

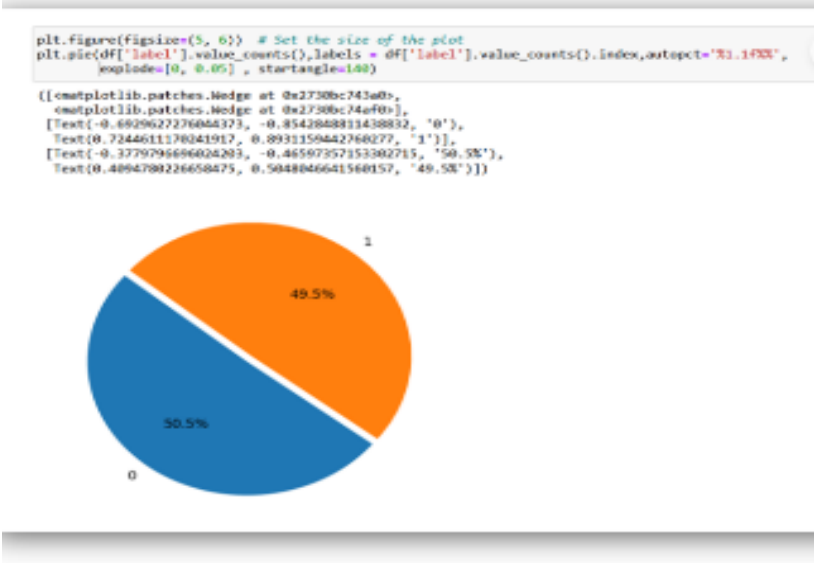
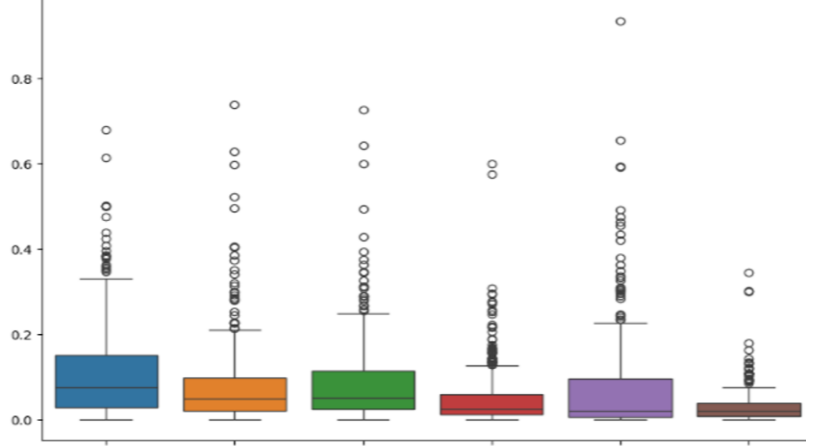
## Data Collection and Preprocessing Phase

Date	Nov 30, 2024
Team ID	739891
Project Title	Unlocking the Minds: Analyzing Mental Health with NLP
Maximum Marks	6 Marks

## Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description
Data Overview	<div> <div>textlabel</div> <div> <div>0dear american teens question dutch person hear...</div> <div>0</div> </div> <div> <div>1nothing look forward lifei dont many reasons k...</div> <div>1</div> </div> <div> <div>2music recommendations im looking expand playli...</div> <div>0</div> </div> <div> <div>3im done trying feel betterthe reason im still ...</div> <div>1</div> </div> <div> <div>4worried year old girl subject domestic physic...</div> <div>1</div> </div> <div> <div>...</div> <div>...</div> </div> <div> <div>27972posting everyday people stop caring religion ...</div> <div>0</div> </div> <div> <div>27973okay definetly need hear guys opinion ive pret...</div> <div>0</div> </div> <div> <div>27974cant get dog think ill kill myselfthe last thi...</div> <div>1</div> </div> <div> <div>27975whats point princess bridei really think like ...</div> <div>1</div> </div> <div> <div>27976got nudes person might might know snapchat do ...</div> <div>0</div> </div> <div>27977 rows × 2 columns</div> </div>

Univariate Analysis	<pre>plt.figure(figsize=(5, 6)) # Set the size of the plot plt.pie(df['label'].value_counts(), labels = df['label'].value_counts().index, autopct='%1.1f%%',         explode=[0, 0.05], startangle=100)  ([matplotlib.patches.Wedge at 0x2730bc743a08],  matplotlib.patches.Wedge at 0x2730bc74af08],  [Text(-0.6629627276044373, -0.8542848811438832, '0'),   Text(0.7244611178041917, 0.8931159442760277, '1')],  [Text(-0.3779796666024293, -0.46597357153302715, '50.5%'),   Text(0.4094780226658475, 0.5048046641560157, '49.5%')])</pre> 
Bivariate Analysis	<p>To find the relation between two features we use bivariate analysis. Here we don't have 2 numerical columns in our dataset.</p>
Multivariate Analysis	<p>In simple words, multivariate analysis is to find the relation between multiple features. Here we don't have multiple features in our dataset.</p>
Outliers and Anomalies	
Data Preprocessing Code Screenshots	

