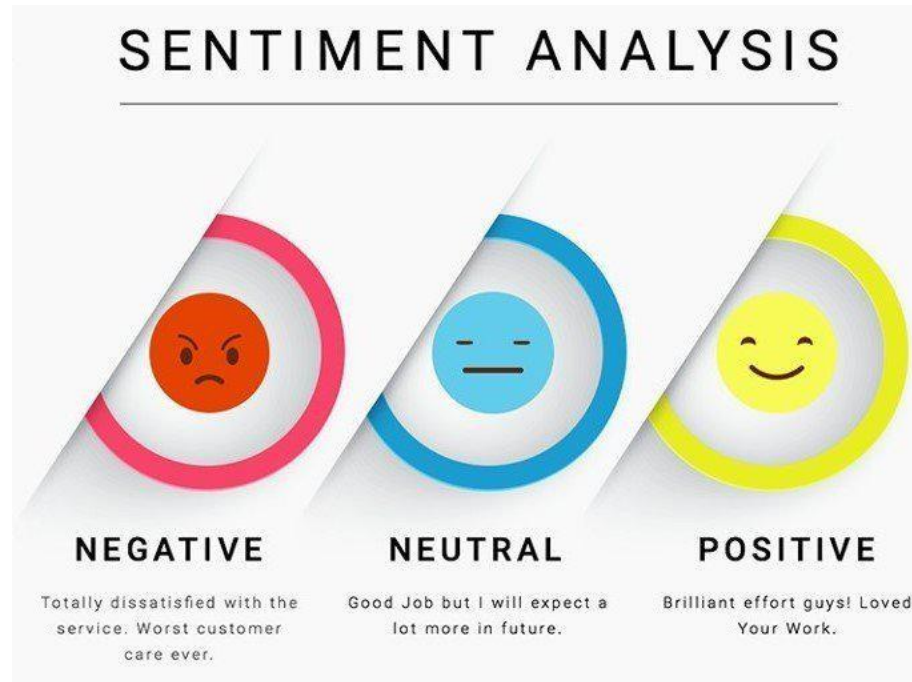


Introduction:



Hello everyone, and welcome to this end-to-end project on sentiment analysis using Jupyter Notebook in Colab!

My name is Shaikh Mohammad Mujammil, and I'm delighted to be your guide throughout this End to End Sentiment Analysis Project. Today, we will embark on a journey to create a sentiment analysis model using a pre-built state-of-the-art model from Hugging Face, all within the powerful environment of Jupyter Notebook.

Now, you might be wondering, what is the purpose of this project? Well, our goal is to delve deep into the realm of sentiment analysis, an essential technique in natural language processing. Sentiment analysis allows us to understand and analyze the sentiment or emotion expressed in a given text, enabling us to gain valuable insights from vast amounts of textual data.

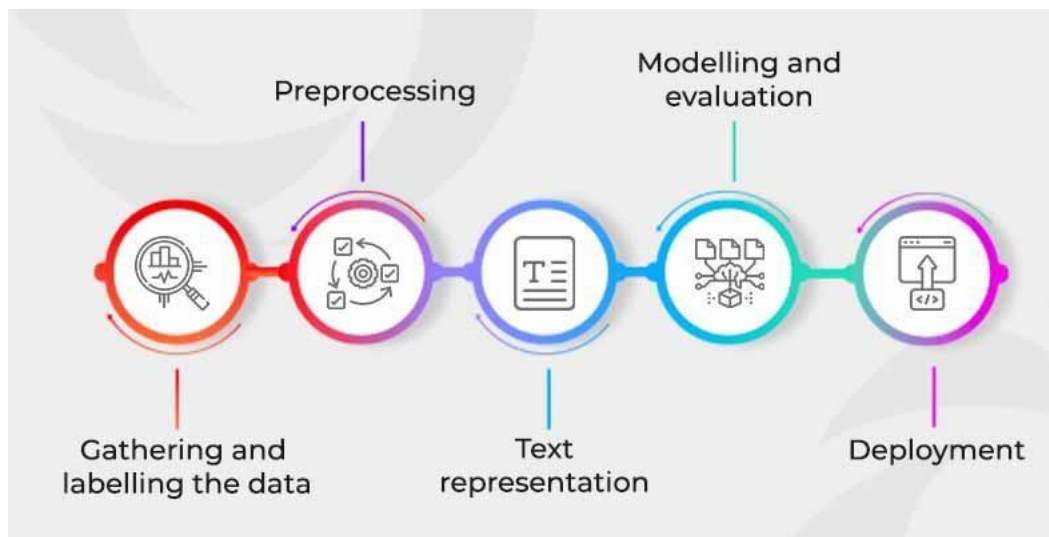
Unlike conventional sentiment analysis tutorials, we're taking it a step further. This Project is designed for those with an advanced level of Python knowledge. Together, we will explore the intricacies of sentiment analysis, from data collection and preprocessing to model selection, evaluation, and visualization.

Throughout this project, i will employ Jupyter Notebook, a versatile and interactive tool that empowers us to write and execute code, annotate our findings, and create compelling visualizations. Colab, Google's cloud-based platform, will provide us with the necessary computational resources to handle the complexities of sentiment analysis.

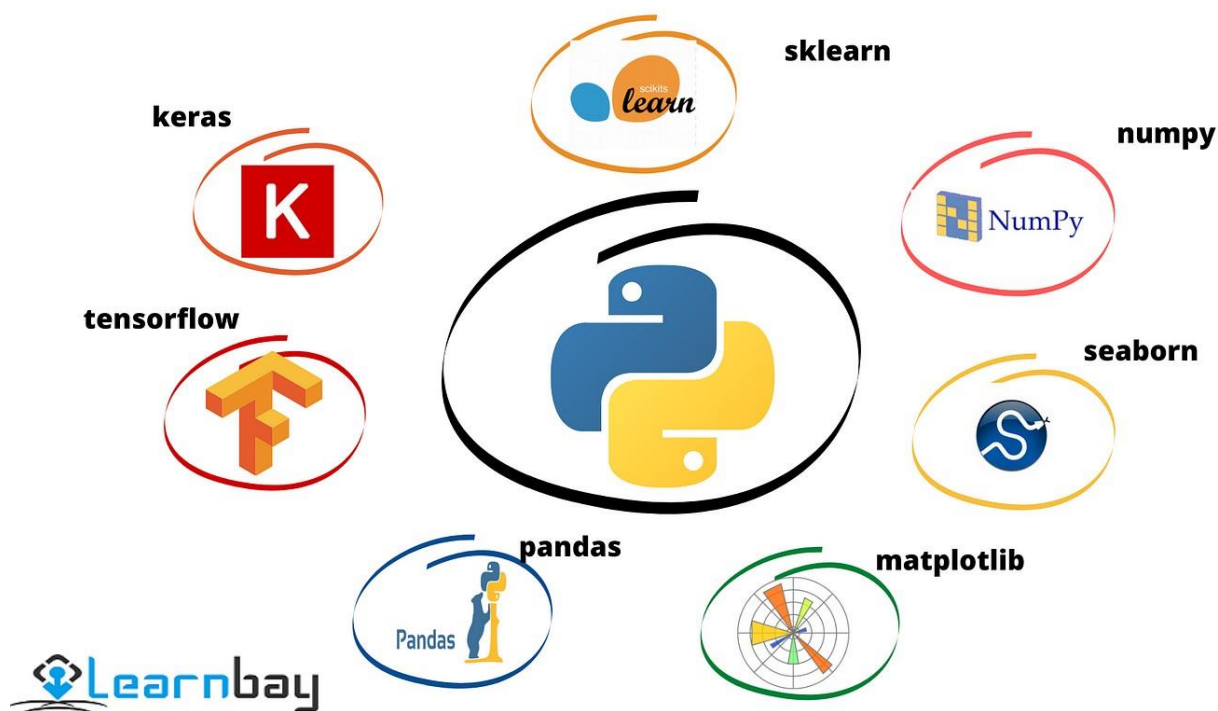
Sentiment analysis is a thriving field, with numerous real-world applications. It plays a crucial role in social media monitoring, customer feedback analysis, market research, and many other domains. By mastering the art of sentiment analysis, you can unlock the potential to extract valuable insights from text data and make data-driven decisions.



Now, let's take a moment to set our expectations for this project. By the end of this tutorial, you will have a solid understanding of how to build an end-to-end sentiment analysis model using Jupyter Notebook. We will explore the key steps involved in this process, from collecting relevant data to preprocessing it, selecting an appropriate model, evaluating the results, and visualizing the outcomes in a comprehensive dashboard.



To embark on this sentiment analysis journey, it's essential to have an advanced level of Python knowledge. Familiarity with libraries such as pandas, numpy, transformers, and torch will greatly benefit you throughout this project. However, even if you're new to sentiment analysis, fear not! We will provide detailed explanations and code snippets to help you grasp the concepts and techniques involved.



So, are you ready to dive into the fascinating world of sentiment analysis? Excellent! Together, we will uncover the sentiment behind texts, discover hidden patterns, and gain actionable insights from textual data.

Throughout this Project, feel free to ask questions, experiment with the provided code, and explore additional resources for a deeper understanding. We're here to learn and grow together.

Without further ado, let's begin our journey into sentiment analysis using Jupyter Notebook in Colab. Get ready to unlock the power of text analysis and take your data exploration to new heights!

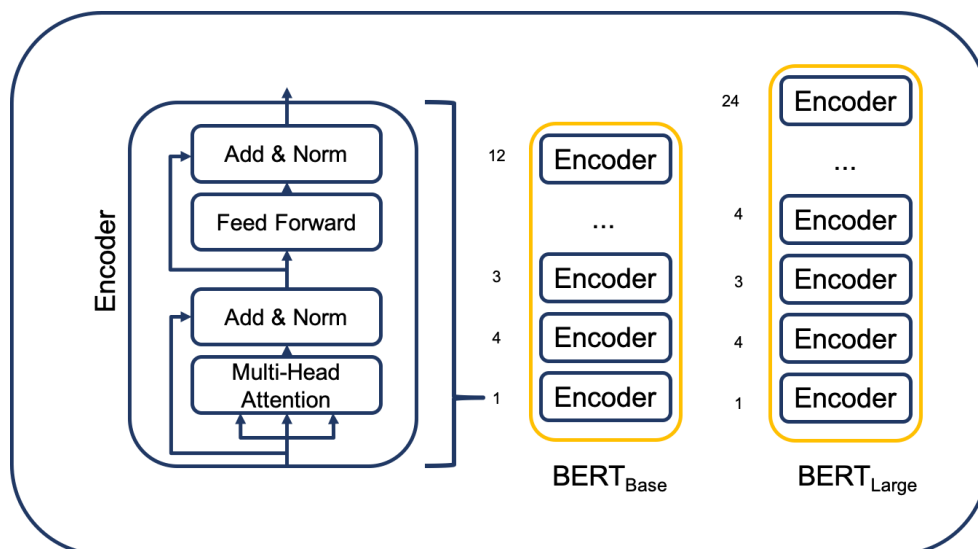
Background Information:

Before we delve into the project details, let's take a moment to understand the background of sentiment analysis and its significance in the realm of text analysis.

