



# DEEP LEARNING

## CSL 4020

Assignment **3** Report

---

### FFN and RNN for sentiment analysis

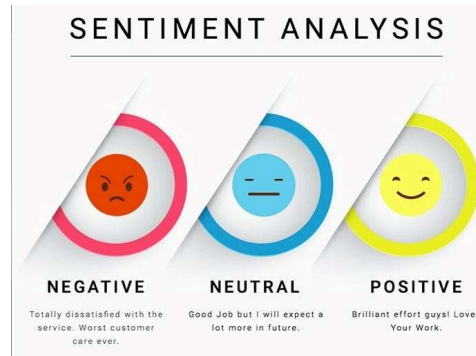
By - PRANAV PANT  
B21CS088



# FFN and RNN for Sentiment Analysis

## Introduction:

**Sentiment analysis** is a fundamental task in natural language processing (NLP) that focuses on identifying and categorizing sentiments expressed in textual data. It plays a vital role in applications such as opinion mining, customer feedback analysis, and social media monitoring.



The goal of this assignment was to implement two DL models, a Feed-Forward Neural Network (FFNN) and a Recurrent Neural Network (RNN), for binary and multi-class sentiment analysis tasks. These models were applied to understand how well they can classify sentiments based on textual inputs.

## Dataset details:

**1. IMDB Dataset:** The IMDB dataset contains movie reviews labeled as either **positive or negative**, making it suitable for **binary sentiment classification**. This dataset is widely used due to its large volume of text data, providing a variety of writing styles, opinions, and contexts.

- The dataset includes both testing and training set, total **50k** sentences with their labels, **25k** (12k positive, 12k negative) for training and testing each.
- Each review varies in length, making it challenging for models to handle both short and long text inputs.

**2. SemEval-2013 Dataset:** The SemEval-2013 Twitter dataset is used for **multi-class sentiment classification**, where tweets are categorized as **positive, negative, or neutral**. Tweets are typically short and contain informal language, adding complexity to the classification task.

- The dataset consists of short tweets, often containing slang, abbreviations, and emoticons.
- It is designed for multi-class classification, with tweets labeled into three sentiment classes: positive, negative, and neutral. The distribution of data is **imbalanced**. The dataset has around 10k samples in training set, 1.5k in validation and 3.5k in the test set.

## Preprocessing Methodology:

the preprocessing steps applied to the IMDB and SemEval datasets to prepare the textual data for sentiment analysis using both a Feed-Forward Neural Network (FFNN) and a Recurrent Neural Network (RNN).

### 1. Text Cleaning:

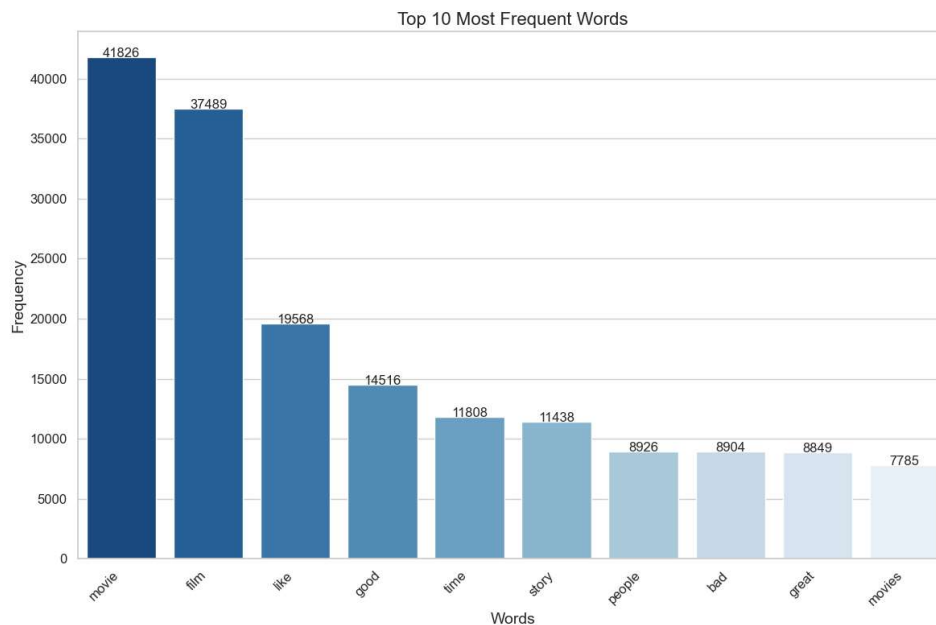
- Each review was cleaned by removing **HTML tags** and **non-alphabetic characters** using regular expressions.
- All text was converted to lowercase to ensure uniformity.

### 2. Tokenization:

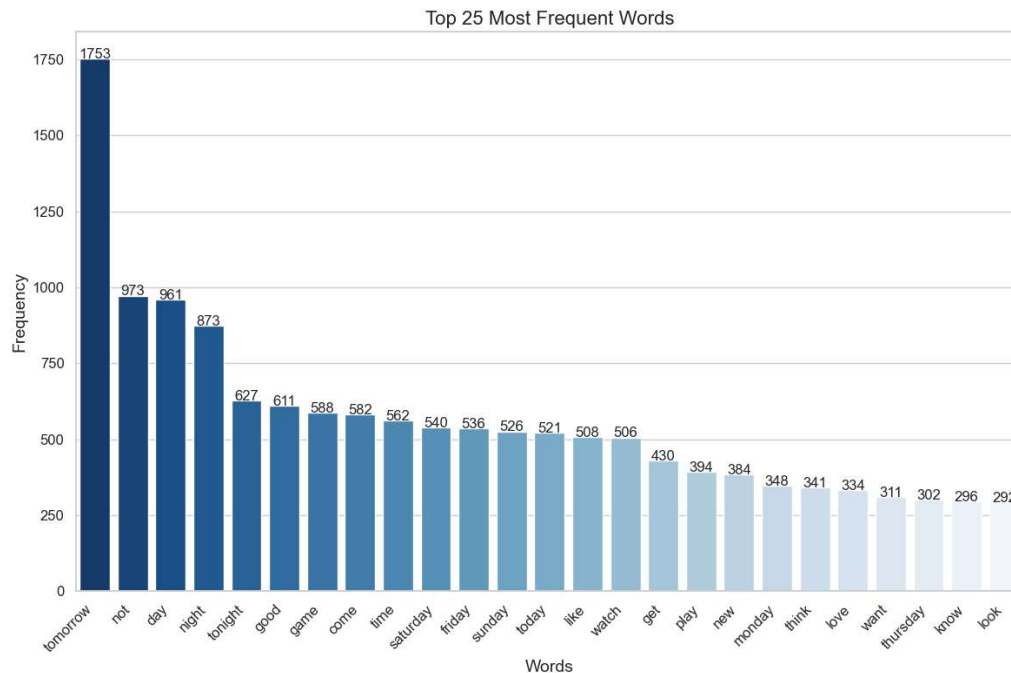
- Tokenization and lemmatization were done using the spaCy library (en\_core\_web\_sm model). During this step:
  - Stop words (common words like "the" or "and") and punctuation were removed.
  - Only alphabetic characters were retained for further analysis.

