

# Toxic Comment Classification

---

## **Problem Statement**

Online platforms and social media websites often face the challenge of dealing with toxic and offensive comments. These comments can be harmful, abusive, or contain hate speech, which creates a negative environment for users. Detecting and filtering such toxic comments manually is a time-consuming and daunting task. Therefore, there is a need to develop automated systems that can classify toxic comments accurately.

**Dataset:** The dataset used in this project is a collection of comments from various sources, labeled with different categories of toxicity. The dataset contains two main components:

- 1) **Train Data:** This dataset consists of comments along with their corresponding labels for six categories of toxicity: toxic, severe\_toxic, obscene, threat, insult, and identity\_hate. Each comment is labeled with 0 or 1, indicating whether it belongs to a specific toxic category.
- 2) **Test Data:** This dataset contains comments without any labels. The goal is to classify these comments into the toxic categories based on the trained model.

**Problem Statement:** The problem at hand is to develop a machine learning model that can accurately classify toxic comments into different categories of toxicity. Given a comment, the model should predict the probabilities of it belonging to each toxicity category.

**Approach:** The code provided implements a toxic comment classification model using a recurrent neural network (RNN) with LSTM layers. The following steps are performed:

- 1) **Data Loading:** The train and test datasets are loaded from CSV files.
- 2) **Data Preprocessing:** Null values in the datasets are checked and filled to ensure clean data for modeling.

- 3) Tokenization:** The comment texts are tokenized using the Keras Tokenizer class, which converts the comments into sequences of numerical indices. Each unique word is assigned a unique index.
- 4) Padding:** The tokenized sequences are padded or truncated to a fixed length using the Keras pad\_sequences function. This ensures that all input sequences have the same length, necessary for inputting into the RNN model.
- 5) Model Architecture:** The model architecture is defined using Keras functional API. It consists of an embedding layer, LSTM layer, global max pooling layer, and dense layers with dropout regularization. The model takes the padded tokenized sequences as input and outputs a probability distribution over the toxicity categories using a sigmoid activation function.
- 6) Model Compilation and Training:** The model is compiled with the binary cross-entropy loss function and the Adam optimizer. It is then trained on the training data with a specified batch size and number of epochs.

**Objective:** The objective of this project is to train a model that can accurately classify toxic comments into different categories of toxicity. The trained model can be used by online platforms to automatically identify and filter toxic comments, creating a safer and more positive environment for users.

# **Framework**

## **1. Importing Required Libraries**

- Import the necessary libraries such as numpy, pandas, matplotlib, and the Keras modules for building the model.

## **2. Loading the Dataset**

- Read the train and test data from CSV files using the `pd.read_csv()` function.
- Display a sneak peek of the training dataset using the `train.head()` function.

## **3. Data Preprocessing**

- Check for null values in the train and test datasets using the `isnull().any()` function.
- Fill any null values with appropriate values to ensure clean data for modeling.

## **4. Tokenization and Sequencing**

- Create a list of toxic comment categories (`list_classes`) that we want to classify.
- Extract the comment texts from the train and test datasets into separate lists (`list_sentences_train` and `list_sentences_test`).
- Set the maximum number of unique words to consider (`max_features`).
- Create a `Tokenizer` object and fit it on the training data to generate a word index.
- Convert the comment texts into sequences of numerical indices using the `texts_to_sequences()` function of the tokenizer.
- Store the tokenized sequences in `list_tokenized_train` and `list_tokenized_test`.

## **5. Padding**

- Set the maximum sequence length (`maxlen`) for the comments.
- Pad or truncate the tokenized sequences to the specified maximum length using the `pad_sequences()` function.

## **6. Exploratory Data Analysis (EDA)**

- Calculate the distribution of the total number of words in the comment texts using a histogram.
- Visualize the histogram using `matplotlib` to understand the length distribution of the comments.