Sentiment analysis, also known as opinion mining, is a powerful technique that enables us to understand and analyze the sentiment or emotion expressed in a given text. It involves the use of natural language processing and machine learning algorithms to automatically determine the sentiment conveyed by a piece of text, whether it's positive, negative, or neutral.

The importance of sentiment analysis cannot be overstated. In today's digital age, where vast amounts of textual data are generated through social media, customer reviews, surveys, and more, sentiment analysis provides valuable insights into public opinion, customer satisfaction, market trends, and brand reputation. By accurately gauging the sentiment behind texts, businesses can make data-driven decisions, tailor their marketing strategies, and improve customer experiences.

Now, let's talk about the availability of pre-built state-of-the-art models for sentiment analysis. Thanks to advancements in natural language processing and the efforts of researchers and developers, we now have access to powerful models that have been trained on massive amounts of text data. These models, such as BERT (Bidirectional Encoder Representations from Transformers), provide a solid foundation for sentiment analysis tasks. They have learned to capture the nuances of language, contextual understanding, and sentiment dynamics, allowing us to achieve accurate sentiment analysis results.



In our project, we will leverage one of these pre-built state-of-the-art models from Hugging Face. This model has been fine-tuned specifically for sentiment analysis tasks and is known for its high performance and accuracy. By utilizing such models, we can save time and effort in training our own sentiment analysis model from scratch and focus on the other crucial steps involved in the end-to-end process.

Now, let's talk about the tools we will use for this project. Jupyter Notebook in Colab provides an excellent environment for conducting this end-to-end project. Jupyter Notebook allows us to write and execute code in an interactive manner, making it easy to experiment, iterate, and annotate our findings. Colab, being a cloud-based platform, offers us the advantage of accessing computational resources, including GPUs, that can handle the computational demands of sentiment analysis.

The combination of Jupyter Notebook and Colab provides a seamless and efficient workflow for sentiment analysis projects. It empowers us to leverage pre-built models, manipulate and analyze data using Python libraries like pandas and numpy, visualize our results, and document our process in a single, comprehensive notebook.

With this background information in mind, let's dive into our end-to-end sentiment analysis project. We will explore the key steps involved, from data collection and preprocessing to model selection, evaluation, and visualization. Get ready to unravel the sentiments behind texts and gain valuable insights from our text analysis journey.

Let's embark on this exciting project and unlock the power of sentiment analysis using Jupyter Notebook in Colab!

Data Collection:

The first crucial step in our sentiment analysis project is data collection. Collecting a suitable dataset is essential for training and evaluating our sentiment analysis model. In this section, we will discuss the process of collecting data and the dataset we have chosen for our project.

As a starting point, I extensively researched and explored various sources to find relevant datasets for sentiment analysis. I utilized AI-powered web searches and platforms like Kaggle, which host a wide range of datasets contributed by the community. These sources offered a wealth of options to choose from and ensured that our dataset is diverse and representative.

After careful consideration, I selected the Hasoc 2019 dataset as our primary dataset for this project. The Hasoc 2019 dataset is a well-known dataset in the field of hate speech and

offensive language identification. It has been used in multiple research studies and competitions. Let's take a closer look at the features of the Hasoc 2019 dataset.

HASOC_ID	TEXT	TYPE_1	TYPE_2	TYPE_3
hasoc_en_1	#DhoniKeepsTheGlove WATCH: Sports Minister Kiren Rijiju issues statement backing MS Dhoni over 'Balidaan Badge', tells BCCI to take up the matter with ICC and keep government in the know as nation's pride is involved https://t.co/zuo5335Rjr	NOT	NONE	NONE
hasoc_en_10	By wearing the #BalidaanBadge over his gloves @msdhoni has shown his love & respect for the forces. @icc should understand that this is not related to any political/religious/racial activities This is about our #NationalPride #DhoniKeepsTheGlove	NOT	NONE	NONE
hasoc_en_100	Marlon Craft - Gang Shit (Official Music Video) https://t.co/jh7QJds6n0 via @YouTube @GOP #FuckTrump #CopsAreTheGang	NOT	NONE	NONE
hasoc_en_1000	ICC is trying their best to stop dhoni from showing his love to army.Dhoni should retaliate by donating his one year income to indian army. #CWC19 #DhoniKeepsTheGlove #DhoniKeepTheGloves #DhoniKeSaathDesh	NOT	NONE	NONE
hasoc_en_1001	#DoctorsFightBack Parents reaction when a doctor returns home safely in Bengal...🤔🤔🤔 https://t.co/8Ar617t80U	NOT	NONE	NONE
hasoc_en_1002	@NatalieHarp @JoeBiden @realDonaldTrump #FridayFeeling Blaming Joe Biden? Proves you're truly #Deplorable. Best of luck. #FuckTrump, #FucktheGOP and #FuckTrumpSupporters. Democrats have a #healthcare plan. Republicans have a #dontcare plan. https://t.co/3lc0t6nZmj	HOF	PRFN	TIN
hasoc_en_1003	We stand with @msdhoni infact #indiastandswithdhoni #DhoniKeepsTheGlove #IndiaStandsWithDhoni https://t.co/58GaSwrSuK	NOT	NONE	NONE
hasoc_en_1004	👉 Tripura doctors protest in solidarity with doctors in West Bengal #Justice4DrParibaha #DoctorsFightBack 👉 https://t.co/gYuHm8hxfF @ATGDAofficial @sanjayswadesh @ChoudhuryKanak @Bismita14 @talk_anderson @JasBJP @Satyanewshi @biswajitroy2009 @AD_BJP @trunilss @jayhind22090440	NOT	NONE	NONE
hasoc_en_1005	@saquibIM @sardesaiarajdeep @msdhoni What rules you are referring? Rules of having namaz in the play ground? We can understand your feeling. #biasedicc #DhoniKeepTheGlove #IndiaStandsWithDhoni #DhoniKeSaathDesh #DhoniKeepsTheGlove #dhonigloves	HOF	HATE	TIN
hasoc_en_1006	@jyotika35246268 @narendramodi This man completely nailed Mamata, her politics and her hypocrisy. #DoctorsFightBack https://t.co/yvS1nr5Zp	NOT	NONE	NONE
hasoc_en_1007	@MedCrisis Thanks @MedCrisis for covering this. For years the working conditions at the government run hospitals remained the same. Being an alumni from Kolkata I feel sorry for my medical friends, seniors and juniors. @arjunk @SauravChMD we support you #DoctorsFightBack	NOT	NONE	NONE
hasoc_en_1008	yo, fuuuuuuuuuuuuuuuuuuuuuuuuuuuuck trump. y'all's president is TRASH!!!!!! #ImpeachTrump #fucktrump https://t.co/QOf71Fr5rT	HOF	PRFN	TIN
hasoc_en_1009	Belgian backpacker reveals how Guntree #rapist lured her into shed and kept her #hostage - Jun 16 @ 10:50 AM ET https://t.co/O3tWQJtrB2	NOT	NONE	NONE
hasoc_en_101	@BorisJohnson If we want to unite our country- we should commit that anyone who us scared to face a bit of public scrutiny shouldn't be allowed to run for Prime Minister... #BorisJohnsonShouldNotBePM #BorisJohnson #ToryLeadership #ToryLeadershipDebate #ToriesOut	NOT	NONE	NONE
hasoc_en_1010	Dear ICC, please don't forget that BCCI is your father..... So,mind on your business.... Do better umpiring.... #DhoniKeepsTheGlove #DhoniKeepsTheGlove	HOF	HATE	TIN
hasoc_en_1011	A billion Cricket fans waited four years to watch how it rains in England #shameonicc https://t.co/n3enbfwl78	HOF	HATE	TIN
hasoc_en_1012	This is the man @JamesCleverly @MattHancock @trussliz that you think will make a great #PM for fucks sake, God help us #BorisJohnsonShouldNotBePM #ToryLeadershipContest https://t.co/j0GoJgY10A	HOF	PRFN	TIN
hasoc_en_1013	@ICC @ICCCAD @BCCI Poorest management of #CWC19 Need to change venues Please take a look of geographical And environmental conditions Not only on gloves..... #ShameOnICC	HOF	HATE	TIN
hasoc_en_1014	And the protest is on @ArdPgimer #DoctorsFightBack #DoctorsStrike #doctors_against_assault @httweets https://t.co/nLTk7KvTjF	NOT	NONE	NONE

The Hasoc 2019 dataset is well-organized, labeled, and provides text samples in three different languages: English, Hindi, and German. This multilingual aspect of the dataset allows us to analyze sentiments expressed in different languages and explore the dynamics of hate speech and offensive language across diverse linguistic contexts.

The dataset offers three sub-tasks, each focusing on different aspects of hate speech and offensive language identification. These sub-tasks provide us with a comprehensive view of the dataset and enable us to analyze sentiments at different levels of granularity.

Type A focuses on hate speech and offensive language identification. It involves a coarsegrained binary classification, where participating systems classify tweets into two classes: Hate and Offensive (HOF) and Non-Hate and Offensive (NOT). This classification is based on the presence or absence of hate speech, offensive content, and profanity in the tweets.

Type B represents a fine-grained classification. It further classifies the hate speech and offensive posts from Type A into three categories: Hate speech (HATE), Offensive (OFFN), and Profane (PRFN). This fine-grained classification allows us to understand the type of content within the hate speech and offensive categories in more detail.

Type C, available only for English and Hindi, considers the type of offense in the hate speech and offensive posts labeled as HOF in Type A. It classifies them into two categories: Targeted Insult (TIN) and Untargeted (UNT). This classification helps us identify whether the offensive content is directed towards a specific individual or group or if it contains non-targeted profanity.

By leveraging the Hasoc 2019 dataset and its sub-tasks, we can gain valuable insights into the dynamics of hate speech and offensive language across different languages and analyze sentiments at various levels of specificity.

In the next section, we will delve into the data preprocessing steps required to prepare the dataset for analysis. Stay tuned as we explore the journey from raw data to actionable insights.

Data Preprocessing:

Once we have collected our dataset, the next crucial step is data preprocessing. Data preprocessing involves preparing the dataset for analysis by cleaning and organizing it. In this section, we will discuss the steps involved in preprocessing the dataset for our sentiment analysis project.