### 3. Building Vocabulary:

- After tokenization, a vocabulary was built based on word frequency, with *a minimum frequency threshold of 5 words*. Words that occurred fewer times were discarded to reduce noise.
- After testing for bigger vocabulary sizes (the kernel crashes for big vocabulary size), the vocabulary size was capped (**17,500** vocab size for **IMDB** dataset) to get unique tokens for efficient model training. (+2 for UNK and PAD tokens).
- Words that were **too short** (1-2 characters), the unwanted tokens were also removed.



Top 10 frequent words in **IMDB** dataset.



Top 25 frequent words in **Semeval dataset**.

- Special tokens like PAD (for padding shorter sentences) and UNK (for unknown words) were added to handle variable-length inputs and out-of-vocabulary words.

#### 4. **Padding and Truncation:**

- To ensure all reviews and tweets are of equal length, we calculated the **average sentence length** from the training data.

*Avg sentence length of IMDB dataset: **105 tokens***

*Avg sentence length of Semeval dataset: **11 tokens***

- Sentences longer than this average were truncated, and shorter sentences were padded with the **PAD** token to achieve the same length for all inputs.

#### 5. **Data Export:**

- Preprocessed and tokenized reviews were saved into separate **CSV files** for training and testing datasets, which contained both tokenized and padded versions of the data.

This preprocessing pipeline ensured that the text data was cleaned, tokenized, and padded to a standard format, ready for input into the deep learning models.

## Model Architectures:

## 1. Feed-Forward Neural Network (FFNN)

### Architecture Overview:

The proposed FFNN consists of two hidden layers, following the architecture outlined below:

### Input Layer:

The input layer takes in a flattened one-hot encoded sequence of words from the IMDB dataset.

The input size is  $105 \times 17502 = 1,837,710$ , where:

- 105 is the sequence length (padded or truncated)
- 17,502 is the vocabulary size

### Hidden Layer 1:

- **Size:** 256 neurons
- **Activation Function:** ReLU (Rectified Linear Unit)
- **Dropout:** 50% dropout applied after this layer to prevent overfitting.

### Hidden Layer 2:

- **Size:** 128 neurons
- **Activation Function:** ReLU
- **Dropout:** 50% dropout applied after this layer.

### Output Layer:

- **Size:** 1 neuron (for binary classification: positive vs. negative)
- **Activation Function:** Sigmoid activation function (for binary classification)

## 2. Recurrent Neural Network (RNN)

**Architecture Overview:** The RNN consists was tried and tested with different architectures, with the first layer as a linear layer followed by a single RNN layer. The other RNN variant was when we used the first layer as an embedding layer followed by the RNN layer. The architecture is given below:

### Embedding Layer:

- **Size:** Embedding dimension of 128
- This layer converts each word in the padded sequence (of max length 105) into a 128-dimensional embedding vector.
- In another variant, it was replaced by a simple linear layer of 256 dimensions.

### RNN Layer:

- **Type:** Standard RNN with tanh nonlinearity
- **Size:** 256 hidden units
- **Number of Layers:** 1 (no stacking of RNN layers)
- **Output:** The RNN processes the entire sequence of 105 words and outputs the hidden state at each time step.

### Last Hidden State:

After processing the entire sequence, the hidden state at the final time step is used as the representation for the entire sentence. This is passed to the next layer for classification.

#### Output Layer:

- **Size:** 2 neurons (for binary classification: positive vs. negative)
- **Activation Function:** Softmax or CrossEntropyLoss

#### Loss Function for Binary Classification (IMDB dataset):

For the Feed-Forward Neural Network (FFNN) and RNN we use Binary Cross-Entropy (BCE) loss since the task is **binary classification** (positive or negative sentiment).

The BCE loss is defined as:

$$\mathcal{L}_{BCE}(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

Where:

- $y_i$  is the true label (0 or 1).
- $\hat{y}_i$  is the predicted probability of the positive class (output from the Sigmoid function).
- $N$  is the total number of samples.

This loss function works by penalizing the model based on how far its predicted probabilities ( $\hat{y}_i$ ) are from the true binary labels ( $y_i$ ). When the prediction is close to the true label, the log term is small, minimizing the loss. If the prediction is far from the true label, the log term becomes large, increasing the loss.

#### Loss Function for Multiclass Classification (Semeval Dataset):

For the **SemEval dataset**, the task involves **multi-class classification** with three classes (positive, neutral, negative sentiments). The appropriate loss function for this setting is **Cross-Entropy Loss**, which works well for multi-class classification problems.

The **Cross-Entropy Loss** is defined as:

$$\mathcal{L}_{CE}(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log(\hat{y}_{i,c})$$

Where:

- $y_i(c)$  is the true label for class  $c$  (encoded as a one-hot vector).
- $\hat{y}_i(c)$  is the predicted probability of class  $c$  (output from the Softmax function).
- $N$  is the total number of samples.
- $C$  is the total number of classes (in this case, **C=3**).

In this case, for the SemEval dataset, the three possible classes are:

1. Positive (Class 1)
2. Neutral (Class 2)
3. Negative (Class 0)

The **Softmax function** ensures that the output probabilities sum to 1 across the three classes, distributing the probability mass across the possible sentiment categories.

Cross-Entropy Loss calculates how far the predicted probabilities ( $y^i, c$ ) are from the true one-hot encoded labels ( $y_i, c$ ) and penalizes incorrect predictions accordingly.

This loss is particularly effective for **multi-class classification** because it encourages the model to not only assign a high probability to the correct class but also to reduce the probabilities for the incorrect classes.

