



- **Stemming and lowercasing:** we reduce words to their base or root form by removing suffixes and sometimes prefixes and then we lowercase them. The main goal is to reduce the vocabulary size.

### 3 Model selection

Now we will analyze the performance of different machine learning models for our task by using different scores with different meanings. First, we will analyze the sentiment analysis task using **Macro F1-score** as a metric, then we will explore how to use our models to select the most important words for the sentiment of the sentence and measure the performance using **Jaccard scores**<sup>3</sup>.

We will explore two possible approaches for solving our task:

- **Non Deep approach** using a *Logistic Regression* as a baseline model
- **Deep approach** finetuning *DistilBERT* pretrained transformer model.

#### 3.1 Training Procedure

We split our dataset into train-validation-test (70%, 20%, 10%) maintaining the same balance between classes as in the starting data. We use the validation dataset to check for overfitting and to tune the hyperparameters of our models.

Then, we retrain our chosen model on both train and validation and we evaluate on the test set to see how well our model generalized.

#### 3.2 Logistic Regression

We start by training a simple baseline logistic regression model.

##### 3.2.1 Feature Extraction

We start with a simple bag of words model. A bag of words simply assigns an index to each word in our vocabulary, and embeds each sentence as a list of 0s, with a 1 at each index corresponding to a word present in the sentence.

In order to make the vocabulary size manageable, preprocessing the tweets is crucial, so

we use all the techniques described in section 2.2.

##### Example of features extraction:

- **Original tweet:** Hmm..You can't judge a book by looking at its cover
- **Processed tweet:** ['hmm', '..', 'judg', 'book', 'look', 'cover']

##### 3.2.2 Model validation

As we can see in Figure 1, our data set is unbalanced, so we expect that taking this into account within the *logistic regression* by adjusting the weights inversely proportional to the class frequencies in the input data will help our model.

We then perform cross-validation to choose the value for the L2 regularization term, and we obtain that the best model is with  $C=0.1$ <sup>4</sup>, which obtains a **F1-score** of **0.687** on the validation set.

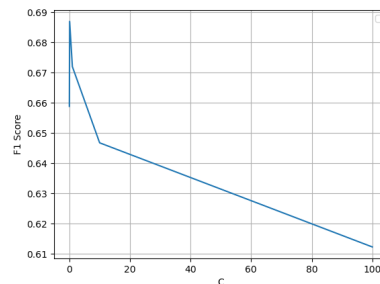


Figure 3: F1-score obtained using different values of C

##### 3.2.3 Model explainability

A big advantage of using non deep models is that it is easy to explain why a model makes certain decisions. In this case, just by looking at the weights of the decision function of our model, we can see which are the most important words for each class (*Positive*, *Neutral*, *Negative*). It's interesting to notice that, as we can see in figure 4, when it comes to the neutral class, many of the least important words are among the most important words for the positive and negative classes.

<sup>3</sup>Intersection over union of two sets of words

<sup>4</sup>Inverse of regularization strength, smaller values specify stronger regularization



This allows us to look at this interesting example, which clearly shows the different behavior between transformer models and simple Bag of Words models. If we take the sentence “*The dinner yesterday was not so amazing*” we can see in Figure 6 how the model is able to select “*was not so amazing*” and correctly label the sentence as negative because of the presence of *not*. On the other hand, our LR-BoG model misclassifies as positive the sentence because the features extracted from the tweet are only the words *dinner* and *amazing*, which, without the *not*, is mainly a positive word.



Figure 6: phrase explained by *Captum* (*transformers-interpret*)

### 3.3.3 Keywords extraction

We now try to predict the most important words to describe the sentiment of the sentence according to our model and measure its performance using **Jaccard scores**. We use the token weights given in the output of the *transformers-interpret* explainer, and we use a threshold (fine-tuned over a validation set) to decide whether a token should be considered important or not.

We get a very low **Jaccard score** of **0.289**. The main problem is that, even though we use the same tokenizer, some words are divided in chunks by the explainer, leading to a different word parsing between predicted text and the ground truth, thus to low Jaccard scores.

As a future improvement, it’s possible to try to use metrics that are less susceptible to tokenization problems to evaluate the selected text, like using a pre-trained Bert model to convert both the true selected and our predicted one to the embedding space and then compare them using a cosine similarity score (given the properties of the space, two semantically similar strings should be close). This process takes a lot more resources, but it’s more accurate.

## 3.4 Model choice

	BoG-LR	DistilBERT
<b>F1 on test</b>	0.695	0.810

We now compare our two models and see that the deep model clearly outperforms our baseline logistic regression. This is to be expected, as pre-trained transformer models have proven to be very effective in NLP tasks such as sentiment analysis.

## 4 Conclusions

In this report, we explored sentiment analysis using both a non-deep approach with Logistic Regression (BoG-LR) and a deep learning approach with DistilBERT. While Logistic Regression offers advantages in model interpretability and ease of explaining decisions based on word importance, DistilBERT significantly outperforms it in terms of classification performance, achieving a higher F1-score. Future work could focus on improving keyword extraction techniques and exploring alternative evaluation metrics to address tokenization issues.

## References

- [1] cdpierse. ‘transformers-interpret repository’.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [3] Google. Google Scholar. Last accessed 6 June 2024.
- [4] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020.
- [5] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020.