To streamline the process and ensure the dataset is in a suitable format, I utilized Excel to format and combine the datasets. Excel provides a user-friendly interface for data manipulation and allows us to perform various operations efficiently.

The first step in preprocessing the dataset was to merge multiple datasets into a single cohesive dataset. I used Excel's functionality to combine datasets, ensuring that the resulting dataset contains all the necessary information for our analysis.

During this merging process, I performed some specific operations to enhance the dataset's usability. One of these operations involved splitting the ID column into a language column.

By extracting the language information from the ID column, we can differentiate between English, Hindi, and German texts, which is essential for our multilingual sentiment analysis.

Additionally, I dropped the ID column from the dataset as it was not necessary for our analysis. This step helped to simplify the dataset structure and focus on the relevant information.

Ensuring data quality and accuracy is crucial for any analysis. To achieve this, I removed empty rows, duplicates, and rows with NaN (blank) values from the dataset. Removing empty rows and duplicates helps eliminate inconsistencies and ensures that we are working with clean and reliable data. NaN values, which indicate missing or empty cells, were also removed to avoid any issues during the analysis.

INDEX	LANG	TEXT	TYPE_A	TYPE_B	TYPE_C
1	en	Young Eritreans are risking death to migrate. Here's why https://wef.ch/2Gz5xK0 #Refugees #Migration	NOT	NONE	NONE
2	en	Thanks @LRMNetwork @LewishSanctuary @valewisham @kevinbonavia for bringing us together to discuss to how to make #Lewisham a #sanctuaryborough #RefugeesAreWelcome #RefugeeWeek2019	NOT	NONE	NONE
3	en	Fuck. Is this really happening! You guys are awesome!!! Please make a vlog there! Copyright my ass ICC!	HOF	PRFN	TIN
4	en	Former Vice President of Aligarh Muslim University Students' Union (AMUSU) Hamza Sufiyan rusticated for 5 years for allegedly storming office of the Vice Chancellor with armed outsiders and debarred students. #AMU #AMUSU	NOT	NONE	NONE
5	en	BBC News - Waltham Forest Pride: Police probe 'homophobic abuse' video https://www.bbc.com/news/uk-england-london-49147605 ... She is a disgrace to Islam	NOT	NONE	NONE
6	en	'We Should've Built A Wall' Upset your racist relatives next Thanksgiving. https://www.teepublic.com/t-shirt/5467819-we-shouldve-built-a-wall?store_id=123265 ... GET 35% OFF EVERYTHING!! #NativeAmerican #fucktrump #Immigration #immigrants #immigrantsgetthejobdone #progressive #socialist #nowall #dontbuildthewall #DoubleStandards #Hypocrisy	NOT	NONE	NONE
7	en	Dr. Qanta Ahmed: Rep. Omar is a disgrace to Islam https://youtu.be/4o-RL7yqV6Y via @YouTube YOU CAN SAY THAT AGAIN!	HOF	HATE	TIN
8	en	Good economy my ass! One of our Lowes this morning eliminated 2 departments and fired 6 people without notice. One was there 11 years. #lowes #fucktrump	HOF	PRFN	TIN
9	en	Yes Halala girl, How about having jhatka meat up yours ?	NOT	NONE	NONE
10	en	I wouldn't bother mate - guy clearly has an agenda. Finds no issue with @BorisJohnson being a huge and open racist yet this offends him? Logic is clearly not in abundance here	HOF	OFFN	TIN
11	en	If you really care what people think & you believe the majority of people would back this reckless non-plan, then why not put it to the people with a #PeoplesVote, instead of using meaningless twitter pledges to mine data? Or would that be too democratic for @BorisJohnson ?	NOT	NONE	NONE
12	en	I do not think the British people will take being lied to again. Can Nigel Farage see into the future I was saying give him a chance Nigel but I do not like the rumours I am hearing today and I hope it is not true.	NOT	NONE	NONE
13	en	Unless your marrying someone who's not Muslim or someone of the same sex, 99% all parents will accept it at some point	HOF	HATE	TIN
14	en	Fuck icc for firing him part 3 will be a flop won't be watching	HOF	PRFN	UNT
15	en	On our way to El Paso #FamiliesBelongTogether #michaelavenatti #washingtondc #dc #impeachtrump #whitehouse #washingtondc #dc #25thamendment #25thamendmentnow #impeach #resist #familiesbelongtogether #vote #nomuslimban #nomuslimbanever #keepfamiliesaltogether #bluewave #enough #eno	NOT	NONE	NONE
16	en	Great read. Dont Miss at all. Each and every point has been explained with every bit of detail. #ICC #WTC21 https://twitter.com/krick3r/status...	NOT	NONE	NONE
17	en	Disagreeing with a person of color isn't RACIST Disagreeing with a woman isn't SEXIST Disagreeing with LGBT isn't HOMOPHOBIA Disagreeing with a Muslim isn't ISLAMOPHOBIA Disagreeing with an immigrant isn't XENOPHOBIA	NOT	NONE	NONE
18	en	Couldn't agree more. N.Ireland voted to remain.	NOT	NONE	NONE
19	en	Used to take my parents to speakers corner for a fun day out. Now it's a national disgrace, soaked in diversity not diverse opinions. Islam is peace, or else.	NOT	NONE	NONE
20	en	Pig 3 really having 'fun' with #EuropeanCricketLeague standards, because of course. I've seen Indians lately argue that it's okay if cricket becomes IPL only, and if anyone expects the ICC will open up the game, they're stupid. What a sorry bunch of people.	HOF	OFFN	TIN

By performing these preprocessing steps, we prepared the dataset for further analysis and ensured that it is in a suitable format for our sentiment analysis project. In the next section, we will explore the exploratory data analysis (EDA) process to gain insights into the dataset's characteristics and distribution of sentiments.

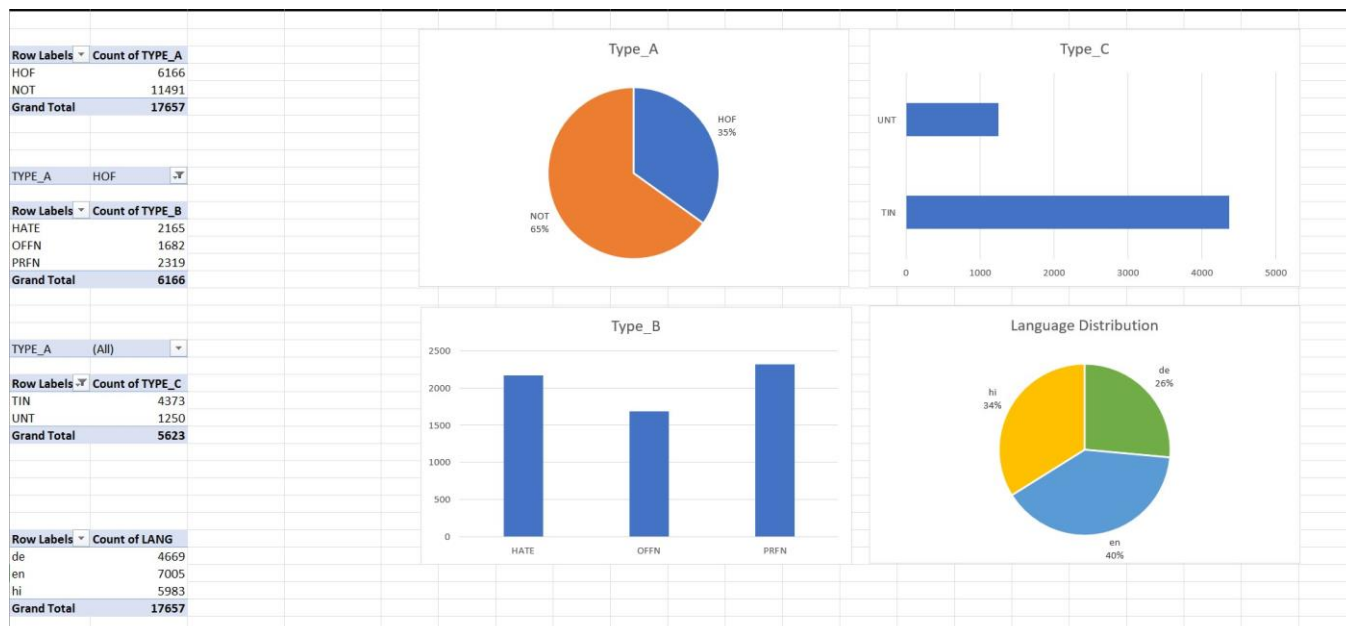
Exploratory Data Analysis (EDA):

Once we have preprocessed the dataset, the next step is to perform Exploratory Data Analysis (EDA). EDA helps us gain insights and understand the characteristics of the dataset, providing a foundation for further analysis. In this section, we will explore the key findings from our EDA and share interesting insights.

The first aspect we examined during the EDA was the language distribution within the dataset. Analyzing the language distribution helps us understand the composition of texts in different languages and their significance for sentiment analysis.

Based on our analysis, we found that approximately 40% of the dataset consists of English texts, 34% consists of Hindi texts, and 26% consists of German texts. This distribution highlights the diversity of languages in the dataset, allowing us to perform sentiment analysis across multiple languages.

Our analysis revealed that 65% of the dataset was labeled as Non-Hate and Offensive (NOT), while the remaining 35% fell under the Hate and Offensive (HOF) category. This distribution showcases the distribution of sentiments within the dataset and serves as a basis for our sentiment analysis.



Another interesting aspect we investigated during the EDA was the uniqueness of the data. Analyzing the percentage of unique data helps us understand the diversity and richness of the dataset, contributing to the reliability of our analysis.

Based on our analysis, we found that the dataset contains approximately 17,657 unique data points. This substantial number of unique data points highlights the diversity and richness of the dataset, providing a solid foundation for our sentiment analysis.

These key findings from our EDA provide valuable insights into the dataset's characteristics and distribution of sentiments. The language distribution, label distribution, and percentage of unique data enable us to understand the composition of the dataset and its relevance for our sentiment analysis.

Language Translation:

One of the challenges we encountered in our sentiment analysis project was the presence of a mixed-language dataset. Working with a mixed-language dataset can introduce complexities and hinder the analysis process. To ensure compatibility and consistency in our analysis, we needed to convert the text to a single language, preferably English.

Initially, we attempted to address this challenge by writing a Python program for language translation. However, we found that the process was time-consuming, especially considering the size of the dataset and the complexity of language translation. It took approximately 4 hours of runtime to translate the text, and some data couldn't be translated, leaving blank strings in their place.