### Note:

Both architectures use non-linearities (ReLU in FFNN and tanh in RNN). Since in initial testing, the model did overfit the model so I included **dropout** and **early stopping** to reduce overfitting, while the choice of loss functions ensures that the models can handle binary classification effectively.

## Implementation:

A variety of ways were tried to get the best model for each of the datasets by varying the features of the datasets and the models.

### 1. Loading the Data:

#### 1.1 Custom Dataset Class

The IMDB dataset, which has been preprocessed and saved into CSV files. A custom SentimentDataset class is created, inheriting from torch.utils.data.Dataset. This class:

- Loads the padded reviews and corresponding labels from the dataset.
- Clamps word indices to ensure they do not exceed the specified vocabulary size (17,502).
- Converts each word index into a **one-hot encoded vector** for the FFNN model, and returns the processed review and label when accessed.

#### 1.2 DataLoader Creation

PyTorch DataLoader objects are created for batching and iterating through the training, validation (dev), and test datasets. Each DataLoader batches data into smaller chunks (batch size = 16) to make the training process more efficient. The number of batches is printed for each dataset.

---

## 2 Training the FFNN

The FFNN is trained using the **Adam optimizer** and **Binary Cross-Entropy Loss** (BCELoss()), which is suited for binary classification. The training process:

- Iterates over multiple epochs. (no of epoch = **10**).
- Computes the training loss and accuracy after every batch.
- At the end of each epoch, validation loss and accuracy are calculated.

### Early Stopping and Validation

To prevent overfitting, **early stopping** is implemented. If the validation loss does not improve for **3** consecutive epochs, training halts, and the best-performing model (based on validation loss) is restored. This ensures the model does not overfit to the training data.

### FFNN Evaluation and Plotting Results

Once the FFNN is trained, it is evaluated on the test set using metrics such as **accuracy**, **precision**, **recall**, and **F1-score**.

The model's performance is printed, including class-wise metrics for positive and negative sentiments. Training and validation losses, as well as accuracies, are plotted using matplotlib to visualize the model's progress over time.

---

## 3 Training the RNN

Similar to the FFNN, the RNN is trained using the Adam optimizer. The training process involves:

- Forward propagation through the embedding and RNN layers.
- Backward propagation using Cross-Entropy Loss.
- Updating weights with the Adam optimizer.

### Early Stopping and Validation

Early stopping is applied during RNN training as well, based on the validation accuracy. If the model does not show improvement for 5 consecutive epochs, training stops early, and the best model is restored.

### RNN Evaluation and Plotting Results

Once training completes, the RNN is evaluated on the test set, calculating metrics like **accuracy**, **precision**, **recall**, and **F1-score**. Class-wise performance for both sentiment categories is also provided. Finally, training and validation losses, along with accuracies, are plotted for the RNN to visualize how the model's performance evolved over the training process.



# Results

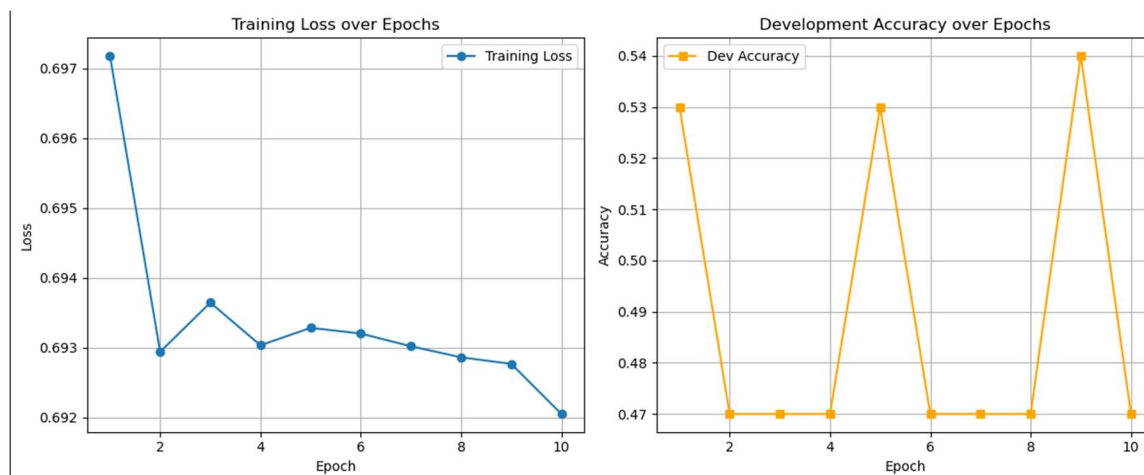
Here are the models with their configurations that were tried to get the best results.

## IMDB dataset:

-> Using subset of full vocabulary and front 100 words with one hot encoding.

The vocab size was taken as 10k (from 30k) and the input was truncated. Avg sentence length: 268 words reduced to top 100 words (to avoid crashing and memory error).

(i) Results for FFN model:



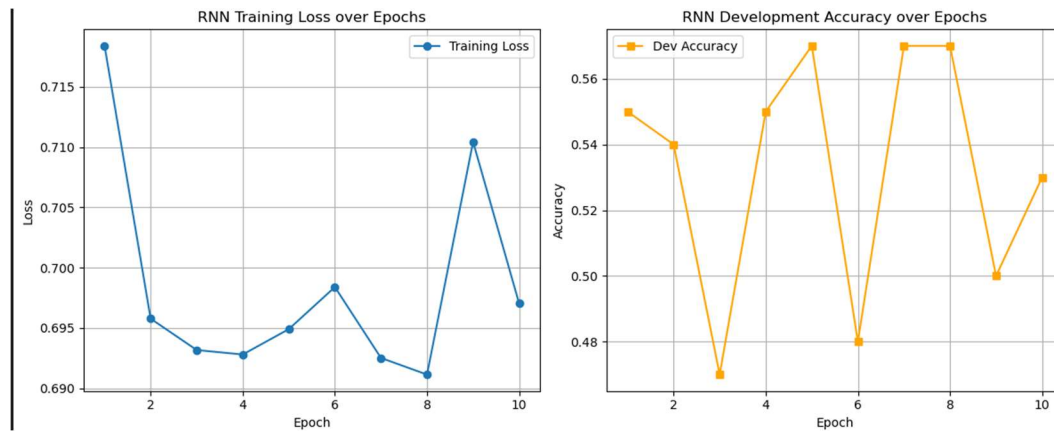
Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Accuracy: 0.5410
Test Precision: 0.5354
Test Recall: 0.6620
Test F1-Score: 0.5920

