

# ExCode-Mixed: Explainable Approaches towards Sentiment Analysis on Code-Mixed Data using BERT models

Aman Priyanshu, Aleti Vardhan\*, Sudarshan Sivakumar\*, Supriti Vijay\*, Nipuna Chhabra\*

Manipal Institute of Technology

{aman.priyanshu, aleti.vardhan, sudarshan.sivakumar, supriti.vijay, nipuna.chhabra}@learner.manipal.edu

## Abstract

The increasing use of social media sites in countries like India has given rise to large volumes of code-mixed data. Sentiment analysis of this data can provide integral insights into people’s perspectives and opinions. Developing robust explainability techniques which explain why models make their predictions becomes essential. In this paper, we propose an adequate methodology to integrate explainable approaches into code-mixed sentiment analysis.

## 1 Introduction

Code-mixing is the mixing of two or more languages in speech. English words, for instance, are often mixed into Hindi sentences to form what is colloquially known as *Hinglish* (Bali et al., 2014). Social media users from multilingual countries like India are seen to express themselves using code-mixed language. Sentiment analysis is a technique which can label this data as positive, negative, or neutral. The rise in the use of deep learning models for sentiment analysis has led to more accurate predictions. However, these complex models lack transparency on model predictions as they behave like black boxes of information. Explainable AI is an evolving area of research that provides a set of methods to help humans understand and interpret a model’s decisions (Miller, 2019). Specifically for the task of sentiment analysis of textual data, these techniques allow us to understand how a word or phrase influences the sentiment of the text (Bodria et al., 2020). This paper discusses the importance and practicality of model-agnostic explainable methods for sentiment analysis of code-mixed data.

## 2 Methodology and Preliminaries

**SAIL 2017** provides a sentiment classification dataset on Hindi-English code-mixed data that has been split into training, validation, and test sets with 10055, 1256, and 1257 instances, respectively (Sarkar, 2018). Every sample has been labeled as positive, negative, or neutral in sentiment.

**XLM-RoBERTa** is a transformer-based cross-lingual language model trained on 100 different languages (Conneau et al., 2020). We fine-tuned the model — pre-trained on Hindi-English code-mixed data — on our dataset. We feed the outputs of the last hidden layer in XLM-RoBERTa to a softmax layer to get the final probability-class distribution.

For this work we will be discussing two model-agnostic methods, LIME (Ribeiro et al., 2016) and SHAP (Lundberg and Lee, 2017).

**LIME** stands for Local Interpretable Model-Agnostic Explanations and examines the model’s behavior around an instance of the dataset rather than looking at the entire dataset. LIME works based on input perturbations and their respective changes on the model output.

**SHAP** or SHapley Additive exPlanations leverages the idea of Shapley values of game theory for model feature influence scoring. Unlike LIME, SHAP can be used to interpret feature dependency on the entire model.

LIME and SHAP use the trained model to extract predictions for the original as well as perturbed sentences. LIME interprets the individual sentences locally, while SHAP interprets the model as a whole. The explanations are extracted for the validation and test sets. We use these explanations to visualize and quantify our experiments.

\* All authors have contributed equally to the work.

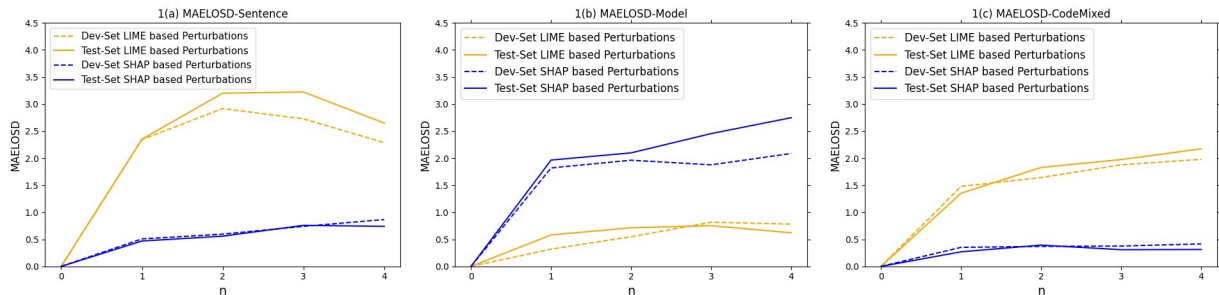


Figure 1: Metric MAELOS D for Sentence, Model, and CodeMixed is presented in this figure.

