

# Project 2 - Emotion

Henrik Daniel Christensen<sup>[hench13@student.sdu.dk]</sup>  
Frode Engtoft Johansen<sup>[fjoha21@student.sdu.dk]</sup>

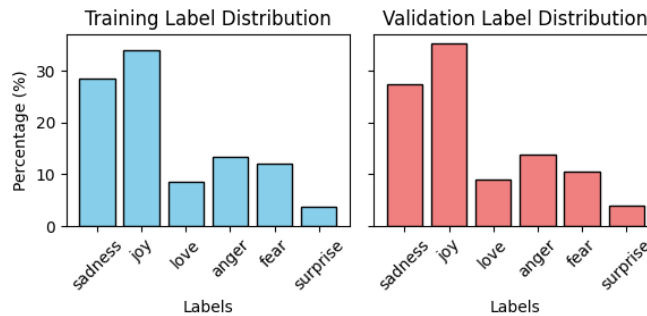
DM873: Deep Learning  
University of Southern Denmark, SDU  
*Department of Mathematics and Computer Science*

## 1 Introduction

The aim of this project is to design and implement two deep learning models for text/sentiment classification, capable of identifying the emotion conveyed in a given text. The model is trained on the *dair-ai/emotion* dataset, a widely used collection of 20,000 labeled tweets available through Hugging Face. This dataset is divided into six emotion categories: sadness (0), joy (1), love (2), anger (3), fear (4), and surprise (5). Of the total dataset, 2,000 tweets are reserved exclusively for final testing.

## 2 Label Distribution

We started by analyzing the distribution of the labels in the dataset, see Figure 1.



**Fig. 1.** Distribution of the labels in the dataset.

From the figure, it is evident that the dataset is highly imbalanced, with the majority of tweets labeled as either joy or sadness. However, the validation set exhibits a roughly similar distribution of labels, suggesting that the test set is likely to follow a comparable pattern.

This imbalance in the training data is expected to impact the model's performance, particularly for less represented labels, as sentences labeled with joy are less likely to be classified correctly. This challenge must be carefully considered when evaluating the model's performance and interpreting the results.

## 3 Preprocessing

To enable a machine learning model to process textual data effectively, it is essential to convert the text into numerical representations. The first step in this process is tokenizing the text, which involves splitting sentences into smaller units, such as words. For this task, we utilize a custom `RegexTokenizer` from the NLTK library, which filters the text to retain only words, numbers, and a limited set of special characters. This ensures that irrelevant symbols are excluded while preserving the meaningful components of the text.



Then, we convert the text data to integers by mapping each word in the text data to the corresponding integer in the vocabulary and pad the sequences to the maximum sequence length. This is done for both the training, validation and test datasets. The text data is now ready to be used for a model.

## 4 Models

Two different models were developed for the sentiment analysis task. The first model is a Recurrent Neural Network (RNN) with Long Short-Term Memory (LSTM) cells, and the second model is a Transformer model.

### 4.1 Model 1: RNN-LSTM

**Model Architecture** We chose to begin with a Recurrent Neural Network (RNN) model using LSTM (Long Short-Term Memory) cells due to their proven effectiveness in capturing long-term dependencies in sequential data, such as text. LSTMs are particularly suited for this task because they address the vanishing gradient problem often encountered in standard RNNs.

The architecture of the model, shown in Table ??, was designed to balance complexity and performance, resulting in a total of 1,117,734 trainable parameters. Key considerations behind the design are as follows:

- **Embedding Layer:** The vocabulary size of 10,336 and embedding dimension of 75 were chosen to compactly represent the input text while capturing semantic relationships between words effectively. This dimensionality ensures the embeddings are rich enough for the sentiment classification task without being computationally excessive.
- **LSTM Layer:** A single LSTM layer with 256 hidden units was selected to provide sufficient capacity to capture sequential patterns in the text. This configuration was chosen to avoid overfitting while retaining the ability to model complex dependencies in the data.
- **Dropout Layer:** A dropout rate of 50% was applied to the LSTM outputs to reduce the risk of overfitting by preventing the model from becoming overly reliant on specific neurons. This regularization technique ensures better generalization to unseen data.
- **Output Layer:** The final linear layer maps the 256 features to the 6 output classes, corresponding to the emotions in the dataset.