## **7. Model Architecture**

- Define the input shape for the model based on the maximum sequence length.
- Set the embedding size (embed\_size) for the word embeddings.
- Create an input layer (inp) with the specified shape.
- Add an embedding layer that maps the words to their respective dense vectors.
- Add an LSTM layer with a specified number of units and return sequences.
- Apply global max pooling to extract the most relevant features.
- Add dropout regularization to prevent overfitting.
- Add dense layers with activation functions to make the final predictions.
- Create a model using the Keras functional API by specifying the input and output layers.

## **8. Model Compilation and Training**

- Compile the model by specifying the loss function, optimizer, and evaluation metrics.
- Set the batch size and number of epochs for training.
- Fit the model on the training data, using a portion of it for validation.

## **9. Model Evaluation**

- Evaluate the trained model on the test data to measure its performance.
- Compute metrics such as accuracy, precision, recall, and F1-score to assess the model's effectiveness.

## **10. Prediction**

- Use the trained model to make predictions on new comments or unseen data.
- Apply thresholding techniques to classify the comments into different toxic categories based on the predicted probabilities.

## **11. Model Deployment**

- Deploy the trained model into a production environment, such as an online platform or API, for real-time toxic comment classification.

## **12. Conclusion**

- Summarize the project, highlighting the importance of automated toxic comment classification and the effectiveness of the developed model.
- Discuss potential future improvements and enhancements to the model or the overall system.

## **Code Explanation**

**1. Importing Required Libraries** The code starts by importing the necessary libraries. These libraries provide various functions and tools to perform different tasks in the project. For example, numpy and pandas are used for data manipulation and analysis, matplotlib is used for data visualization, and the Keras modules are used for building the neural network model.

**2. Loading the Dataset** In this section, the code loads the training and testing datasets from CSV files. The `pd.read_csv()` function is used to read the files and store them in separate variables. The `train.head()` function is then used to display a preview of the training dataset, showing the first few rows of data.

**3. Data Preprocessing** Data preprocessing is an essential step before training a machine learning model. In this code, the first step is to check for any null values in the train and test datasets using the `isnull().any()` function. This helps us identify if there are any missing values in the data. If null values are present, they need to be handled or filled before proceeding further.

**4. Tokenization and Sequencing** To process the text data, it needs to be converted into numerical format that can be understood by the machine learning model. Tokenization is the process of splitting text into individual words or tokens. The code creates a list of toxic comment categories and extracts the comment texts from the train and test datasets. It then uses the `Tokenizer` class from Keras to convert the text into sequences of numerical indices. This step helps in representing the words in a way that the model can process.

**5. Padding In natural language** processing tasks, sequences can have varying lengths. However, for efficient training of neural networks, it's important to have fixed-length input sequences. In this code, the sequences obtained from tokenization are padded or truncated to a fixed length using the `pad_sequences()` function. This ensures that all sequences have the same length, allowing the model to process them uniformly.

**6. Exploratory Data Analysis (EDA)** EDA is an important step to understand the data better and gain insights. In this code, a histogram is created to visualize the distribution of the total number of words in the comment texts of the training set. The histogram

helps in understanding the length distribution of the comments, which can be useful for further analysis and model design decisions.

**7. Model Architecture** This section defines the architecture of the neural network model. It starts by setting the input shape based on the maximum sequence length. An embedding layer is added to map the words to dense vectors. An LSTM layer with a specific number of units is added to capture the sequence information. Global max pooling is applied to extract the most important features. Dropout regularization is used to prevent overfitting. Finally, dense layers with activation functions are added to make the final predictions.

**8. Model Compilation and Training** Before training the model, it needs to be compiled with appropriate settings. This involves specifying the loss function, optimizer, and evaluation metrics. In this code, binary cross-entropy is used as the loss function since it is a multi-label classification problem. The model is then trained using the training data. The batch size and number of epochs are set to control the training process.

**9. Model Evaluation** After training, the model is evaluated on the test data to measure its performance. Metrics such as accuracy and loss are computed to assess how well the model is performing. These metrics provide insights into the model's ability to classify toxic comments correctly.

