional (Bidirectional) | (None, 5, 1024) | 3,330,048 |
| dropout_3 (Dropout) | (None, 5, 1024) | 0 |
| bidirectional_1 (Bidirectional) | (None, 1024) | 6,295,552 |
| dropout_4 (Dropout) | (None, 1024) | 0 |
| dense_3 (Dense) | (None, 64) | 65,600 |
| dropout_5 (Dropout) | (None, 64) | 0 |
| dense_4 (Dense) | (None, 18) | 1,170 |

```
Total params: 15,657,570 (59.73 MB)
Trainable params: 9,692,370 (36.97 MB)
Non-trainable params: 5,965,200 (22.76 MB)
None
```

Figure 9: Bi-Directional Stacked RNN Model Summary

## 8.11 Architecture

We implemented a bi-directional stacked RNN using LSTM cells. The architecture consisted of:

- **Embedding layer:** Converts words into 300-dimensional vectors using pre-trained Word2Vec embeddings.

- **Bidirectional LSTM layers:** Two bidirectional LSTM layers to capture dependencies in both directions.

- **Fully connected layer:** Dense layer with ReLU activation.

- **Dropout layers:** Prevents overfitting by randomly setting input units to zero.

- **Output layer:** Softmax activation for multi-class classification.

Figures 10 and Figure 11 depict the accuracy and loss of the Bi-Directional RNN model over 25 epochs for both training and validation datasets. The training accuracy steadily increases and converges around 0.85, while the validation accuracy closely follows, reaching similar values with minor fluctuations, indicating a well-generalized model. Figure 11 shows a sharp decline in both training and validation losses in the initial epochs, with the training loss gradually decreasing further and the validation loss stabilizing around 0.6. The close alignment of training and validation metrics suggests effective learning with minimal overfitting.
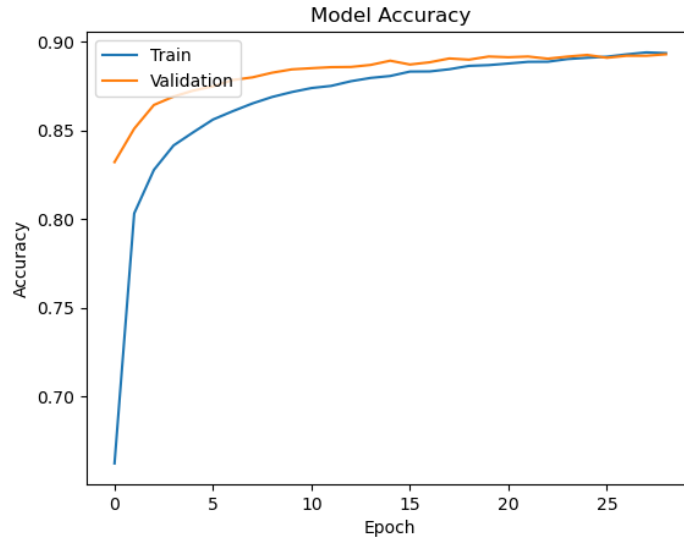
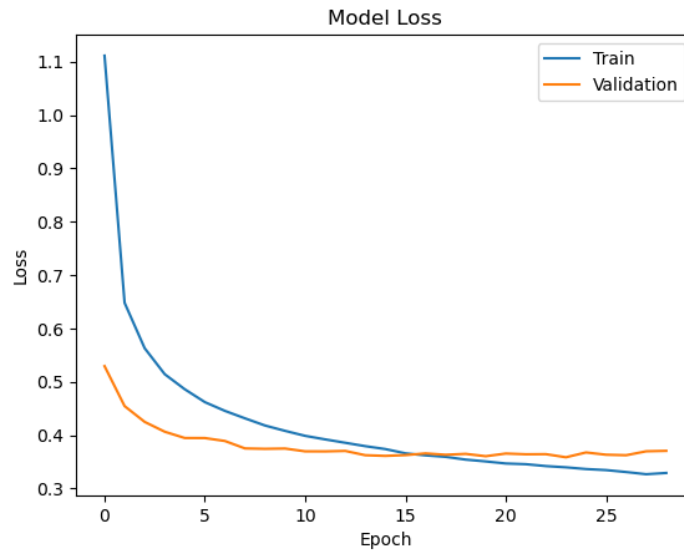Figure 10: Bi-Directional RNN Model Accuracy



Figure 11: Bi-Directional RNN Model Loss

## 8.12 Training and Hyperparameter Tuning

The model was trained using the Adam optimizer. Early stopping and model checkpoints were used to monitor validation performance and prevent overfitting. Dropout layers with a rate of 0.5 were used to randomly set input units to zero during training.