Realizing the limitations and time constraints of our initial approach, we explored alternative solutions. That's when we discovered an efficient and effective method using Google Sheets and its built-in translation formula.

formula " =googletran(text,source,target)"

Here's how we approached language translation using Google Sheets:

1. We created a new column specifically for translation in the Google Sheets document.
2. We leveraged the language column and text column in a translation formula. By utilizing the language column, we could identify the language of each text, and by using the text column, we could provide the text to be translated.
3. We performed batch translation in Google Sheets, ensuring that the translation process was efficient and didn't overload the system. By applying the translation formula to a batch of data, we could handle the translation process in a systematic and manageable manner.

INDEX	LANG	TEXT	TRANSLATED	TYPE_A	TYPE_B	TYPE_C
1	en	Young Eritreans are risking death to migrate. Here's why https://wef.ch/2Gz5xK0 #Refugees #Migration	Young Eritreans are risking death to migrate. Here's why https://wef.ch/2Gz5xK0 #Refugees #Migration	NOT	NONE	NONE
2	en	Thanks @LRMNetwork @LewishSanctuary @valewisham @kevinbonavia for bringing us together to discuss to how to	Thanks @LRMNetwork @LewishSanctuary @valewisham @kevinbonavia for bringing us together to	NOT	NONE	NONE
3	en	Fuck. Is this really happening! You guys are awesome!!! Please make a vlog there! Copyright my ass ICC!	Fuck. Is this really happening! You guys are awesome!!! Please make a vlog there! Copyright my ass ICC!	HOF	PRFN	TIN
4	en	Former Vice President of Aligarh Muslim University Students' Union (AMUSU) Hamza Sufiyan rusticated for 5 years for	Former Vice President of Aligarh Muslim University Students' Union (AMUSU) Hamza Sufiyan rusticated	NOT	NONE	NONE
5	en	BBC News - Waltham Forest Pride: Police probe 'homophobic abuse' video https://www.bbc.com/news/uk-england-	BBC News - Waltham Forest Pride: Police probe 'homophobic abuse' video	NOT	NONE	NONE
6	en	'We Should've Built A Wall' Upset your racist relatives next Thanksgiving. https://www.teepublic.com/t-shirt/5467819-	'We Should've Built A Wall' Upset your racist relatives next Thanksgiving. https://www.teepublic.com/t-	NOT	NONE	NONE
7	en	Dr. Qanta Ahmed: Rep. Omar is a disgrace to Islam https://youtu.be/4o-RL7yqV6Y via @YouTube YOU CAN SAY THAT	Dr. Qanta Ahmed: Rep. Omar is a disgrace to Islam https://youtu.be/4o-RL7yqV6Y via @YouTube YOU	HOF	HATE	TIN
8	en	Good economy my ass! One of our Lowes this morning eliminated 2 departments and fired 6 people without notice.	Good economy my ass! One of our Lowes this morning eliminated 2 departments and fired 6 people	HOF	PRFN	TIN
9	en	Yes Halala girl, How about having jhatka meat up yours ?	Yes Halala girl, How about having jhatka meat up yours ?	NOT	NONE	NONE
10	en	I wouldn't bother mate - guy clearly has an agenda. Finds no issue with @BorisJohnson being a huge and open racist yet	I wouldn't bother mate - guy clearly has an agenda. Finds no issue with @BorisJohnson being a huge and	HOF	OFFN	TIN
11	en	If you really care what people think & you believe the majority of people would back this reckless non-plan, then why	If you really care what people think & you believe the majority of people would back this reckless non-	NOT	NONE	NONE
12	en	I do not think the British people will take being lied to again. Can Nigel Farage see into the future I was saying give him a	I do not think the British people will take being lied to again. Can Nigel Farage see into the future I was	NOT	NONE	NONE
13	en	Unless your marrying someone who's not Muslim or someone of the same sex, 99% all parents will accept it at some	Unless your marrying someone who's not Muslim or someone of the same sex, 99% all parents will	HOF	HATE	TIN
14	en	Fuck icc for firing him part 3 will be a flop won't be watching	Fuck icc for firing him part 3 will be a flop won't be watching	HOF	PRFN	UNT
15	en	On our way to El Paso #FamiliesBelongTogether #michaelavenatti #washingtondc #dc #impeachtrump #whitehouse	On our way to El Paso #FamiliesBelongTogether #michaelavenatti #washingtondc #dc #impeachtrump	NOT	NONE	NONE
16	en	Great read. Dont Miss at all. Each and every point has been explained with every bit of detail. #ICC #WTC21	Great read. Dont Miss at all. Each and every point has been explained with every bit of detail. #ICC	NOT	NONE	NONE
17	en	Disagreeing with a person of color isn't RACIST Disagreeing with a woman isn't SEXIST Disagreeing with LGBT isn't	Disagreeing with a person of color isn't RACIST Disagreeing with a woman isn't SEXIST Disagreeing with	NOT	NONE	NONE
18	en	Couldn't agree more. N.Ireland voted to remain.	Couldn't agree more. N.Ireland voted to remain.	NOT	NONE	NONE
19	en	Used to take my parents to speakers corner for a fun day out. Now it's a national disgrace, soaked in diversity not	Used to take my parents to speakers corner for a fun day out. Now it's a national disgrace, soaked in	NOT	NONE	NONE
20	en	Pig 3 really having 'fun' with #EuropeanCrickettLeague standards, because of course. I've seen Indians lately argue that	Pig 3 really having 'fun' with #EuropeanCrickettLeague standards, because of course. I've seen Indians	HOF	OFFN	TIN
21	en	This fellow is worse than a stone pelter. He should immediately resign. Most Intolerant like his sister Mamata Banerji	This fellow is worse than a stone pelter. He should immediately resign. Most Intolerant like his sister	HOF	HATE	TIN
22	en	Let this sink in... please repost #Farrakhan #realwords #realtalk #yougoback #fucktrump #notrump2020 #lockhimup	Let this sink in... please repost #Farrakhan #realwords #realtalk #yougoback #fucktrump #notrump2020	NOT	NONE	NONE
23	en	What the fuck is going on here @ICC !! Bunch of fuckwits controlling cricket. @ECB_criccket , @BCCI are oppressing the	What the fuck is going on here @ICC !! Bunch of fuckwits controlling cricket. @ECB_criccket , @BCCI are	HOF	PRFN	TIN
24	en	Good on #MayorPete "...the thing you <GOP loads> will be remembered for is whether in this moment, with this	Good on #MayorPete "...the thing you <GOP loads> will be remembered for is whether in this moment,	NOT	NONE	NONE
25	en	#ThoughtForTheDay #freepress #ParisAttacks Empowering site:Should FreeExpression have limit http://bit.ly/Upv9XI	#ThoughtForTheDay #freepress #ParisAttacks Empowering site:Should FreeExpression have limit	NOT	NONE	NONE
26	en	@ZomatoIN #BycottZomato HALAL AND HALALA Are two most cruel possible things possible to an animal or a woman.	@ZomatoIN #BycottZomato HALAL AND HALALA Are two most cruel possible things possible to an	HOF	HATE	TIN
27	en	You know we Sri lankans are biggest cheaters and chucks... we cant raise a finger on anyone. We lankans need to get	You know we Sri lankans are biggest cheaters and chucks... we cant raise a finger on anyone. We	HOF	PRFN	TIN
28	en	Bengal doctors crisis: A state which has pulled out of Ayushman Bharat, the poor are the worst affected of Mamata's	Bengal doctors crisis: A state which has pulled out of Ayushman Bharat, the poor are the worst affected	NOT	NONE	NONE
29	en	everyday another Mass Shooting in the USA. All on Trumps clock.. And what does he do... Ignore & support the NRA.	everyday another Mass Shooting in the USA. All on Trumps clock.. And what does he do... Ignore &	NOT	NONE	NONE
30	en	See unknown person this is good if this govt ban halala first, this is also haram practice, if you read the process of talaq	See unknown person this is good if this govt ban halala first, this is also haram practice, if you read the	NOT	NONE	NONE
31	en	This Photoshoot Of A Hindu-Muslim, Same Sex Couple Proves That Love Has No Boundaries	This Photoshoot Of A Hindu-Muslim, Same Sex Couple Proves That Love Has No Boundaries	NOT	NONE	NONE
32	en	It's simple, folks. We want a great new deal - but if the EU won't give us one, we'll leave with no deal on October 31st.	It's simple, folks. We want a great new deal - but if the EU won't give us one, we'll leave with no deal on	NOT	NONE	NONE
33	en	Keep in touch with the ones who have forgotten you, & forgive who has wronged you, & do not stop praying for the best	Keep in touch with the ones who have forgotten you, & forgive who has wronged you, & do not stop	NOT	NONE	NONE
34	en	Master of Deception #monoprnt #printmaking #art #artists #fucktrump	Master of Deception #monoprnt #printmaking #art #artists #fucktrump	HOF	OFFN	TIN
35	en	If @POTUS is lying then Pradhan Mantri ji should suspend all relations with the US till they accept they lied	If @POTUS is lying then Pradhan Mantri ji should suspend all relations with the US till they accept they	NOT	NONE	NONE
36	en	Ireland prove that england never won the worldcup just ICC gave them favour fuck icc pigs #ENGvIRE	Ireland prove that england never won the worldcup just ICC gave them favour fuck icc pigs #ENGvIRE	HOF	PRFN	TIN
37	en	@SenateDems @HouseDemocrats DO YOUR FUCKING JOBS!! #FuckTrump And Who Votes For Him! You Need To	@SenateDems @HouseDemocrats DO YOUR FUCKING JOBS!! #FuckTrump And Who Votes For Him!	HOF	PRFN	TIN
38	en	@DnaZeeNews I am a Hindu Hindusthani Nationalist.Down with Mahuwa Moitra and her supporters.She is parrot of	@DnaZeeNews I am a Hindu Hindusthani Nationalist.Down with Mahuwa Moitra and her	NOT	NONE	NONE
39	en	What r Muslims in India doing on #TripleTalaqBill ? Their lack of knowledge about Islam clearly exposed. Is there none	What r Muslims in India doing on #TripleTalaqBill ? Their lack of knowledge about Islam clearly exposed.	HOF	HATE	TIN
40	en	When you enter in, this is gonna be really funny to read hark ! It's like you let negative text lead the whole post	When you enter in, this is gonna be really funny to read hark ! It's like you let negative text lead the	NOT	NONE	NONE

By employing this approach, we were able to translate the entire dataset efficiently and effectively. The built-in translation formula in Google Sheets saved us considerable time and effort compared to the initial Python program.

The utilization of Google Sheets and its translation formula allowed us to overcome the challenges associated with working with a mixed-language dataset. It enabled us to convert the text to a single language, in this case, English, making it compatible with the subsequent steps of our sentiment analysis.