Per Class Metrics:
Class 0 - Precision: 0.5503, Recall: 0.4185, F1-Score: 0.4754
Class 1 - Precision: 0.5354, Recall: 0.6620, F1-Score: 0.5920
```

Evaluation metrics for the model

(ii) Results for RNN model:



Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```

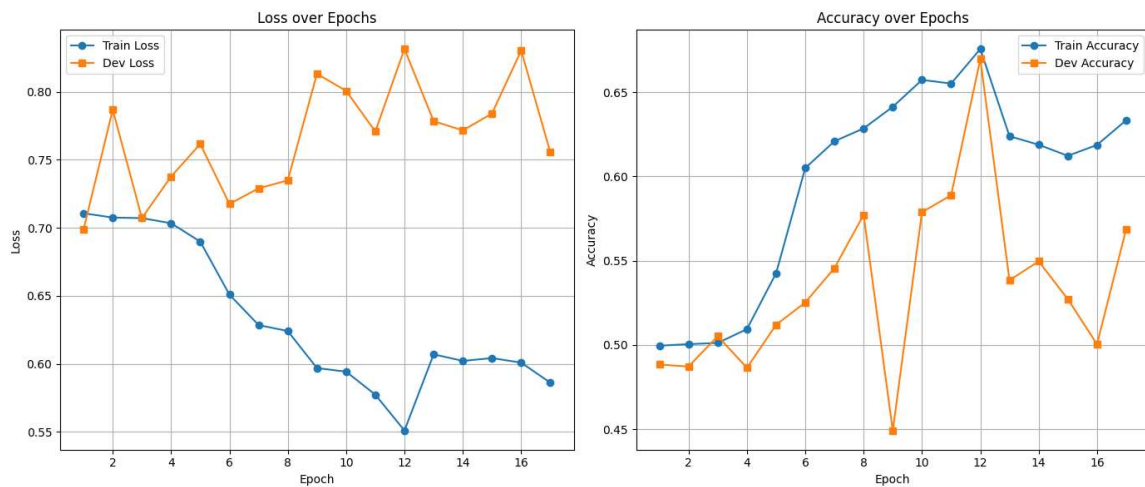
RNN Test Accuracy: 0.5380
RNN Test Precision: 0.5348
RNN Test Recall: 0.6262
RNN Test F1-Score: 0.5769

RNN Per Class Metrics:
Class 0 - Precision: 0.5426, Recall: 0.4487, F1-Score: 0.4912
Class 1 - Precision: 0.5348, Recall: 0.6262, F1-Score: 0.5769
  
```

Evaluation metrics for the model

This did not give good results as the accuracy was around 50% only, but it established a baseline for our model.

\* Tried a version with **Multi Hot encoding**, which reduced the processing time but gave similar results.



Test accuracy: 0.52

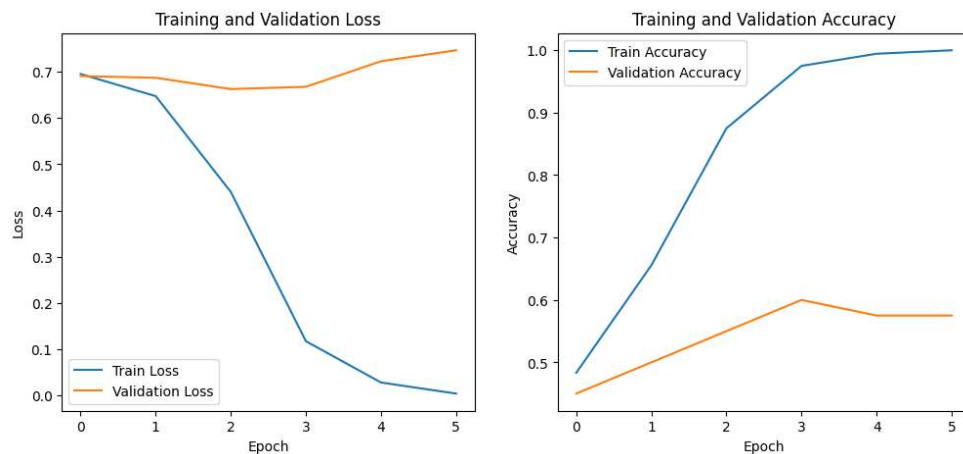
-> **For new IMDB dataset with vocabulary of top frequency words with one hot encoding.**

In this new dataset, the words for the vocabulary were taken to be the top 17500 words with the most frequency and accordingly the data was tokenized. This led to overall better results for the model.

vocab\_size = 17502.

Avg sentence length: 105

(i) **FFN model:**



Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Accuracy: 0.7348
Test Precision: 0.6927
Test Recall: 0.8440
Test F1-Score: 0.7609

Per Class Metrics:
Class 0 - Precision: 0.8004, Recall: 0.6256, F1-Score: 0.7023
Class 1 - Precision: 0.6927, Recall: 0.8440, F1-Score: 0.7609
```

Evaluation metrics for the model

(ii) **RNN model with embedding dim = 128.**

```

Test Loss: 0.7052
Test Accuracy: 0.4976
Test Precision: 0.4987
Test Recall: 0.9220
Test F1-Score: 0.6473