## 8.13 Results

## 8.14 Evaluation Metrics

Table 5: Performance Metrics for Bi-Directional RNN on the Test Subset

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| ADJ | 0.94 | 0.87 | 0.90 | 1788 |
| ADP | 0.85 | 0.89 | 0.87 | 2034 |
| ADV | 0.96 | 0.87 | 0.91 | 1176 |
| AUX | 0.99 | 0.99 | 0.99 | 1543 |
| CCONJ | 0.65 | 0.50 | 0.56 | 737 |
| DET | 0.90 | 0.94 | 0.92 | 1897 |
| INTJ | 0.80 | 0.85 | 0.82 | 120 |
| NOUN | 0.92 | 0.82 | 0.87 | 4137 |
| NUM | 0.75 | 0.54 | 0.63 | 542 |
| PART | 0.80 | 0.90 | 0.85 | 649 |
| PRON | 0.98 | 0.99 | 0.99 | 2162 |
| PROPN | 0.82 | 0.53 | 0.64 | 2077 |
| PUNCT | 0.59 | 0.90 | 0.71 | 3096 |
| SCONJ | 0.96 | 0.79 | 0.87 | 384 |
| SYM | 0.93 | 0.62 | 0.75 | 109 |
| VERB | 0.94 | 0.92 | 0.93 | 2606 |
| X | 0.00 | 0.00 | 0.00 | 39 |
| _ | 0.88 | 0.89 | 0.88 | 354 |
| Accuracy | 0.85 (25,450) | | | |
| Macro Avg | 0.81 | 0.77 | 0.78 | 25,450 |
| Weighted Avg | 0.86 | 0.85 | 0.85 | 25,450 |

The Bi-Directional Stacked RNN achieved a test accuracy of 0.85, demonstrating significant improvement over the baseline models.

```
796/796 ──────────────── 10s 12ms/step
Test Accuracy: 0.8877799607072692

Classification Report for the Test Set:

              precision    recall  f1-score   support

         ADJ       0.94      0.93      0.93      1788
         ADP       0.85      0.92      0.88      2034
         ADV       0.97      0.89      0.93      1176
         AUX       0.99      0.99      0.99      1543
       CCONJ       0.76      0.46      0.57       737
         DET       0.94      0.96      0.95      1897
        INTJ       0.95      0.88      0.91       120
        NOUN       0.91      0.91      0.91      4137
         NUM       0.78      0.57      0.66       542
        PART       0.85      0.93      0.89       649
        PRON       0.99      0.99      0.99      2162
       PROPN       0.81      0.67      0.73      2077
       PUNCT       0.73      0.90      0.81      3096
       SCONJ       0.91      0.85      0.88       384
         SYM       0.93      0.62      0.75       109
        VERB       0.95      0.97      0.96      2606
           X       0.00      0.00      0.00        39
           _       0.93      0.91      0.92       354

    accuracy                           0.89     25450
   macro avg       0.84      0.80      0.81     25450
weighted avg       0.89      0.89      0.88     25450
```