Sentiment Analysis Model:

For our sentiment analysis task, we chose to leverage the power of the BERT-base pretrained model from Hugging Face. BERT (Bidirectional Encoder Representations from Transformers) has been widely recognized for its state-of-the-art performance in various natural language processing tasks, including sentiment classification.

The BERT-base model we used is specifically designed for sentiment analysis on product reviews in six languages: English, Dutch, German, French, Spanish, and Italian. It predicts the sentiment of a review by assigning a number of stars ranging from 1 to 5.



The BERT-based sentiment analysis model was trained using a large number of product reviews in each of the six languages. Here's an overview of the training data:

Language | Number of Reviews

English | 150,000 Dutch | 80,000 German | 137,000 French | 140,000 Italian | 72,000 Spanish | 50,000

During the fine-tuning process, we evaluated the performance of the model on a held-out dataset consisting of 5,000 product reviews in each language. The model achieved the following accuracy metrics:

Language | Accuracy (exact) | Accuracy (off-by-1)

English | 67% | 95% Dutch | 57% | 93% German | 61% | 94% French | 59% | 94% Italian | 59% | 95% Spanish | 58% | 95%

Now let's dive into the code used to implement the sentiment analysis model in Jupyter Notebook.

First, we need to load the tokenizer associated with the BERT-base model. The tokenizer is responsible for breaking down the text into tokens that the model can understand. We accomplish this using the `AutoTokenizer.from_pretrained()` function.

Next, we load the pre-trained sentiment classification model itself. In our case, it is the BERTbase model for sequence classification, which has been fine-tuned specifically for sentiment analysis. We utilize the `AutoModelForSequenceClassification.from_pretrained()` function to load the model.

With both the tokenizer and the sentiment classification model loaded, we can now proceed to the sentiment scoring process. The sentiment scoring function takes in a review as input and calculates the sentiment score using the BERT-base model. We apply this function to each review in the translated text column of our dataset using the

```
df['TRANSLATED'].apply(lambda x: sentiment_score(x[:512])) syntax.
```

The sentiment scores are then added as a new 'sentiment' column in the DataFrame, allowing us to further analyze and visualize the sentiment analysis results.

By leveraging the BERT-base pre-trained model and performing sentiment analysis on our translated dataset, we can gain valuable insights into the sentiment expressed in the text. The model's ability to accurately classify sentiments provides us with reliable results for our sentiment analysis task.


```
In [4]: !pip install torch
!pip install requests
!pip install numpy
!pip install pandas
!pip install transformers
```

```
Requirement already satisfied: torch in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (2.0.1+cu118)
Requirement already satisfied: filelock in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from torch) (3.9.0)
Requirement already satisfied: typing-extensions in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from torch) (4.5.0)
Requirement already satisfied: sympy in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from torch) (1.11.1)
Requirement already satisfied: networkx in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from torch) (3.0)
Requirement already satisfied: Jinja2 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from torch) (3.1.2)
Requirement already satisfied: MarkupSafe>=2.0 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from Jinja2->torch) (2.1.3)
Requirement already satisfied: mpmath>=0.19 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from sympy->torch) (1.2.1)
Requirement already satisfied: requests in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests) (3.2.0)
Requirement already satisfied: idna<4,>=2.5 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests) (2.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests) (1.26.16)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests) (2023.5.7)
Requirement already satisfied: numpy in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (1.24.3)
Requirement already satisfied: pandas in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (2.0.3)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from pandas) (2023.3)
Requirement already satisfied: tzdata>=2022.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from pandas) (2023.3)
Requirement already satisfied: numpy>=1.21.0 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from pandas) (1.24.3)
Requirement already satisfied: six>=1.5 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
Requirement already satisfied: transformers in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (4.30.2)
Requirement already satisfied: filelock in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (3.9.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.14.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (0.16.4)
Requirement already satisfied: numpy>=1.17 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (1.24.3)
Requirement already satisfied: packaging>=20.0 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (23.1)
Requirement already satisfied: pyyaml>=5.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (6.0)
Requirement already satisfied: regex!=2019.12.17 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (2023.6.3)
Requirement already satisfied: requests in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (2.31.0)
Requirement already satisfied: tokenizers!=0.11.3,<0.14,>=0.11.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (0.13.3)
Requirement already satisfied: safetensors>=0.3.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (0.3.1)
Requirement already satisfied: tqdm>=4.27 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from transformers) (4.65.0)
Requirement already satisfied: fsspec in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from huggingface-hub<1.0,>=0.14.1->transformers) (2023.6.0)
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from huggingface-hub<1.0,>=0.14.1->transformers) (4.5.0)
Requirement already satisfied: colorama in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from tqdm>=4.27->transformers) (0.4.6)
Requirement already satisfied: charset-normalizer<4,>=2 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests->transformers) (3.2.0)
Requirement already satisfied: idna<4,>=2.5 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests->transformers) (2.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests->transformers) (1.26.16)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\shaik\appdata\local\programs\python\python311\lib\site-packages (from requests->transformers) (2023.5.7)
```

```
In [2]: import torch import requests import re import numpy as np import pandas
as pd from transformers import AutoTokenizer,
AutoModelForSequenceClassification
```

```
In [3]: # Check if a GPU is available and set the device accordingly
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

print(device)
```

cuda

```
In [4]: !nvidia-smi
```

Thu Jul 20 00:23:03 2023

```
+-----+
| NVIDIA-SMI 536.67                Driver Version: 536.67        CUDA Version: 12.2
|
+-----+-----+-----+-----+
| GPU  Name                      TCC/WDDM  | Bus-Id      Disp.A | Volatile Uncorr.
ECC |
+-----+-----+-----+-----+
| Fan  Temp   Perf              Pwr:Usage/Cap |      Memory-Usage | GPU-Util  Compute M. |
|              |              |          |              |      |          |
+-----+-----+-----+-----+-----+
|   0  NVIDIA GeForce RTX 3070 ... WDDM      | 00000000:01:00:0  On  |              N/A  | | | |
| N/A   41C    P8              14W / 140W   | 1156MiB / 8192MiB |      0%      Defa ult |
|              |              |          |              |      |          |
+-----+-----+-----+-----+-----+

+-----+
| Processes:
| GPU  GI   CI        PID   Type   Process name                      GPU Mem ory |
|      ID   ID              |              | Usage      |
+-----+-----+-----+-----+-----+
|   0  N/A  N/A       4928    C+G    ...GeForce Experience\NVIDIA Share.exe  N/A  |
|
|   0  N/A  N/A      14224    C+G    ...on\114.0.1823.82\msedgewebview2.exe  N/A  |
|   0  N/A  N/A      14236    C+G    ...crosoft\Edge\Application\msedge.exe  N/A  |
|   0  N/A  N/A      14792    C+G    ...5n1h2txyewy\ShellExperienceHost.exe  N/A  |
|   0  N/A  N/A      23516    C+G    ...oogle\Chrome\Application\chrome.exe  N/A  |
|   0  N/A  N/A      24320    C+G    ...GeForce Experience\NVIDIA Share.exe  N/A  |
|   0  N/A  N/A      27244    C+G    ...__8wekyb3d8bbwe\Notepad\Notepad.exe  N/A  |
|   0  N/A  N/A      27996    C+G    ...__8wekyb3d8bbwe\WindowsTerminal.exe  N/A  |
|   0  N/A  N/A      31836    C+G    ...2txyewy\StartMenuExperienceHost.exe  N/A  |
|   0  N/A  N/A      32136    C+G    ...ekyb3d8bbwe\PhoneExperienceHost.exe  N/A  |
|   0  N/A  N/A      33656    C+G    ...siveControlPanel\SystemSettings.exe  N/A  |
|
|   0  N/A  N/A      34468    C+G    C:\Windows\explorer.exe                  N/A  |
|   0  N/A  N/A      34840    C+G    ...CBS_cw5n1h2txyewy\TextInputHost.exe  N/A  |
|   0  N/A  N/A      34992    C+G    ...nt.CBS_cw5n1h2txyewy\SearchHost.exe  N/A  |
+-----+
+-----+
```

```
In [5] # Read the Excel file
input_file_path =
r"C:\Users\shaik\OneDrive\Desktop\Jupyter\Hasoc_Dataset_Translate df =
pd.read_excel(input_file_path, sheet_name="Sheet1")
```

```
In [6]: # Print the information about the
DataFrame print("\nInfo:")
print(df.info())

# Print the summary statistics of the
DataFrame print("\nDescribe:")
print(df.describe())

# Print the tail of the
DataFrame print("\nTail:")
print(df.tail())

# Print the shape of the
DataFrame print("\nShape:")
print(df.shape)
```

```
Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17657 entries, 0 to 17656
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  ---
0    INDEX      17657 non-null  int64
1    LANG        17657 non-null  object
2    TEXT        17657 non-null  object
3    TRANSLATED  17657 non-null  object
4    TYPE_A      17657 non-null  object
5    TYPE_B      17657 non-null  object  6    TYPE_C      12988 non-null  object dtypes: int64(1), object(6) memory usage: 965.7+ KB
None
```

```
Describe:
              INDEX
count  17657.000000 mean
8829.000000 std      5097.281187 min
1.000000 25%      4415.000000
50%      8829.000000 75%
13243.000000 max    17657.000000

Tail:
              INDEX LANG                                TEXT \ 17652  17653   de  #Erdogan ruft seine Mops wegen
Untauglichkeit ...
17653      17654   de  Ich würde der Linken Chefin Kipping die 1.000 ...
17654      17655   de  Ausländer sind Top & Trumpf Schon aufgefal...
17655      17656   de  Wie der Regierungsjet die Regierungsflieger wi... 17656  17657   de  #guteMütter sind heutzutage gut
ausgebildete F...
```

```
              TRANSLATED TYPE_A TYPE_B TYPE_C
17652  #Erdogan calls his pug back for unsuitability....  HOF  OFFN  NaN
17653  I would give the left boss Kipping the 1,000 e...  HOF  OFFN  NaN
17654  Foreigners are TOP & AMP; Trumpf noticed that ...  HOF  OFFN  NaN
17655  How the government jet reflects the government...  HOF  OFFN  NaN  17656  #Goods nowadays are well -trained specialists
...  NOT  NONE  NaN
Shape:
(17657, 7)
```