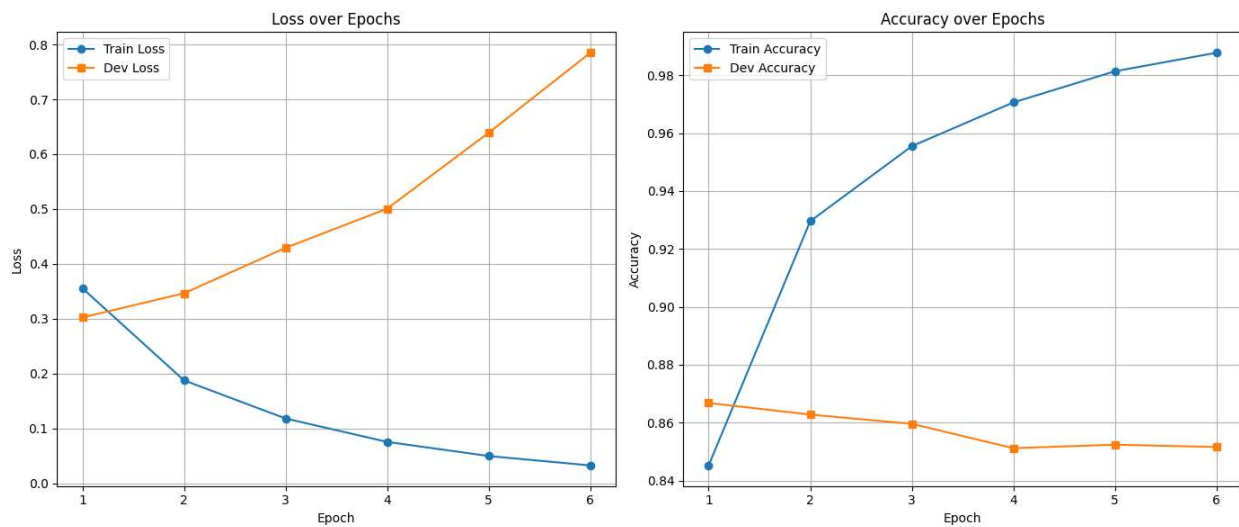
Per Class Metrics:
Class 0 - Precision: 0.4841, Recall: 0.0732, F1-Score: 0.1272
Class 1 - Precision: 0.4987, Recall: 0.9220, F1-Score: 0.6473

```

Evaluation metrics for the model

The RNN model with embedding dimensions did not give very good results. So, I replaced it with a linear layer.

### (iii) RNN model with linear layer



Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```

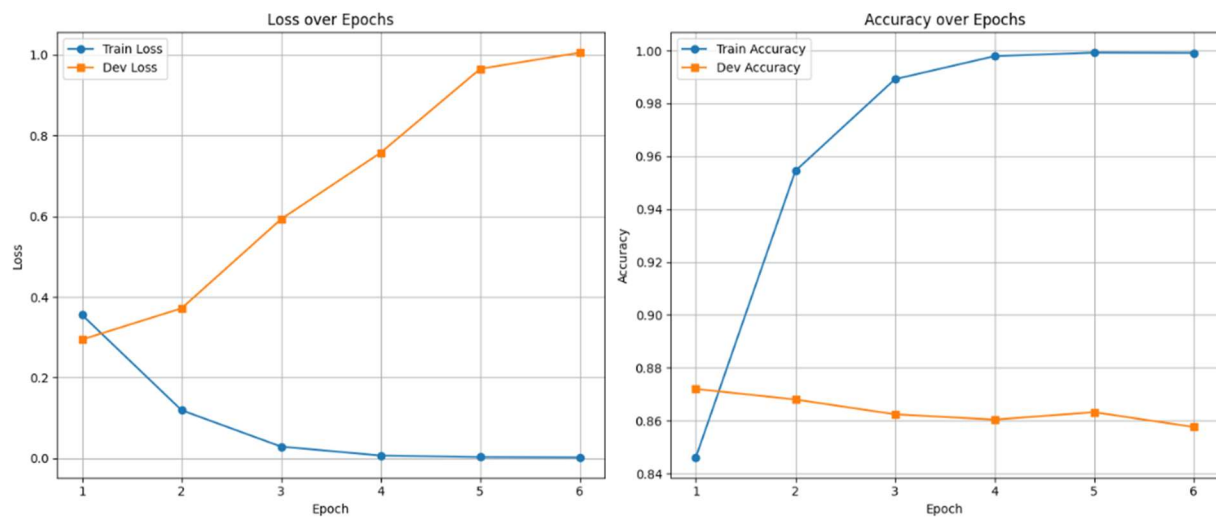
Test Loss: 0.3473
Test Accuracy: 0.8478
Test Precision: 0.8557
Test Recall: 0.8366
Test F1-Score: 0.8460

Per Class Metrics:
Class 0 - Precision: 0.8401, Recall: 0.8590, F1-Score: 0.8494
Class 1 - Precision: 0.8557, Recall: 0.8366, F1-Score: 0.8460

```

Evaluation metrics for the model

(iv) LSTM model with linear layer

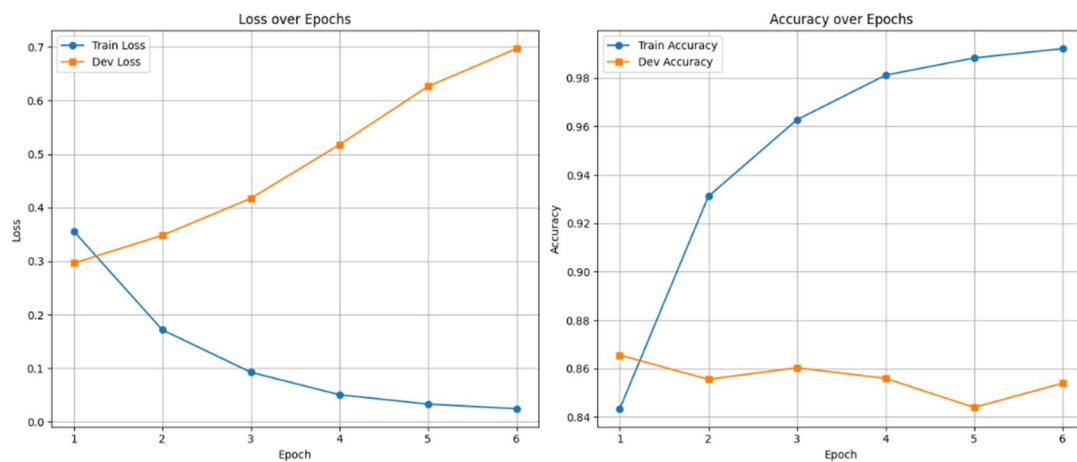


Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Loss: 0.3454
Test Accuracy: 0.8475
Test Precision: 0.8561
Test Recall: 0.8354
Test F1-Score: 0.8456
```

Evaluation metrics for the model

(v) GRU model with linear layer and multi-hot



## Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Loss: 0.3405
Test Accuracy: 0.8485
Test Precision: 0.8573
Test Recall: 0.8362
Test F1-Score: 0.8466

Per Class Metrics:
Class 0 - Precision: 0.8402, Recall: 0.8608, F1-Score: 0.8504
Class 1 - Precision: 0.8573, Recall: 0.8362, F1-Score: 0.8466
```

## Evaluation metrics for the model

### Example statement:

```
Statement : this was a good movie after all
Predicted label for the input sentence: Positive
```

We can see that the model predicts very well on the given statements.

