

Data Collection and Preprocessing Phase

Date	26 November 2024
Team ID	SWTID1726490119
Project Title	Toxic Comment Classification for Social Media
Maximum Marks	6 Marks

Preprocessing Template

The text data will be preprocessed by cleaning, normalizing, tokenizing, augmenting, denoising, embedding, applying batch normalization, and, if necessary, whitening. These steps will improve data quality, support model generalization, and enhance convergence during neural network training, ensuring effective and reliable performance across various natural language processing tasks.

Section	Description
Data Overview	Overview of the text data used in the project, including the train.csv and test.csv files, with features like comment_text and multi-label targets (toxic, severe_toxic, obscene, threat, insult, and identity_hate).
Text Cleaning	Lowercasing: converting all text to lowercase to maintain uniformity. Removing special characters: Eliminating symbols, punctuation, and numbers. Removing single characters: Removes isolated single characters (e.g., "a", "b"). Removing multiple spaces: Consolidates multiple spaces into a single space.
Tokenization	Split sentences into individual tokens (words or sub-words) to allow neural network processing. Using a tokenizer to convert comment_text into sequences for further input into models.
Padding	Padding standardizes all input sequences by adding zeros to shorter sequences or truncating longer ones at the end, creating uniform input shapes for the model.
Handling imbalanced data	Calculates class weights to handle class imbalance by converting multi-label data to single-label, computing balanced class weights, and storing them in a dictionary for model

	training.
Data splitting	Split the data into training and validation sets to monitor model performance and prevent overfitting.
Data Preprocessing Code Screenshots	
Data Overview	<pre>[] train_data.describe()</pre> 
Text Cleaning	<pre>[] def clean_text(text): text = text.lower() # Convert to lowercase text = re.sub(r'\W', ' ', text) # Remove special characters text = re.sub(r'\s+[a-zA-Z]\s+', ' ', text) # Remove single characters text = re.sub(r'\s+', ' ', text) # Remove multiple spaces return text [] train_data['comment_text'] = train_data['comment_text'].apply(clean_text)</pre>
Tokenization	<pre>[] num_words=20000 tokenizer = Tokenizer(num_words) # Only the top 20,000 words tokenizer.fit_on_texts(train_data['comment_text']) X_train = tokenizer.texts_to_sequences(train_data['comment_text']) X_test = tokenizer.texts_to_sequences(test_data['comment_text']) [] import pickle # Assuming `tokenizer` is your fitted Tokenizer instance with open('tokenizer.pkl', 'wb') as handle: pickle.dump(tokenizer, handle, protocol=pickle.HIGHEST_PROTOCOL)</pre>
Padding	<pre>[] maxlen = 200 # Maximum length of input sequences X_train= pad_sequences(X_train, padding='post', maxlen=maxlen,truncating='post') X_test = pad_sequences(X_test, padding='post', maxlen=maxlen,truncating='post') [] # Output the shapes to verify print(f'Train data shape: {X_train.shape}, Train labels shape: {y_train.shape}') print(f'Test data shape: {X_test.shape}')</pre> 

Handling imbalanced data	<pre>[] # Handle class imbalance # Calculate class weights from sklearn.utils import class_weight y_train_argmax = np.argmax(y_train, axis=1) # Convert to single-label for class weight calculation class_weights = class_weight.compute_class_weight('balanced', classes=np.unique(y_train_argmax), y=y_train_argmax) class_weights_dict = {i: class_weights[i] for i in range(len(class_weights))}</pre>
Data splitting	<pre>[] # Split the data for training and validation X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=0.2, random_state=42)</pre>