**10. Prediction Once the model** is trained and evaluated, it can be used to make predictions on new comments or unseen data. The model takes in the input text, applies preprocessing steps, and generates predictions for the different toxic comment categories. Thresholding techniques can be applied to classify the comments into different categories based on the predicted probabilities.

**11. Model Deployment** After building and testing the model, it can be deployed into a production environment. This could involve integrating the model into an online platform or creating an API that allows real-time toxic comment classification. Deploying the model makes it accessible to users who can utilize its functionality for their own purposes.

**12. Conclusion To summarize**, the provided code performs the following steps: loading the dataset, preprocessing the data, tokenization and sequencing, padding, exploratory data analysis, defining the model architecture, compiling and training the model, evaluating the model's performance, making predictions, and potentially deploying the

model. These steps collectively aim to build a toxic comment classification model that can automatically identify different types of toxic comments.



## **Concept Explanation**

Imagine you're a comment police officer, and your job is to identify toxic comments on the internet. But hey, there are so many comments out there, and you can't manually read them all day long. That's where the algorithm comes to your rescue!

This algorithm is like your super-smart partner who knows how to spot those nasty, toxic comments. Let's break it down step by step, shall we?

**Step 1: Loading the Training Data** First, you gather a bunch of comments that are already labeled as toxic or not. It's like collecting evidence against those toxic comment culprits. You make two piles: one for the training comments and another for testing comments.

**Step 2: Preprocessing** Now, you need to prepare the comments for analysis. It's like putting on your detective hat and getting everything in order. You check if any comments are missing, like a piece of evidence disappearing. If there are any missing comments, you fill in the blanks so that your algorithm doesn't get confused later.

**Step 3: Tokenization and Sequencing** To teach your algorithm how to understand the comments, you need to convert them into a language it understands — numbers! Each word in the comments is like a clue, and the algorithm needs to know their secret numerical codes. So, you assign a unique code to each word in the comments, creating a sequence of numbers that represents the comments.

**Step 4: Padding** In real life, comments can have different lengths. Some are short and snappy, while others are long and rambling. But your algorithm needs consistency, like keeping all your detective notes in neat and tidy folders. So, you pad the comments with extra zeros or trim them down to a fixed length, ensuring that they all have the same length.

**Step 5: Exploratory Data Analysis (EDA)** Now it's time to dive deeper into the comments. You want to know how long they are on average. It's like measuring the size of the criminal files. You create a histogram, which is like a fancy graph, showing you how many comments have certain lengths. This helps you understand the patterns and plan your investigation strategy.

**Step 6: Building the Model** Alright, detective, it's time to create your secret weapon! You design a powerful neural network that will learn how to spot toxic comments. It's like training a highly skilled detective who can detect even the slightest hint of toxicity. The network has special layers like LSTM, which can remember important details and make accurate predictions.

**Step 7: Training the Model** Now comes the exciting part — training your detective network! You feed it the training comments and let it learn from them. It's like sending your detective to a training camp where they analyze and understand the toxic comments. The network adjusts its internal parameters to improve its detective skills, based on the feedback you provide.

**Step 8: Evaluation** Once your detective network completes its training, it's time to test its skills. You give it some test comments and see how well it performs. It's like putting your detective to the test in a real case scenario. You measure its accuracy and other performance metrics to see if it can correctly identify toxic comments.

**Step 9: Prediction** Congratulations, detective! Your network is ready to tackle new comments. You give it fresh, unseen comments, and it predicts whether they are toxic or not. It's like having a super-intelligent partner who can quickly scan through comments and flag the toxic ones. You can now take appropriate action and keep the internet a safer place!

