

# Shared Task Information Science - Baseline Design

Murali Manohar  
S5397294

Karlo Slot  
S4031970

Sanne Weering  
S4127579

Amber Chen  
S3333302

## 1 Introduction

SemEval 2023, more precisely the task of explainable Detection of Online Sexism (EDOS) (Guest et al., 2021) is a shared task focused on online sexism. We will be participating in this shared task and in this article we present our baselines for the SemEval 2023 task and discuss various preprocessing and feature selection techniques while implementing them.

## 2 Preprocessing

To ensure that the SemEval 2023 training data is used to its full potential, we have used several methods to clean the data for our baselines.

First, we used Python’s `.lower()` method to change all characters in the data entries to lowercase. This was done to reduce possible noise in the data and to ensure that tokens with capitalization, i.e. words at the start of a sentence, are treated the same as its non-capitalized version. One drawback is that named entities are treated the same as normal tokens. For example, “Trump” as a last name and as a verb would be mapped to the same token.

Second, we interpret this task as a textual problem focused on hate speech. Therefore, the choice has been made to remove non-words, such as URLs, interpunction and non-alphanumeric characters. This was accomplished using the regular expressions `"http\S+"` (URLs) and `"[^\w\s]"` (non-word/non-alphanumeric characters, interpunction).

Third, we tokenized and lemmatized the texts using NLTK<sup>1</sup>. For the lemmatization, we used NLTK’s `WordNetLemmatizer()`<sup>2</sup>, while we used NLTK’s `word_tokenize()`<sup>3</sup> function for tokenization.

As a final step, we removed the data entries with category ‘none’ in Task B. This was done due to the fact that the goal of Task B is classifying sexist messages in unique categories and the ‘none’ label only appears with data entries where the label of Task A is ‘not sexist’.

## 3 Baselines

To estimate the difficulty of SemEval 2023 task, we start by implementing the baselines, most frequent classifier and SVM classifier. In addition to the descriptions of the baselines, we briefly talk about the data splitting and the features used.

1. **Most frequent Classifier:** We start with a simple baseline, where we always predict an instance as the most frequent class label in the training data. While this approach may seem redundant and non-robust, this is useful in the

imbalanced scenarios. For example, if a class occurs 90% of the time in the training data, predicting all the test instances as the same class would result in 90% accuracy. This proves to be a strong baseline to beat for the upcoming complex machine learning approaches.

2. **Support Vector Machines:** As a second baseline, we evaluate Support Vector Machine (SVM). (Cortes and Vapnik, 1995) The basic notion of this algorithm is to divide the data points into different classes by identifying a hyperplane which maximizes the distance between the nearest data points. We conjecture that the n-grams can help the model to distinguish between named entities, negation, etc. Therefore, we convert comments into n-grams (unigrams and bigrams) and convert them into features by either counting the term frequencies or TF-IDFs. (Wikipedia, 2022)

The dataset is split into 70:15:15 ratio for train, dev and test splits. We used `Train Test Split`<sup>4</sup> to split the data after shuffling it. As we were not performing hyperparameter search for SVM, we merged train and dev splits into a bigger training split with 85 percent instances of the whole dataset.

Task	Baseline	Prec.	Rec.	F1	Acc.
A	Maj.	37.98	50.00	43.17	75.95
	SVM + TF	78.36	71.86	74.10	83.00
	SVM + TFIDF	80.30	71.93	<b>74.60</b>	83.71
B	Maj.	15.11	20.00	17.22	75.57
	SVM + TF	53.25	44.63	<b>47.57</b>	51.76
	SVM + TFIDF	58.80	41.46	44.86	53.92

Table 1: Table depicting the macro-average precision (Prec.), macro-average recall (Rec.), macro-average F1-score (F1) and accuracy (Acc.) of all baseline models per task. Note that, TF refers to the term frequencies and TFIDF refers to the term-frequency-inverse-document-frequency score

## 4 Results and Conclusion

In table 1, it can be seen that the majority-based baselines performs poorly, since the macro-average F1-score for Task A is 43.17, where the majority label is ‘not sexist’, and for Task B is 17.22, where the majority label is ‘2. Derogation’.

In the aforementioned table, it can also be seen that the SVM trained on TF-IDF performs better with a macro-average F1-score of 74.60 on Task A than its term frequency (TF) equivalent. And on Task B, SVM trained on TF performs slightly better than TF-IDF with a F1-score of 47.57 on Task B respectively.

<sup>1</sup>NLTK Documentation

<sup>2</sup>NLTK’s `WordNetLemmatizer()` Documentation

<sup>3</sup>NLTK’s `word_tokenizer()` Documentation

<sup>4</sup>sklearn’s `Train Test Random Splitter`

The scores of the baselines, specifically Task B, shows that there's room for improvement. From the predictions, we observed that narrations of abuse are sometimes labeled as being sexist, while they are in fact commentaries on the concept and occurrences of sexism, rather than being sexist themselves. We hypothesize that contextual models like BERT (Devlin et al., 2018) or LSTM (Hochreiter and Schmidhuber, 1997) can help in distinguishing the categories better.

## References

- Corinna Cortes and Vladimir Vapnik. 1995. [Support-vector networks](#).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.
- Ella Guest, Bertie Vidgen, Alexandros Mittos, Nishanth Sastri, Gareth Tyson, and Helen Margetts. 2021. [An expert annotated dataset for the detection of online misogyny](#).
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Wikipedia. 2022. [tf-idf](#).