The specific parameters mentioned above was found using a combination of empirical testing (grid search) and theoretical considerations to achieve a balance between model complexity and performance.

**Training and Evaluation** The model training was conducted using the AdamW optimizer with a learning rate of 0.001 and a batch size of 16. AdamW was chosen due to its effectiveness in handling sparse gradients and its built-in weight decay mechanism, which helps prevent overfitting. The batch size of 16 strikes a balance between computational efficiency and the stability of gradient updates.

For the loss function, we used CrossEntropyLoss, which is well-suited for multi-class classification problems like this one, where the task involves predicting one of six emotion categories. To further regularize the model and reduce the risk of overfitting, an L2-norm penalty (weight decay) of 0.0001 was applied to the weights, encouraging smaller parameter values and promoting a simpler model. The model was trained for 4 epochs, as using more epochs lead to overfitting. The final test accuracy and F1-score is summarized in Table ??.

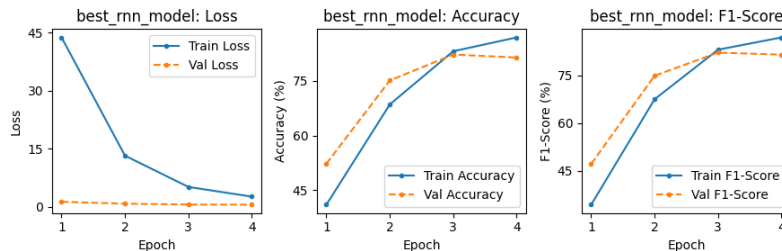


Fig. 4. Training and validation loss, accuracy, and F1-score for the RNN-LSTM model.

Label	Precision	Recall	F1-Score
sadness	0.91	0.86	0.89
joy	0.85	0.86	0.85
love	0.62	0.69	0.65
anger	0.84	0.81	0.82
fear	0.81	0.85	0.83
surprise	0.65	0.65	0.65
accuracy			0.83
weighted avg	0.84	0.83	0.83

**Table 1.** RNN-LSTM model performance on test set.

From the above table, we can see that the model performs well on sadness, joy, anger and fear, but not as well on love and surprise. This is likely due to the imbalanced distribution of the classes in the dataset as discussed in Section 2, see Figure 1.

#### 4.2 Model 2: Transformer

### 5 Analysis and Final Prediction

### 6 Pretrained Model (DistilBERT)

To evaluate how our model compares to a larger pretrained model, we utilized the DistilBERT base model (uncased) from Hugging Face. DistilBERT is a distilled version of BERT, pretrained on a large corpus comprising 11,038 books and the English Wikipedia. It features a vocabulary size of 30,000 tokens and contains approximately 67 million parameters, making it significantly larger and more complex than our models. To adapt DistilBERT to our sentiment analysis task, we fine-tuned the model for 5 epochs using a small initial learning rate of  $5 \times 10^{-6}$ . This low learning rate ensures stable and effective fine-tuning of the pretrained weights without overfitting to our dataset. The performance of DistilBERT on our task is summarized in Table 2.

Label	Precision	Recall	F1-Score
sadness	0.92	0.93	0.92
joy	0.90	0.92	0.91
love	0.74	0.71	0.72
anger	0.90	0.89	0.89
fear	0.87	0.89	0.88
surprise	0.82	0.61	0.70
accuracy			0.89
weighted avg	0.89	0.89	0.89

**Table 2.** DistilBERT model performance on test set.

From the above results, we see that the pretrained model performs generally better than our two models. Notice that the performance on the less presents classes love and surprise is also lower as seen in the previous models.

## 7 Conclusion

### 7.1 Individual Contributions

	Henrik Daniel Christensen	Frode Engtoft Johansen
<b>Code</b>	Task 1, 2, 5	Task 3, 4
<b>Report</b>	Section 1, 2, 3, 4, 7	Section 5, 6, 8

**Table 3.** Individual contributions.