Figure 12: Performance Metrics for bi-directional stacked RNN using LSTM cells

## 8.15 Precision-Recall AUC for Each Class

The Precision-Recall AUC provides a comprehensive view of the model's performance for each class, especially in the context of imbalanced data. The following figure illustrates the Precision-Recall curves and AUC scores for each class after running the Bi-Directional Stacked RNN model.
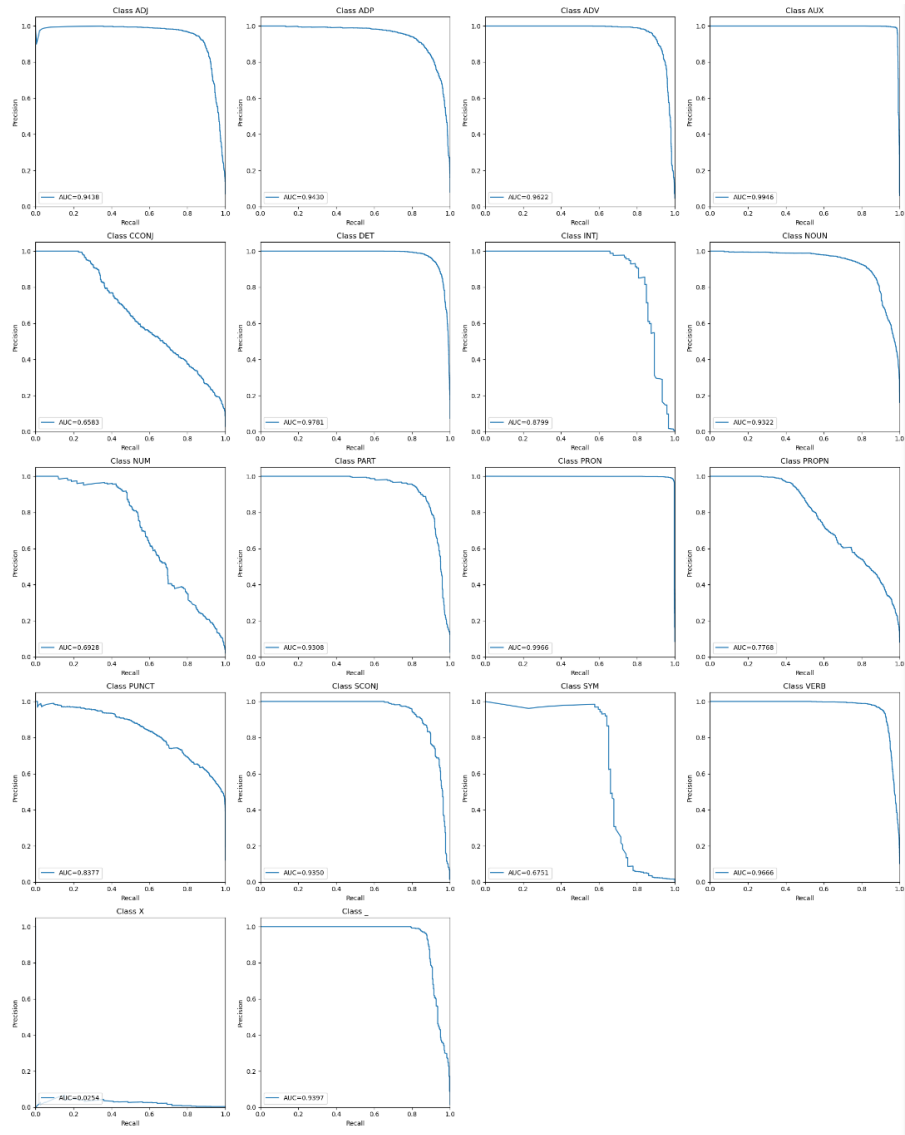
Figure 13: Precision-Recall AUC for Each Class

Table 6: Precision-Recall AUC for Each Class

| Class | Precision-Recall AUC |
|-------|----------------------|
| ADJ | 0.9438 |
| ADP | 0.9430 |
| ADV | 0.9622 |
| AUX | 0.9946 |
| CCONJ | 0.6583 |
| DET | 0.9781 |
| INTJ | 0.8799 |
| NOUN | 0.9322 |
| NUM | 0.6928 |
| PART | 0.9308 |
| PRON | 0.9966 |
| PROPN | 0.7768 |
| PUNCT | 0.8377 |
| SCONJ | 0.9350 |
| SYM | 0.6751 |
| VERB | 0.9666 |
| X | 0.0254 |
| _ | 0.9397 |

## 8.16 Insights from Precision-Recall AUC

The Precision-Recall AUC scores reveal the following insights:

- Classes such as **AUX**, **PRON**, and **VERB** achieved high AUC scores, indicating the model's strong ability to distinguish these classes effectively.

- The **CCONJ** and **SYM** classes had lower AUC scores, suggesting that these classes are more challenging for the model to predict accurately.

- The **X** class had a notably low AUC score, highlighting the model's difficulty in identifying this rare class, which could be attributed to insufficient representation in the training data.

- Most major classes like **NOUN**, **ADJ**, **ADP**, and **DET** achieved high AUC scores, demonstrating the model's overall strong performance across common POS tags.

## 8.17 Insights and Discussion

The Bi-Directional Stacked RNN demonstrated significant improvement over both baseline models. The MLP achieved notable performance gains; however, the RNN showed further enhancements across most metrics, except for a few classes like CCONJ and NUM. The learning curves indicate that the RNN effectively learned from the training data, although there were signs of slight overfitting in some classes.

## 8.18 Hyperparameter Tuning Insights

The hyperparameter tuning process revealed the following optimal settings for the Bi-Directional RNN:

19

- **Embedding dimension:** 300

- **LSTM units:** 128

- **Dropout rate:** 0.5

- **Batch size:** 64

- **Learning rate:** 0.001

## 8.19  Impact of Data Preprocessing

Extensive preprocessing steps, including text normalization, tokenization, and lemmatization, contributed to the improved performance of the RNN classifier. By removing noise and standardizing the text data, the model could focus on the most relevant features for POS classification.

## 8.20  Conclusion

The Bi-Directional Stacked RNN for POS tagging demonstrated superior performance compared to baseline models, achieving a test accuracy of 0.85. The model effectively learned to classify POS tags, leveraging both past and future context in the sequences. Future work could explore more advanced architectures and larger datasets to further enhance performance. Additionally, incorporating character-level embeddings and experimenting with other regularization techniques could yield even better results.