## Conclusion:

We can see that using Linear Layer with RNN, GRU and LSTM's on the new dataset (with top words) gave the best accuracy of **84.85%** over **25k** samples test set. The model also achieved a high F1-score of **84.66%**, showing a good balance between precision and recall. This means the RNN's, GRU's and LSTM's were effective at understanding the relationships between words in the text, making it the best-performing model for this dataset.

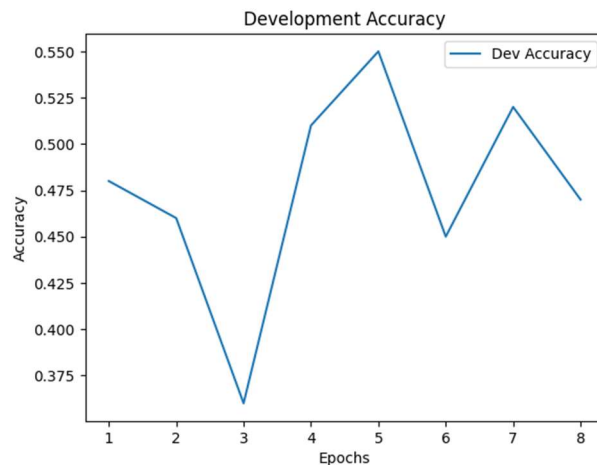
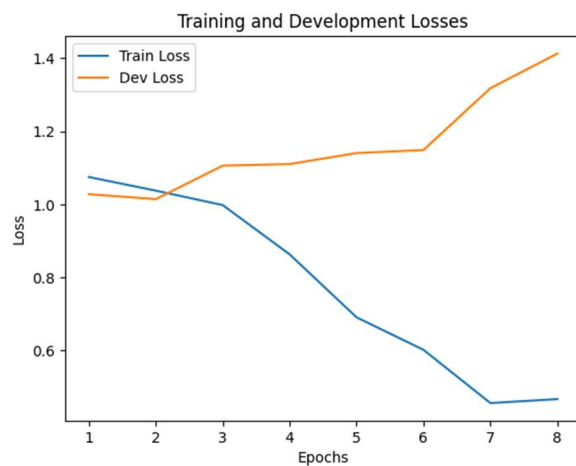
## Semeval dataset:

-> **Full vocabulary with one hot encoding.**

The vocab size for this dataset was 3083 and the input was truncated.

Avg sentence length of the tweets: 26.

**(i) Results for FFN model:**



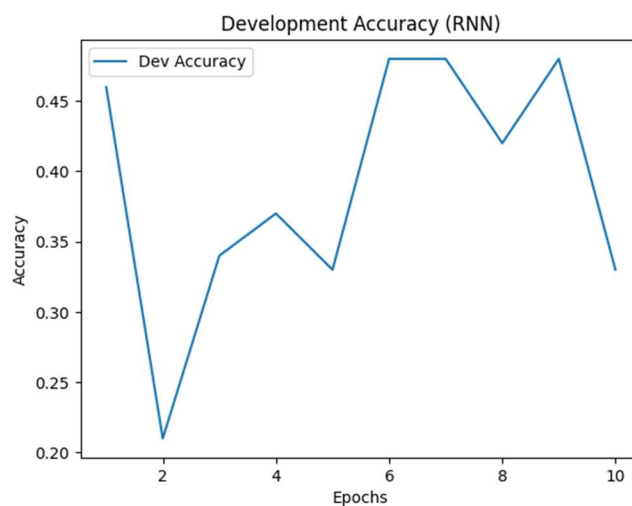
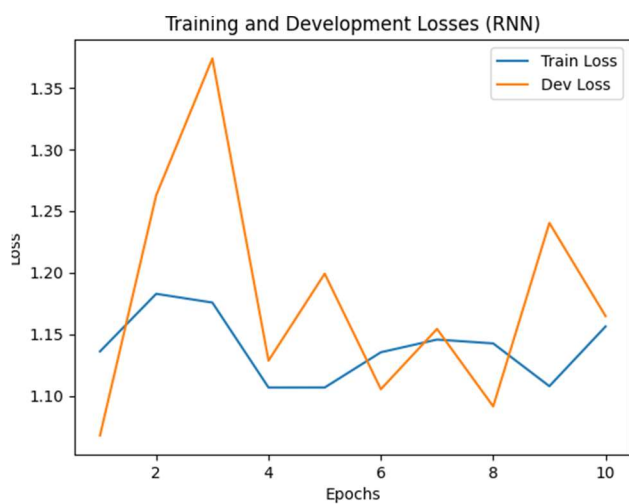
Training loss vs epoch graph and Dev Accuracy vs epoch graphs

Test Accuracy: 0.3950  
Test Precision: 0.2516  
Test Recall: 0.2995  
Test F1-Score: 0.2655

Per Class Metrics:  
Class 0 - Precision: 0.0000, Recall: 0.0000, F1-Score: 0.0000  
Class 1 - Precision: 0.4214, Recall: 0.6484, F1-Score: 0.5108  
Class 2 - Precision: 0.3333, Recall: 0.2500, F1-Score: 0.2857

Evaluation metrics for the model

**(ii) Results for RNN model:**



Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Accuracy: 0.3600
Test Precision: 0.3175
Test Recall: 0.3170
Test F1-Score: 0.3172

