

# **IMPLEMENTATION OF HATE-SPEECH USING** **TRANSFORMERS**

Name – Deekshant Nandeshwar  
Email – deekshantnandeshwar@gmail.com  
Country – United Kingdom  
College – University of Surrey  
Specialization – Master of Science (Data Science)

## Problem Description–

Any form of communication—verbal, written, or nonverbal—that targets or employs disparaging or discriminatory language against an individual or group because of who they are—their religion, ethnicity, nationality, race, color, ancestry, sex, or another identity characteristic—is considered hate speech. We'll walk you through a model that detects hate speech in this issue.

The task of sentiment categorization is typically involved in hate speech detection. Therefore, by using data that is often used to categorize attitudes, a model that can identify hate speech from a specific piece of text may be trained. As a result, in order to identify tweets that include hate speech, we will use the Twitter tweets.

## Project lifecycle along with deadline –

Weeks 07	Problem description, Project lifecycle along with deadline, Data Report
Weeks 08	Data Analysis
Weeks 09	different featurization technique
Weeks 10	EDA performed on the data
Weeks 11	EDA presentation for business users
Weeks 12	Model Selection and Model Building
Weeks 13	Final Project Report and Code

## Data Collection

The data on Twitter hate speech was extracted from Kaggle and includes 3 characteristics and 31962 observations. It was a dataset created using Twitter data and used to study the identification of hate speech. The text is categorized as either offensive language, hate speech, or neither. It is crucial to be aware that this dataset includes language that could be seen as objectionable, racist, sexist, or homophobic due to the nature of the study.

## Data Analysis

- Text Cleaning - We cleaned our text because the data was so disorganized.
- Remove Punctuation - Punctuation should be deleted since it is unnecessary. Because of this, we utilize regular expression to eliminate the punctuation.

- Lowercase - Lowercase word conversion . Although terms like racism and racism have the same meaning when written in lower case, the vector space model represents them as two distinct meanings (resulting in more dimensions). As a result, we change all content to lower case letters.
- Remove URL - In this section, URLs are removed since we are working on a hate speech program that detects hate and free speech and because in order to acquire the result, we can only provide text and not URLs.
- Remove @ and Special Character - We eliminate the @tags, which were essentially used when mentioning someone. Which doesn't matter. The group of symbols known as Remove Special Characters essentially has no meaning.