In [8]:

```
# Define the sentiment score function
def sentiment_score(review, index):
    tokens = tokenizer.encode(review, return_tensors='pt').to(device) # Move
input    if index <= 5:
        print(f"Tokens ({index}):", tokens)    result =
model(tokens.to(device)) # Move model to GPU    if
index <= 5:
        print(f"Model output ({index}):",
result)    return
int(torch.argmax(result.logits)) + 1
# Convert reviews to a list reviews = df['TRANSLATED'].tolist()

# Load the tokenizer and model onto the GPU
tokenizer = AutoTokenizer.from_pretrained('nlptown/bert-base-multilingual-
uncased-s model =
AutoModelForSequenceClassification.from_pretrained('nlptown/bert-base-multi
```

In [9]:

```
# Print the model architecture
print("\nModel Architecture:")
print(model)
```

```
Model Architecture:
BertForSequenceClassification(
  (bert): BertModel(
    (embeddings): BertEmbeddings(
      (word_embeddings): Embedding(105879, 768, padding_idx=0)
      (position_embeddings): Embedding(512, 768)
      (token_type_embeddings): Embedding(2, 768)
      (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): BertEncoder(
      (layer): ModuleList(
        (0-11): 12 x BertLayer(
          (attention): BertAttention(
            (self): BertSelfAttention(
              (query): Linear(in_features=768, out_features=768, bias=True)
              (key): Linear(in_features=768,
out_features=768, bias=True)
              (value): Linear(in_features=768, out_features=768, bias=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (output): BertSelfOutput(
              (dense): Linear(in_features=768, out_features=768, bias=True)
              (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
          )
          (intermediate): BertIntermediate(
            (dense): Linear(in_features=768, out_features=3072, bias=True)
            (intermediate_act_fn): GELUActivation()
          )
          (output): BertOutput(
            (dense): Linear(in_features=3072, out_features=768, bias=True)
            (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
    )
    (pooler): BertPooler(
      (dense): Linear(in_features=768, out_features=768, bias=True)
      (activation): Tanh()
    )
  )
  (dropout): Dropout(p=0.1, inplace=False)
  (classifier): Linear(in_features=768, out_features=5, bias=True) )
```

```
In [10]: # Calculate sentiment scores for each review
df['sentiment'] = df.apply(lambda row: sentiment_score(row['TRANSLATED'],
row['INDE
# Print the top 5 reviews, tokens, and model output print("\nTop 5 Reviews,
Tokens, and Model Output:") for _, row in df.head(5).iterrows():
print("Review:", row['TRANSLATED']) print("Tokens:",
tokenizer.encode(row['TRANSLATED'], return_tensors='pt').to(de
```

```
print("Model output:", model(tokenizer.encode(row['TRANSLATED'], return_tensors
print("\n"))
```

```
Tokens (1): tensor([[ 101, 11803, 52350, 10933, 10320, 21249, 10285, 11901, 10114, 17019,

20664, 119, 14048, 112, 161, 18469, 14540, 131, 120, 120,
11312, 10481, 119, 14879, 120, 123, 10251, 10311, 11301, 10661,
10167, 10995, 108, 56032, 108, 41406, 102]], device='cuda:0')
Model output (1): SequenceClassifierOutput(loss=None, logits=tensor([[ 1.0675, 0.38
51, -0.1189, -0.5166, -0.7552]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)
Tokens (2): tensor([[ 101, 47530, 137, 64647, 42871, 36057, 20894, 137, 14015, 21470,
84404, 44574, 10158, 137, 23172, 93300, 12947, 137, 15672, 62884,
14487, 10139, 39684, 10763, 13627, 10114, 66681, 10114, 12548, 10114,
12696, 108, 14015, 12947, 143, 108, 47300, 38709, 108, 56032,
11623, 24414, 21567, 10111, 108, 87890, 67919, 22734, 49669, 102]], device='cuda:0')
Model output (2): SequenceClassifierOutput(loss=None, logits=tensor([[ -1.1487, -1.28
85, -0.1772, 0.9982, 1.1996]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)
Tokens (3): tensor([[ 101, 69338, 119, 10127, 10372, 25165, 61330, 10285, 106, 10855,
67922, 10320, 37079, 42279, 10688, 106, 106, 106, 38881, 12696,
143, 27144, 12908, 10768, 106, 46921, 11153, 13967, 67351, 106,
102]], device='cuda:0')
Model output (3): SequenceClassifierOutput(loss=None, logits=tensor([[ 2.6249, -0.18
02, -0.8666, -1.4106, 0.0873]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)
Tokens (4): tensor([[ 101, 11876, 13111, 11250, 10108, 11435, 84684, 19918, 10435, 13157,
112, 11168, 113, 10345, 27485, 114, 13222, 10601, 87982, 15866,
51905, 53325, 10163, 10139, 126, 10868, 10139, 72393, 18122, 10285,
12027, 10108, 10103, 13111, 43757, 10171, 24365, 16751, 11082, 10110,
46348, 32492, 10163, 13157, 119, 108, 10345, 10136, 108, 10345,
27485, 102]], device='cuda:0')
Model output (4): SequenceClassifierOutput(loss=None, logits=tensor([[ 1.1311, 0.34
19, -0.2421, -0.5088, -0.5824]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)
Tokens (5): tensor([[ 101, 11896, 11636, 118, 24357, 12947, 14958, 29152, 131, 13202,
50528, 112, 25204, 26284, 49563, 10261, 41294, 112, 11379, 14540,
131, 120, 120, 10561, 119, 11896, 119, 10241, 120, 11636,
120, 10419, 118, 11915, 118, 10912, 118, 41743, 11124, 11444,
47173, 11301, 100, 10572, 10127, 143, 23145, 33711, 10421, 10114,
13689, 102]], device='cuda:0')
Model output (5): SequenceClassifierOutput(loss=None, logits=tensor([[ 2.0429, 0.49
56, -0.1504, -0.9026, -1.1270]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)

Top 5 Reviews, Tokens, and Model Output:
Review: Young Eritreans are risking death to migrate. Here's why https://wef.ch/2Gz5 xK0 #Refugees #Migration
Tokens: tensor([[ 101, 11803, 52350, 10933, 10320, 21249, 10285, 11901, 10114, 17019, 9,
20664, 119, 14048, 112, 161, 18469, 14540, 131, 120, 120,
11312, 10481, 119, 14879, 120, 123, 10251, 10311, 11301, 10661,
10167, 10995, 108, 56032, 108, 41406, 102]], device='cuda:0')
Model output: SequenceClassifierOutput(loss=None, logits=tensor([[ 1.0675, 0.3851, -0.1189, -0.5166, -0.7552]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)

Review: Thanks @RMNetwork @LewishamSanctuary @valewishesam @kevinbonavia for bringing u s together to discuss to how to make #Lewisham a #sanctuaryborough #RefugeesAreWelco me #RefugeeWeek2019
Tokens: tensor([[ 101, 47530, 137, 64647, 42871, 36057, 20894, 137, 14015, 21470, 0,
84404, 44574, 10158, 137, 23172, 93300, 12947, 137, 15672, 62884,
14487, 10139, 39684, 10763, 13627, 10114, 66681, 10114, 12548, 10114,
12696, 108, 14015, 12947, 143, 108, 47300, 38709, 108, 56032,
11623, 24414, 21567, 10111, 108, 87890, 67919, 22734, 49669, 102]], device='cuda:0')
Model output: SequenceClassifierOutput(loss=None, logits=tensor([[ -1.1487, -1.2885,
-0.1772, 0.9982, 1.1996]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)

Review: Fuck. Is this really happening! You guys are awesome!!! Please make a vlog t here! Copyright my ass ICC!
Tokens: tensor([[ 101, 69338, 119, 10127, 10372, 25165, 61330, 10285, 106, 10855, 5,
67922, 10320, 37079, 42279, 10688, 106, 106, 106, 38881, 12696,
143, 27144, 12908, 10768, 106, 46921, 11153, 13967, 67351, 106,
102]], device='cuda:0')
Model output: SequenceClassifierOutput(loss=None, logits=tensor([[ 2.6249, -0.1802,
-0.8666, -1.4106, 0.0873]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)

Review: Former Vice President of Aligarh Muslim University Students' Union (AMUSU) H amza Sufiyan rusticated for 5 years for allegedly storming office of the Vice Chance llor with armed outsiders and
debarred students. #AMU #AMUSU
Tokens: tensor([[ 101, 11876, 13111, 11250, 10108, 11435, 84684, 19918, 10435, 13157, 7,
112, 11168, 113, 10345, 27485, 114, 13222, 10601, 87982, 15866,
51905, 53325, 10163, 10139, 126, 10868, 10139, 72393, 18122, 10285,
12027, 10108, 10103, 13111, 43757, 10171, 24365, 16751, 11082, 10110,
46348, 32492, 10163, 13157, 119, 108, 10345, 10136, 108, 10345,
27485, 102]], device='cuda:0')
Model output: SequenceClassifierOutput(loss=None, logits=tensor([[ 1.1311, 0.3419,
-0.2421, -0.5088, -0.5824]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)

Review: BBC News - Waltham Forest Pride: Police probe 'homophobic abuse' video http s://www.bbc.com/news/uk-england-london-49147605 - She is a disgrace to Islam
Tokens: tensor([[ 101, 11896, 11636, 118, 24357, 12947, 14958, 29152, 131, 13202, 2,
50528, 112, 25204, 26284, 49563, 10261, 41294, 112, 11379, 14540,
131, 120, 120, 10561, 119, 11896, 119, 10241, 120, 11636,
120, 10419, 118, 11915, 118, 10912, 118, 41743, 11124, 11444,
47173, 11301, 100, 10572, 10127, 143, 23145, 33711, 10421, 10114,
13689, 102]], device='cuda:0')
Model output: SequenceClassifierOutput(loss=None, logits=tensor([[ 2.0429, 0.4956, -0.1504, -0.9026, -1.1270]], device='cuda:0',
grad_fn=<AddmmBackward0>), hidden_states=None, attentions=None)
```


In [11]:

```
# Print the top 5 reviews
print("\nTop 10 Sentiment:")
print(df['sentiment'].head(10))
```

Top 10 Sentiment:

```
0    1
1    5
2    1
3    1
4    1
5    1
6    1
7    5
8    1
9    2
```

Name: sentiment, dtype: int64

In [13]:

```
# Save the DataFrame with sentiment scores to output_file.xlsx
output_file_path =
r"C:\Users\shaik\OneDrive\Desktop\Jupyter\Hasoc_Dataset_Analysis
df.to_excel(output_file_path, index=False) print(f"DataFrame saved to
{output_file_path}")
```

DataFrame saved to

C:\Users\shaik\OneDrive\Desktop\Jupyter\Hasoc_Dataset_Analysis.xlsx

Evaluation and Analysis:

Now that we have performed sentiment analysis on our dataset and obtained sentiment scores, let's proceed with the evaluation and analysis of the results.