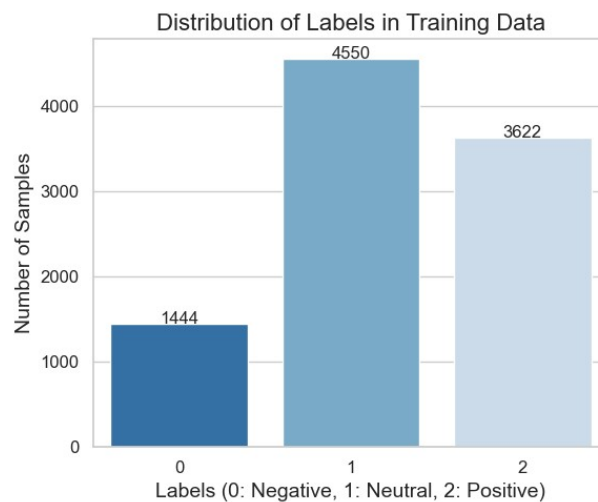
Per Class Metrics:
Class 0 - Precision: 0.1786, Recall: 0.1724, F1-Score: 0.1754
Class 1 - Precision: 0.4194, Recall: 0.4286, F1-Score: 0.4239
Class 2 - Precision: 0.3544, Recall: 0.3500, F1-Score: 0.3522
```

### Evaluation metrics for the model

This did not give good results as the accuracy was around 36% only. but it established a baseline for our model. Also, our model did not predict class 0, which implies there may be a **class imbalance** in the dataset.

Test accuracy: **0.36**

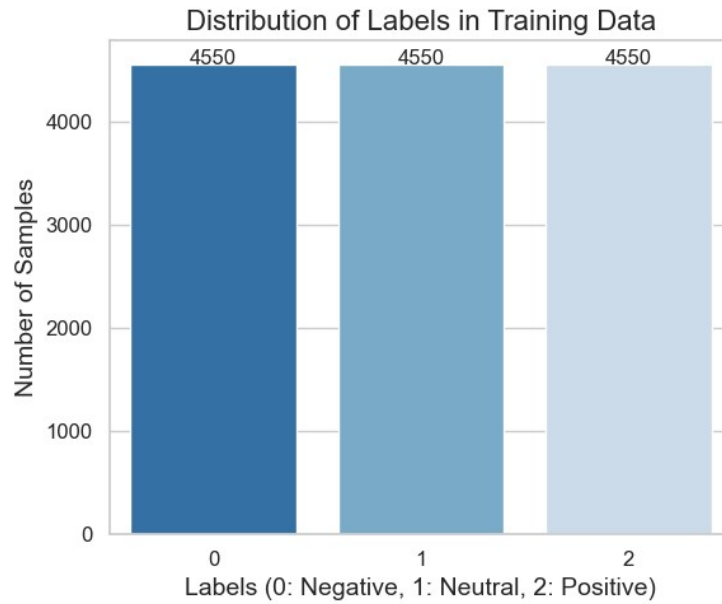
### # Removing Imbalance from the data:



From the figure, we can clearly see an imbalance in the number of samples of class 0 due to which our model did not give good predictions. Hence I applied two methods for correcting it:

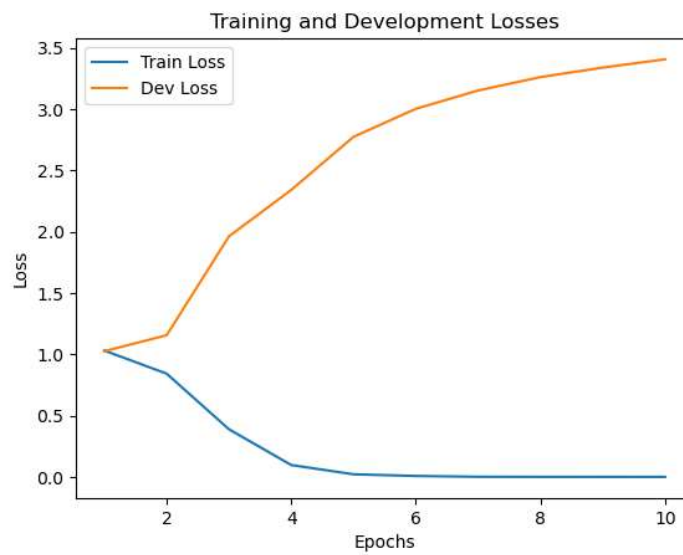
- (a) Giving more weights to undersampled class.
- (b) **Oversampling** the class with less samples.





Distribution of training data after oversampling.

-> FFN with weights to imbalanced class:



Training and dev losses vs epoch

```

Test Accuracy: 0.5180
Test Precision: 0.4513
Test Recall: 0.4496
Test F1-Score: 0.4500

Per Class Metrics:
Class 0 - Precision: 0.1852, Recall: 0.2069, F1-Score: 0.1954
Class 1 - Precision: 0.5588, Recall: 0.5676, F1-Score: 0.5631
Class 2 - Precision: 0.6098, Recall: 0.5742, F1-Score: 0.5915

```

### Evaluation Metrics

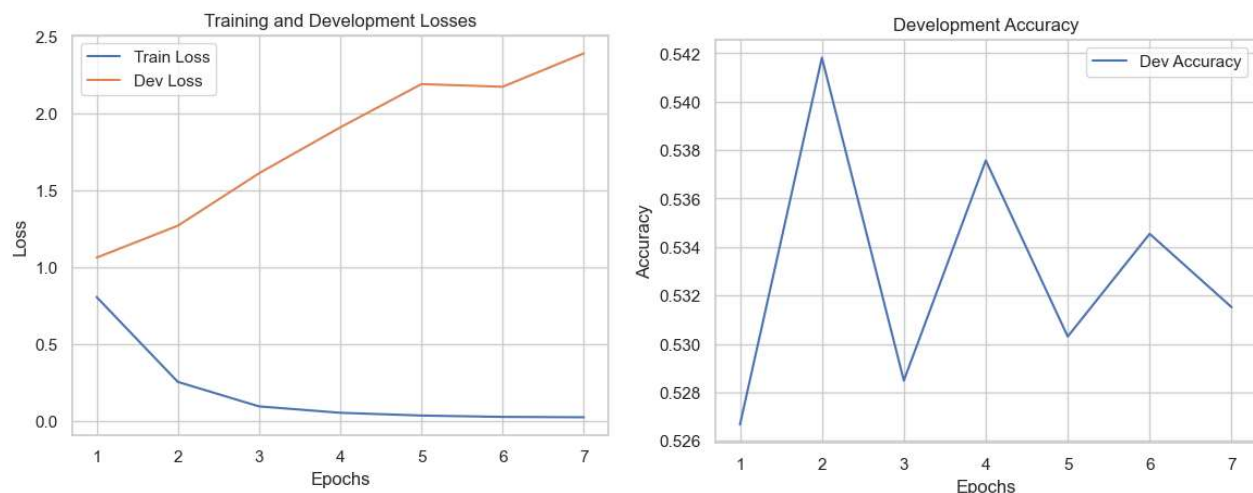
We can see that after giving weights to the imbalanced class, we get a much better model with better accuracy and f1 score than before.

-> **For new Semeval dataset with vocabulary of top frequency words with one hot encoding and oversampled training set.**

In this new dataset, the words for the vocabulary were taken to be the top **1502** words with the most frequency ( $\geq 5$ ) and accordingly the data was tokenized. This led to overall better results for the model.

Avg sentence length: 12

(i) **FFN model:**



Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```

Test Loss: 1.1761, Test Accuracy: 0.5563
Class 0: Precision: 0.3740, Recall: 0.3506, F1-Score: 0.3620
Class 1: Precision: 0.5697, Recall: 0.6345, F1-Score: 0.6004
Class 2: Precision: 0.6108, Recall: 0.5540, F1-Score: 0.5810

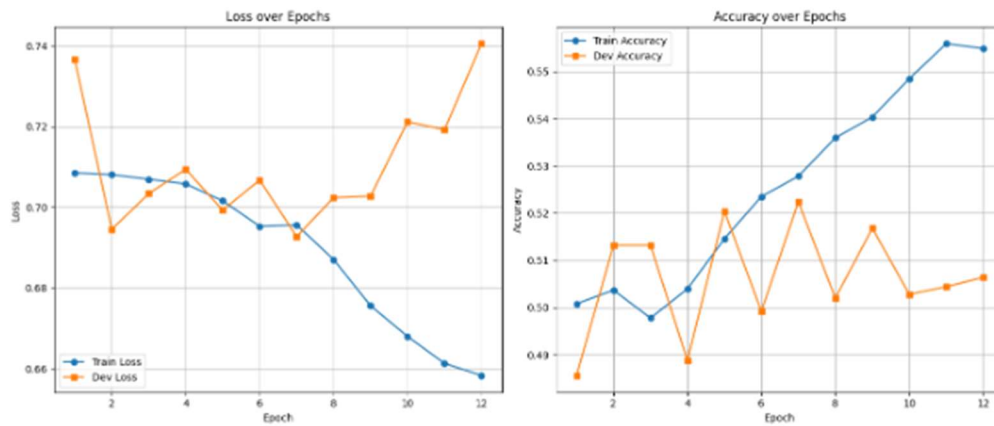
```