Loading Data	<pre>i4]: df=pd.read_csv("mental_health.csv") df</pre> <table><tr><th></th><th>text</th><th>label</th></tr><tr><td>0</td><td>dear american teens question dutch person hear...</td><td>0</td></tr><tr><td>1</td><td>nothing look forward lifei dont many reasons k...</td><td>1</td></tr><tr><td>2</td><td>music recommendations im looking expand playli...</td><td>0</td></tr><tr><td>3</td><td>im done trying feel betterthe reason im still ...</td><td>1</td></tr><tr><td>4</td><td>worried year old girl subject domestic physic...</td><td>1</td></tr><tr><td>...</td><td>...</td><td>...</td></tr><tr><td>27972</td><td>posting everyday people stop caring religion ...</td><td>0</td></tr><tr><td>27973</td><td>okay definetly need hear guys opinion ive pret...</td><td>0</td></tr><tr><td>27974</td><td>cant get dog think ill kill myselfthe last thi...</td><td>1</td></tr><tr><td>27975</td><td>whats point princess bridei really think like ...</td><td>1</td></tr><tr><td>27976</td><td>got nudes person might might know snapchat do ...</td><td>0</td></tr></table> <p>27977 rows × 2 columns</p>		text	label	0	dear american teens question dutch person hear...	0	1	nothing look forward lifei dont many reasons k...	1	2	music recommendations im looking expand playli...	0	3	im done trying feel betterthe reason im still ...	1	4	worried year old girl subject domestic physic...	1	...	...	...	27972	posting everyday people stop caring religion ...	0	27973	okay definetly need hear guys opinion ive pret...	0	27974	cant get dog think ill kill myselfthe last thi...	1	27975	whats point princess bridei really think like ...	1	27976	got nudes person might might know snapchat do ...	0
	text	label																																			
0	dear american teens question dutch person hear...	0																																			
1	nothing look forward lifei dont many reasons k...	1																																			
2	music recommendations im looking expand playli...	0																																			
3	im done trying feel betterthe reason im still ...	1																																			
4	worried year old girl subject domestic physic...	1																																			
...	...	...																																			
27972	posting everyday people stop caring religion ...	0																																			
27973	okay definetly need hear guys opinion ive pret...	0																																			
27974	cant get dog think ill kill myselfthe last thi...	1																																			
27975	whats point princess bridei really think like ...	1																																			
27976	got nudes person might might know snapchat do ...	0																																			
Handling Missing Data	<pre>i1]: df.describe(include='all')</pre> <table><tr><th></th><th>text</th><th>label</th></tr><tr><td>count</td><td>27972</td><td>27972.000000</td></tr><tr><td>unique</td><td>27972</td><td>NaN</td></tr><tr><td>top</td><td>dear american teens question dutch person hear...</td><td>NaN</td></tr><tr><td>freq</td><td>1</td><td>NaN</td></tr><tr><td>mean</td><td>NaN</td><td>0.494709</td></tr><tr><td>std</td><td>NaN</td><td>0.499981</td></tr><tr><td>min</td><td>NaN</td><td>0.000000</td></tr><tr><td>25%</td><td>NaN</td><td>0.000000</td></tr><tr><td>50%</td><td>NaN</td><td>0.000000</td></tr><tr><td>75%</td><td>NaN</td><td>1.000000</td></tr><tr><td>max</td><td>NaN</td><td>1.000000</td></tr></table>		text	label	count	27972	27972.000000	unique	27972	NaN	top	dear american teens question dutch person hear...	NaN	freq	1	NaN	mean	NaN	0.494709	std	NaN	0.499981	min	NaN	0.000000	25%	NaN	0.000000	50%	NaN	0.000000	75%	NaN	1.000000	max	NaN	1.000000
	text	label																																			
count	27972	27972.000000																																			
unique	27972	NaN																																			
top	dear american teens question dutch person hear...	NaN																																			
freq	1	NaN																																			
mean	NaN	0.494709																																			
std	NaN	0.499981																																			
min	NaN	0.000000																																			
25%	NaN	0.000000																																			
50%	NaN	0.000000																																			
75%	NaN	1.000000																																			
max	NaN	1.000000																																			
Data Transformation	<pre>from sklearn.feature_extraction.text import TfidfVectorizer tf = TfidfVectorizer() data_vec = tf.fit_transform(df['preprocessed_text']) print(data_vec)</pre> <table><tr><td>(0, 35255)</td><td>0.21805636791047633</td></tr><tr><td>(0, 48216)</td><td>0.11455966748601445</td></tr><tr><td>(0, 33204)</td><td>0.07123588556785876</td></tr><tr><td>(0, 23766)</td><td>0.32202234022472354</td></tr><tr><td>(0, 57707)</td><td>0.36934655986537956</td></tr><tr><td>(0, 53107)</td><td>0.33362148721649937</td></tr><tr><td>(0, 32593)</td><td>0.38735609571127644</td></tr></table>	(0, 35255)	0.21805636791047633	(0, 48216)	0.11455966748601445	(0, 33204)	0.07123588556785876	(0, 23766)	0.32202234022472354	(0, 57707)	0.36934655986537956	(0, 53107)	0.33362148721649937	(0, 32593)	0.38735609571127644																						
(0, 35255)	0.21805636791047633																																				
(0, 48216)	0.11455966748601445																																				
(0, 33204)	0.07123588556785876																																				
(0, 23766)	0.32202234022472354																																				
(0, 57707)	0.36934655986537956																																				
(0, 53107)	0.33362148721649937																																				
(0, 32593)	0.38735609571127644																																				
Feature Engineering	Attached the codes in final submission																																				
Save Processed Data	-																																				