To facilitate a better understanding of the sentiment scores, we replace the numeric values with descriptive sentiment categories. This transformation provides more intuitive labels for the sentiment categories, allowing us to interpret the results more easily.

After replacing the numeric values, we can evaluate the sentiment categories and subcategories based on the sentiment scores. This evaluation provides us with insights into the distribution of sentiments within the dataset and helps us analyze sentiment patterns across different categories.

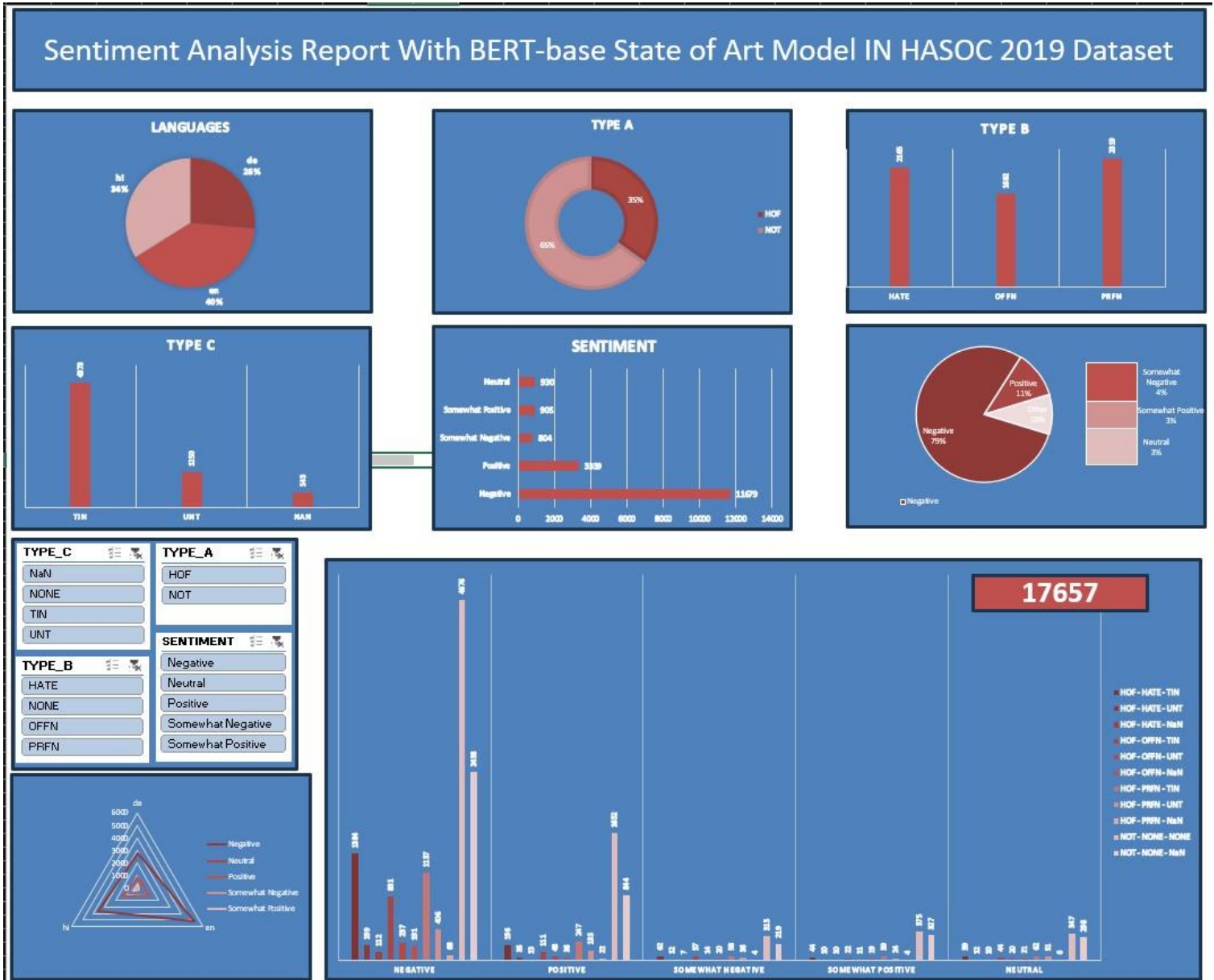
INDEX	LANG	TEXT	TRANSLATED	TYPE_A	TYPE_B	TYPE_C	SENTIMENT
1657	de	#GuteMüter sind heutzutage gut ausgebildete Fachkräfte, die den derzeitigen Fachkräftemangel aufhalten! Durch ihre Einstellung zur Familie, sind Sie Prädestiniert, Führungsaufgaben zu erfüllen! Sie kommen, wenn überhaupt, dem Wunsch des Kinderkriegens nach und...	#Goods nowadays are well-trained specialists who stop the current shortage of skilled workers! Due to their attitude to the family, they are predestined to perform management tasks! If at all, they comply with the desire of child war and...	NOT	NONE	NaN	Somewhat Positive
1656	de	Wie der Regierung ist die Regierungsflietler wieder spiegel. Hatte Probleme mit der Steuerung, legte eine Buchführung hin, kippte nach rechts ab, linker Sumpf, rechter Flügel irreparabel beschädigt. Ein Spiegelbild #Deutschlands, #Berlin's. Der #Ekelwifred https://cof9v7f8z6g	How the government jet reflects the government pilots. Had problems with the control, put on a crash landing, tilted to the right, left & amp; right wing irreparably damaged. A reflection #Germany's, #Berlin's. The #Ekelwifred https://cof9v7f8z6g	HOF	OFFN	NaN	Negative
1655	de	Ausländer sind Top & amp; Trumpf! Schon aufgefalten, dass nur Ausländer was zu sagen haben. Oliver Jarich (Philippinen), Oliver Flesch (Spanien), Hagen's (Ungarn), Ich (Bulgarien, Niederlande) usw. uaf. Nur vers Land verlassen hat, machts Mal auf, und Ausländer in D. @ #Ekelwifred	Foreigners are TOP & AMP, Trumpf noticed that only foreigners have something to say. Oliver Jarich (Philippines), Oliver Flesch (Spain), Hagen's (Hungary), I (Bulgaria, Netherlands) etc. etc. Only who has left land, make up, and foreigners in D. @ #Ekelwifred	HOF	OFFN	NaN	Negative
1654	de	Ich würde der Linken Chef'n Kipping die 1000 Euro Arbeitslosigkeit geben. Die wären bedeutend besser angelegt als die Diät. @ #Der #Ekelwifred, alias das A. Loch #Witewka, Papa #Wilberg von der Capitol und O. #Witke #DerWahrePräsident, auch von #Deutschland https://cof9v7f8z6g	I would give the left boss Kipping the 1000 euros unemployment benefit. They would be significantly better than the diet. @ The #Ekelwifred, alias das A. Loch #Witewka, papa #Wilberg from the Capitol and O. #Witke #Berlin #Lannweg #musweg https://cof9v7f8z6g	HOF	OFFN	NaN	Somewhat Negative
1653	de	#Erdogan ruft seine Mops wegen Unstabilität zurück. Nur wurde bis jetzt noch keiner zurückgeben. @ #Der #Ekelwifred, alias das A. Loch #Witewka, Papa #Wilberg von der Capitol & amp; O. #Witke #DerWahrePräsident, auch von #Deutschland & amp; Kemal-ismus-Fan. https://cof9v7f8z6g	#Erdogan calls his pug back for unstability. Only so far nobody has been returned. O. #Witke #Derwahredress, also from #Deutschland & amp; Kemalism fan. https://cof9v7f8z6g	HOF	OFFN	NaN	Negative
1652	de	Defiedas große E ruft seine Mops wegen Unstabilität zurück. Nur wurde bis jetzt noch keiner zurückgeben. #Witayawilwan, alias #Ekelwifred, A. Loch #Witewka, Papa #Wilberg von der Capitol und O. #Witke #DerWahrePräsident, auch von #Deutschland https://cof9v7f8z6g	The large E is recalling its pugs due to unstability. Only so far nobody has been returned. #Witayawilwan, alias #Ekelwifred, A. Loch #Witewka, Papa #Wilberg from the Capitol and O. #Witke #of the HUCHSPREAGE, also from #Deutschland https://cof9v7f8z6g	HOF	OFFN	NaN	Negative
1651	de	Besser ein rechter #Komiker als #Führungsspitze der #Ukraine. Als ungewillig #falsche Komiker:innen, wie in den meisten anderen #Ländern. @ Das A. Loch #Witewka, alias der #Ekelwifred, Papa #Wilberg von der Capitol, #Witayawilwan und O. #Witke https://cof9v7f8z6g	Better a #rechter #Komiker than #Leadership of #Ukraine. As involuntary #f's comedian, as in most other #countries. @ The A. Loch #Witewka, alias the #Ekelwifred, Papa #Wilberg von der Capitol, #Witayawilwan and O. #Witke https://cof9v7f8z6g	HOF	OFFN	NaN	Somewhat Negative
1650	de	Das hoffen wir alle - zumindest die GUTEN unter uns - Sogar der #Ekelwifred (meiner einer). #Tag24 top, wie RTL. @ #findbeco endlich, sucht endlich nach einer lebenden #trebecca an den richtigen Orten. #freeforian #schwagerwas #reden #Presse #rt https://cof9v7f8z6g	We all hope that - at least the good ones among us - even the #Ekelwifred (one of them). #TAG24 TOP, like RTL. @ #findbeco finally, finally looking for a living #trebecca in the right places. #freeforian #Schwagerwas #reden #Presse #rt https://cof9v7f8z6g	HOF	OFFN	NaN	Positive
1649	de	@Uwe_Junge_Mdl Vergiss es... von den Trotteln überlegt keiner mehr! Nach ausnehmen der Bürger ist im Kopf endet!	@Uwe_Junge_mdl forget it... nobody thinks of the drottel! After taking out the citizens, the end is in the head!	HOF	PRFN	NaN	Negative
1648	de	@Uwe_Junge_Mdl @	@Uwe_Junge_mdl @	NOT	NONE	NaN	Positive
1647	de	@Uwe_Junge_Mdl Ch die AfD als Vertreter der Homosexuellen? Muss ich mir merken.	@Uwe_Junge_mdl ch the AfD as a representative of homosexuals? I have to remember.	NOT	NONE	NaN	Negative
1646	de	@Uwe_Junge_Mdl Ihr streitet es zwar immer ab, aber gibt es in D eine Gruppierung, die mehr gegen Homosexuelle, Juden, Muslime, usw. hetzen als die AfD?	@Uwe_Junge_mdl You always deny it, but is there a group in D that rush more against homosexuals, Jews, Muslims, etc. than the AfD?	NOT	NONE	NaN	Negative
1645	de	@Uwe_Junge_Mdl @Beatrice_Voroch @Evrd @Nicole_Hoehst @Volker_Beck Meine Rede sei langem: https://cof9v7f8z6g	@Uwe_Junge_mdl @Beatrice_voroch @Evrd @Nicole_hoehst @Volker_beck My speech has long been: https://cof9v7f8z6g	NOT	NONE	NaN	Positive
1644	de	@Uwe_Junge_Mdl Ach Junge, hast Du Angst, dass Euch radikale Muslime die Homophobie wegnehmen? Wenn ihr an die Macht kommen solltet (was nie passieren wird), ergeht es Schwulen und Lesben auch nicht besser https://cof9v7f8z6g	@Uwe_Junge_mdl Oh boy, are you afraid that radical Muslims will take away homophobia? If you should get to power (which will never happen), gays and lesbians do not do better https://cof9v7f8z6g	NOT	NONE	NaN	Negative
1643	de	Wer mir mit ihrer Politik die und das Liebeste nimm - Aber nicht nur die Autos - Der bekommt von mir auch nur noch den A. hilt. @ #Der #Ekelwifred, alias das A. Loch #Witewka, Papa #Wilberg von der Capitol, #Witayawilwan und O. #Witke https://cof9v7f8z6g	Whoever takes me and my loved one - but not just the cars - only gets me from me. @ The #Ekelwifred, alias the A. Loch #Witewka, Papa #Wilberg from the Capitol, #Witayawilwan and O. #Witke https://cof9v7f8z6g	HOF	OFFN	NaN	Neutral
1642	de	Ich mag Kim Jong Un. Denn der liebt die Umwelt, fährt mit dem Zug, teils sogar mit sauberem Ökostrom auf die Leute zu, während die meisten Grünen immer fort-fahrend das Umfeld verpesten. Und besser reiten kann er auch. @ Der ekige #Grüne #Ekelwifred https://cof9v7f8z6g	I like Kim Jong un. Because he loves the environment, drives by train, sometimes even with clean green electricity towards people, while most Greens always pollute the environment on the way. And he can ride better too. @ The disgusted #Greens #Ekelwifred https://cof9v7f8z6g	HOF	OFFN	NaN	Somewhat Positive
1641	de	Muss ich mich jetzt mit der Mehrheit oder der Minderheit solidarisieren? Ne aber wer trägt denn die ganze Schuld? Ja wohl die, die nur reden und nichts gegen Illegalität Sumpf, ausfahende Kriminalität tun, und beides legalisieren wollen. @ #Der #Ekelwifred https://cof9v7f8z6g	Do I have to solidarize myself with the majority or the minority now? No but who is the fault? Yes probably those who only talk and nothing against illegality & amp; Do sprawling crime and want to legalize both. @ The #Ekelwifred https://cof9v7f8z6g	HOF	OFFN	NaN	Negative
1640	de	Das wird auch langsam Zeit. Denn wenn eins das Land nicht brauchen kann, dann sind es welche die es wagen die Wahrheit zu sagen. Ob nun dieser Abgründe, der sich Annaeßende oder die Einzelsitzende. Und die kompletten Alternativvollen. @ Der #Ekelwifred https://cof9v7f8z6g	It is slowly getting time. Because if one cannot need the country, then there are those who dare to say the truth. Whether this ravine, the intelligent or the individual sitting end. And the complete alternative. @ The #Ekelwifred https://cof9v7f8z6g	HOF	OFFN	NaN	Somewhat Positive
1639	de	Wenn er wohl auf die Füße getreten ist oder sollte das einem nicht zu denken geben? Mein aufrichtiges Beileid allen oder einigen Opfern von Gewalttaten. Je nach Ideologie oder ob ich dafür jetzt in den Himmel oder die Hölle komme. @ #Der #Ekelwifred https://cof9v7f8z6g	If you probably stepped on your feet or should that not give you think? My sincere condolences to all or some victims of acts of violence. Depending on the ideology or whether I now come to heaven or hell. @ The #Ekelwifred https://cof9v7f8z6g	HOF	OFFN	NaN	Positive

