

DataFrame Column	Description
id_new	Created
text	Unaltered original text
extracted_target	Extracted from original targets
label_hatespeech_binary_offensive_not_included	Offensive not considered hate speech (binary)
label_hatespeech_binary_offensive_included	Offensive included as hate speech (binary)
label_normal_offensive_hatespeech	Normal, offensive, or hate speech (multi-class)
platform	Source of text
original_dataset_title	Original dataset name
original_id	Original dataset ID
original_label	Original label from dataset
original_target	Majority-annotated targets

Table 1: Merged dataset column descriptions

Setting	Value / Description
Seed	42 (for ‘random‘, ‘numpy‘, ‘torch‘, CUDA)
Models	Cardiff (‘twitter-roberta‘) and Facebook (‘roberta-hate-speech‘)
Prediction Method	HuggingFace ‘transformers.pipeline‘
Batch Size	512
<b>LIME Settings</b>	
Number of Features	4
Number of Samples	500
<b>SHAP Settings</b>	
Explainer Type	shap.Explainer
Max Evaluations	100
Masker	shap.maskers.Text with tokenizer

Table 2: Unified experimental settings for LIME and SHAP explanations