## **Exercise Questions**

- 1) Question: What is the purpose of tokenization and sequencing in this project? How does it help in training the toxic comment classification model?**

**Answer:** Tokenization and sequencing play a crucial role in preparing the comments for analysis and training the model. Tokenization involves breaking down the comments into individual words or tokens. This step is like converting the comments into a language the model can understand. Sequencing assigns a unique numerical code to each token, creating a sequence of numbers that represents the comments. This helps in transforming the textual data into a numerical format that the model can process. By tokenizing and sequencing, the model learns the patterns and relationships between different words in the comments, enabling it to make predictions based on those sequences.

- 2) Question: Explain the purpose of padding in the toxic comment classification model. Why is it necessary to have comments of the same length?**

**Answer:** Padding is necessary in the model to ensure that all comments have the same length. In real-life scenarios, comments can vary in length, but for efficient training of the model, it's essential to have consistent input sizes. Padding involves adding extra zeros or truncating the comments to a fixed length. By doing so, we create a uniform structure where each comment has the same length as specified. This allows the model to process the comments in a batch and perform computations efficiently. Padding ensures that all comments are compatible with the model's architecture and guarantees consistent input dimensions for training and prediction.

- 3) Question: What is the purpose of the LSTM layer in the toxic comment classification model? How does it contribute to the model's ability to detect toxic comments?**

**Answer:** The LSTM (Long Short-Term Memory) layer is a special type of recurrent neural network layer used in the model. Its purpose is to capture and remember important information from the comments, considering the contextual dependencies between words. In the context of detecting toxic comments, the LSTM layer helps the model

understand the sequential nature of the text and grasp the meaning behind the comment's structure. By learning long-term dependencies, the LSTM layer enables the model to identify relevant patterns and detect toxic language even when it occurs at different positions within the comment. This layer's ability to retain context over time makes it effective in understanding the nuances of toxic comments.

**4) Question: Explain the activation function used in the dense layers of the toxic comment classification model. Why is it chosen for this task?**

**Answer:** The activation function used in the dense layers of the model is the ReLU (Rectified Linear Unit) activation function. ReLU is a popular choice in many neural network architectures due to its simplicity and effectiveness. It introduces non-linearity into the model, enabling it to learn complex relationships between the input data and the target labels. In the context of toxic comment classification, ReLU activation helps in capturing the nonlinear patterns and features that indicate toxic behavior. It allows the model to discriminate between toxic and non-toxic comments by activating relevant neurons based on the learned weights. ReLU has proven to be effective in various deep learning tasks, including text classification, and thus, it is chosen for this specific task as well.

**5) Question: How can you improve the performance of the toxic comment classification model beyond the provided code?**

**Answer:** There are several ways to improve the performance of the toxic comment classification model beyond the provided code:

- **Model Architecture:** Experiment with different architectures, such as adding more LSTM layers or using convolutional layers in conjunction with LSTMs. These variations can capture different aspects of the comment's structure and potentially enhance performance.
- **Hyperparameter Tuning:** Explore different hyperparameter settings, such as learning rate, batch size, and dropout rates. Conducting a systematic search for optimal hyperparameters using techniques like grid search or random search can help fine-tune the model.
- **Word Embeddings:** Utilize pre-trained word embeddings, such as Word2Vec or GloVe, to initialize the embedding layer. These embeddings capture semantic

relationships between words and can provide a better representation of the comment text.

- Ensemble Learning: Combine predictions from multiple models trained with different configurations to leverage diverse perspectives and improve overall accuracy.
- Data Augmentation: Generate synthetic training examples by applying techniques like word replacement, synonym substitution, or sentence paraphrasing. This approach can increase the diversity of the training data and improve the model's ability to generalize.
- Regularization: Apply regularization techniques like L1 or L2 regularization to prevent overfitting and improve generalization.
- Advanced Architectures: Explore advanced architectures such as Transformer models, which have shown excellent performance in natural language processing tasks. Transformers can capture long-range dependencies and improve the understanding of the comment's context.
- Implementing these strategies and carefully fine-tuning the model can lead to improved performance in classifying toxic comments. Remember, it's all about experimenting, iterating, and finding the best approach for the specific problem at hand!