Evaluation metrics for the model

(ii) RNN model with embedding dim = 128.

```
SentimentRNN(  
    (embedding): Embedding(17502, 128, padding_idx=0)  
    (rnn): RNN(128, 256, batch_first=True)  
    (fc): Linear(in_features=256, out_features=2, bias=True)  
)
```

The RNN model



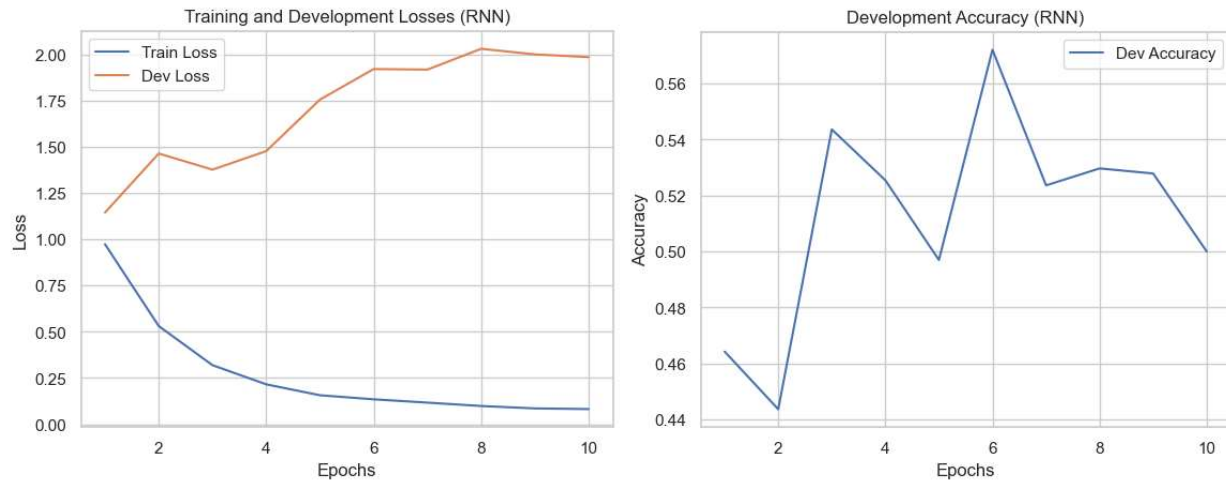
Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Loss: 0.7052  
Test Accuracy: 0.4976  
Test Precision: 0.4987  
Test Recall: 0.9220  
Test F1-Score: 0.6473  
  
Per Class Metrics:  
Class 0 - Precision: 0.4841, Recall: 0.8732, F1-Score: 0.1272  
Class 1 - Precision: 0.4987, Recall: 0.9220, F1-Score: 0.6473
```

Evaluation metrics for the model

The RNN model with embedding dimensions gave good results. So, I also tried replacing it with linear layer.

### (iii) RNN model with linear layer



Training loss vs epoch graph and Dev Accuracy vs epoch graphs

```
Test Accuracy: 0.5822
Test Precision: 0.5656
Test Recall: 0.5235
Test F1-Score: 0.5254

Per Class Metrics:
Class 0 - Precision: 0.4343, Recall: 0.3077, F1-Score: 0.3602
Class 1 - Precision: 0.5519, Recall: 0.7944, F1-Score: 0.6513
Class 2 - Precision: 0.7106, Recall: 0.4684, F1-Score: 0.5646
```

Evaluation metrics for the model

## Conclusion:

Initially, the model struggled with class imbalance, achieving a low accuracy of **36%** and failing to predict class 0. After identifying the imbalance, I applied two techniques to address it: **weighting the undersampled class and oversampling the minority class**.

Furthermore, by modifying the dataset with the top 1502 frequent words and applying tokenization accordingly, both FFN and RNN models showed better results.

By modifying the dataset and removing class imbalance, this significantly improved the model's performance, as seen in the improved accuracy to **49%** and **55%** respectively.

The RNN with an embedding dimension of 128 performed particularly well. Finally, replacing the embedding layer with a linear layer also produced strong results with accuracy of **58%** and f1-score of **52.5%**, demonstrating the effectiveness of the RNN in capturing relationships between tokens for this dataset.

-----Report Ends-----