Let's delve into the findings from our evaluation:

		TYPE_A		HOF			TYPE_A		HOF		
Row Labels		Count of LANG		Row Labels		Count of TYPE_A		Row Labels		Count of TYPE_B	
		de				HOF				TIN	
		en				NOT				UNT	
		hi				Grand Total				NaN	
		Grand Total				17657				Grand Total	
										Grand Tota	
										6166	

Analysis Dashboard:

As part of this end-to-end project, I have created an analysis dashboard to visualize and explore the sentiment analysis results. The dashboard provides an interactive interface that allows you to gain deeper insights into the sentiment patterns within the dataset. Let's take a closer look at the features and components of the dashboard.



The dashboard consists of several interactive visualizations that enable you to explore the sentiment analysis results from different perspectives. These visualizations provide valuable insights into sentiment distribution, sentiment by language, sentiment by sub-category, and more.

```
In [3]: from ipywidgets import Video
# Define the video file path
video_file_path = r"C:\Users\shaik\OneDrive\Desktop\Jupyter\Dash Board - Made
with

# Display the video
Video.from_file(video_file_path)
```

```
Out[3]: Video(value=b'\x00\x00\x00
ftypisom\x00\x00\x02\x00isomiso2avc1mp41\x00\x00\x00\x0
8free...')
```

Here are the key features and components of the analysis dashboard:

1. Language Pie Chart: This chart displays the distribution of languages within the dataset. By interacting with this chart, you can gain insights into the proportion of texts in each language.
2. Type A Donut Chart: This chart represents the distribution of Type A sub-categories, which focus on hate speech and offensive language identification. By hovering over the segments, you can view the percentage and count of each sub-category.
3. Type B and Sentiment Bar Chart: This chart provides a visual representation of the sentiment distribution within different Type B sub-categories. The sentiment values are displayed on the y-axis, while the bars represent the count of texts. This visualization allows you to compare sentiment distribution across different sub-categories.
4. Type A, B, C, and Sentiment Filter Slicers: These interactive filters allow you to select specific sentiment categories (Type A, B, or C) and sentiment levels. By adjusting these filters, you can dynamically update the visualizations and explore sentiment patterns within different sub-categories.
5. Sentiment Percentage Pie Chart: This chart displays the distribution of sentiment percentages based on the selected filters. It provides a visual representation of sentiment percentages for the chosen sentiment categories and levels.
6. Combined Bar Chart: This chart combines all Type A, B, and C sub-categories and sentiment levels into a single visualization. It allows you to compare sentiment distribution across different sub-categories and sentiment levels in a consolidated manner.
7. Radial Chart: This chart visualizes the relationship between language (as the axis) and sentiment (as the range). By exploring this chart, you can identify sentiment patterns within different languages and understand how sentiments vary across languages.
8. Index Count Total Text: This visualization presents the total count of texts available in the dataset. It provides an overview of the dataset's size and serves as a reference for the sentiment analysis results.

By interacting with the dashboard, you can explore and analyze the sentiment analysis results, identify trends, and make data-driven decisions based on the insights gained.

The analysis dashboard offers an intuitive and user-friendly interface for data exploration and visualization. It allows you to delve into the sentiment distribution, sentiment by language, and sentiment by sub-category, enabling you to uncover valuable insights and patterns within the dataset.

With the interactive features and filtering options, you have the flexibility to customize the visualizations based on your specific interests and research questions. This empowers you to explore the sentiment analysis results in a way that best suits your needs and supports informed decision-making.

Overall, the analysis dashboard provides a comprehensive and interactive environment for exploring the sentiment analysis outcomes, making it an invaluable tool for researchers, analysts, and decision-makers in various domains.

Take your time to explore the dashboard, interact with the visualizations, and uncover meaningful insights from the sentiment analysis results.

Conclusion and Next Steps:

In conclusion, we have covered the key steps involved in sentiment analysis using Jupyter Notebook in Colab. Let's summarize what we've learned:

1. **Data Collection:** We researched and explored various sources, including AI-powered web searches and Kaggle datasets, to find relevant datasets. We selected the Hasoc 2019 dataset as our primary dataset, which provided labeled text in English, Hindi, and German.
2. **Data Preprocessing:** We used Excel to format and combine the datasets, merged them into a cohesive dataset, and performed necessary data cleaning steps such as splitting the ID column into a language column, dropping unnecessary columns, and removing empty rows, duplicates, and NaN values.
3. **Exploratory Data Analysis (EDA):** We conducted EDA on the dataset and gained insights into language distribution and label distribution. We observed interesting patterns, such as the percentage of unique data, language distribution percentages, and label distribution percentages.
4. **Language Translation:** We addressed the challenge of working with a mixed-language dataset by using Google Sheets and its built-in translation formula. We created a new column for translation and performed batch translation efficiently using the language and text columns.
5. **Sentiment Analysis Model:** We leveraged the BERT-base pre-trained model from Hugging Face for sentiment analysis. We loaded the tokenizer and the pre-trained

sentiment classification model in Jupyter Notebook and applied the sentiment scoring process to assign sentiment scores to each review.

6. Evaluation and Analysis: We evaluated the sentiment analysis results by replacing numeric values with descriptive sentiment categories. We analyzed sentiment categories and sub-categories based on sentiment scores and derived valuable insights from the evaluation.
7. Analysis Dashboard: We created an analysis dashboard to visualize and explore the sentiment analysis results. The dashboard included interactive visualizations, such as language pie charts, type A donut charts, type B and sentiment bar charts, and more. It provided a user-friendly interface for data exploration, trend identification, and datadriven decision-making.

Through this project, we have gained valuable insights into sentiment analysis, including data collection, preprocessing, model selection, evaluation, and visualization. We have explored the dynamics of sentiment within different languages and sub-categories, providing us with a comprehensive understanding of sentiment patterns.

As the next step, I encourage you to explore the provided Jupyter Notebook and experiment with different datasets. Sentiment analysis has applications in various domains, such as social media monitoring, customer feedback analysis, and more. By further exploring the field of sentiment analysis, you can uncover deeper insights, develop advanced models, and contribute to cutting-edge research.

Thank you for joining me in this end-to-end project on sentiment analysis using Jupyter Notebook in Colab. I hope you found this project informative and insightful. If you have any questions or would like to discuss further, please feel free to ask.

Happy exploring and analyzing sentiments from text data!

References and Resources:

<https://www.perplexity.ai/search> <https://huggingface.co/nlptown/bert-base-multilingual-uncased-sentiment>
<https://hasocfire.github.io/hasoc/2019/index.html>
<https://arxiv.org/pdf/1909.12642.pdf> <https://chat.openai.com>
<https://docs.google.com/spreadsheets>