### 3 Experimental Results

This section discusses the results of both LIME and SHAP explainable approaches on the code-mixed data. We present a comparison between the two methodologies using three metrics, described by Mean Absolute Error of Log-Odds Scores on Deletion (MAELOS D) which is defined as,

$$= \sum \frac{(|\log\_odds_i - \log\_odds_f|)}{n} \quad (1)$$

where  $i$  refers to sentences without word deletions and  $f$  refers to sentences with deleted words (Chen et al., 2020). We define polarizing words as those which have been given the highest weights by the explanations of the LIME and SHAP models.

1. MAELOS D of Sentence-Interpreted Polarizing Words (*MAELOS D-Sentence*): We delete the top  $n$  polarizing words as returned by our explainable model from the sentence and recompute the Log-Odds Scores. We then calculate the Mean Absolute Error (MAE) for all samples.
2. MAELOS D of Model-Interpreted Polarizing Words (*MAELOS D-Model*): We repeat the same computation with summarized weights for the entire vocabulary(English and Code-mixed). For LIME, we take the average of the word weights across each example. The new weights now represent the most polarizing words.
3. MAELOS D of Code-Mixed Words (*MAELOS D-CodeMixed*): This metric calculates sentence-wise MAE upon deletion of top  $n$  polarizing code-mixed words. We consider all words which are not a part of the GloVe(6B) (Pennington et al., 2014) vocabulary as code-mixed. For instance, the word *accha* is a Hindi word that does not

appear in the GloVe vocabulary and is hence considered code-mixed.

We expect the mean absolute error to increase on the deletion of polarizing words since the model would make poor predictions.

We can observe the results of all three metrics in Fig 1. Fig 1(a) — MAELOS D-Sentence — illustrates LIME’s local advantage over SHAP, while Fig 1(b) — MAELOS D-Model — represents SHAP’s global advantage over LIME. We also observe that upon deleting  $n \geq 3$  words, sentences become too short, thus returning random predictions.

We can also see the impact of code-mixed data from our results in Fig 1(c), where the error increases in proportion to the number of words deleted. These words have a similar impact on label prediction to that of English words. Even so, global SHAP explanations do not hold as much impact as local LIME explanations, thus implying that the *Hinglish* vocabulary is much more diverse. This may be due to the deletion of globally important words that may not be present locally.

Our results align with the application of both LIME and SHAP, demonstrating their local and global natures. We can see the application and easy integration of model-agnostic interpretability pipelines on code-mixed data.

### 4 Conclusion and Future Work

Code-mixed data is an integral part of communication in multilingual communities. The use of SHAP and LIME, which quantify global and local model explanations, allows us to display their application and importance for sentiment analysis on code-mixed data. For future work, we would like to formalize and derive a quantifiable metric for global explanations, comparing the impact of code-mixed data with standard English sentences.

We would also like to experiment with different model architectures and datasets. We believe that our work serves as a valuable resource for the code-mixed AI community. The integration of explainable methods in code-mixed data paves a new path towards future research and development.

## Acknowledgments

The authors of the paper are grateful to the reviewers in reviewing the manuscript and their valuable inputs are appreciated. We would also like to thank the Research Society MIT for supporting the project.

## References

- Kalika Bali, Jatin Sharma, Monojit Choudhury, and Yogarshi Vyas. 2014. “i am borrowing ya mixing?” an analysis of english-hindi code mixing in facebook. In *Proceedings of the First Workshop on Computational Approaches to Code Switching*, pages 116–126.
- Francesco Bodria, A. Panisson, A. Perotti, and Simone Piaggese. 2020. Explainability methods for natural language processing: Applications to sentiment analysis. In *SEBD*.
- Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. 2020. Generating hierarchical explanations on text classification via feature interaction detection.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#).
- Scott Lundberg and Su-In Lee. 2017. [A unified approach to interpreting model predictions](#).
- Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “why should i trust you?”: Explaining the predictions of any classifier.
- Kamal Sarkar. 2018. [Ju\\_ks@sail.codemixed-2017: Sentiment analysis for indian code mixed social media texts](#).