

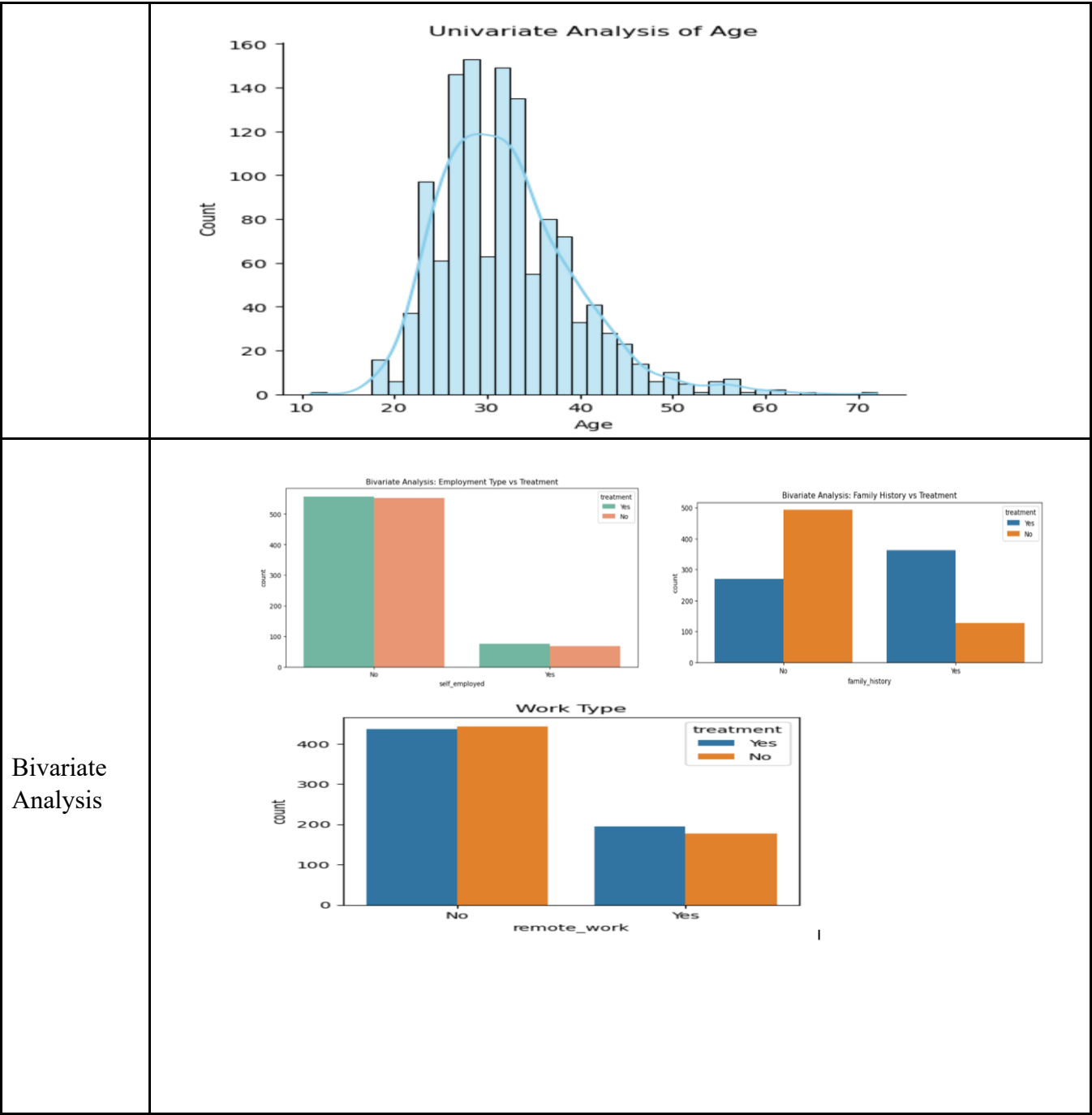
Data Collection and Preprocessing Phase

Date	19 May 2025
Team ID	SWTID1750233055
Project Title	Mental Health Prediction
Maximum Marks	6 Marks

Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description																																																																																																																																				
Data Overview	<u>Dimension:</u> 1259 rows × 27 columns																																																																																																																																				
	<u>Descriptive statistics:</u>																																																																																																																																				
	<table><tr><th></th><th>Age</th><th>Gender</th><th>self_employed</th><th>family_history</th><th>treatment</th><th>work_interfere</th><th>no_employees</th><th>remote_work</th><th>tech_company</th><th>benefits</th></tr><tr><td>count</td><td>1252.000000</td><td>1252</td><td>1252</td><td>1252</td><td>1252</td><td>1252</td><td>1252</td><td>1252</td><td>1252</td><td>1252</td></tr><tr><td>unique</td><td>NaN</td><td>38</td><td>2</td><td>2</td><td>2</td><td>4</td><td>6</td><td>2</td><td>2</td><td>3</td></tr><tr><td>top</td><td>NaN</td><td>Male</td><td>No</td><td>No</td><td>Yes</td><td>Sometimes</td><td>6-25</td><td>No</td><td>Yes</td><td>Yes</td></tr><tr><td>freq</td><td>NaN</td><td>820</td><td>1109</td><td>763</td><td>632</td><td>726</td><td>289</td><td>880</td><td>1026</td><td>473</td></tr><tr><td>mean</td><td>32.059904</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>std</td><td>7.309669</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>min</td><td>11.000000</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>25%</td><td>27.000000</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>50%</td><td>31.000000</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>75%</td><td>36.000000</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>max</td><td>72.000000</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr></table>		Age	Gender	self_employed	family_history	treatment	work_interfere	no_employees	remote_work	tech_company	benefits	count	1252.000000	1252	1252	1252	1252	1252	1252	1252	1252	1252	unique	NaN	38	2	2	2	4	6	2	2	3	top	NaN	Male	No	No	Yes	Sometimes	6-25	No	Yes	Yes	freq	NaN	820	1109	763	632	726	289	880	1026	473	mean	32.059904	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	std	7.309669	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	min	11.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	25%	27.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	50%	31.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	75%	36.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	max	72.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
		Age	Gender	self_employed	family_history	treatment	work_interfere	no_employees	remote_work	tech_company	benefits																																																																																																																										
	count	1252.000000	1252	1252	1252	1252	1252	1252	1252	1252	1252																																																																																																																										
	unique	NaN	38	2	2	2	4	6	2	2	3																																																																																																																										
	top	NaN	Male	No	No	Yes	Sometimes	6-25	No	Yes	Yes																																																																																																																										
	freq	NaN	820	1109	763	632	726	289	880	1026	473																																																																																																																										
	mean	32.059904	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																										
	std	7.309669	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																										
	min	11.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																										
	25%	27.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																										
	50%	31.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																										
75%	36.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																											
max	72.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																											
Univariate Analysis																																																																																																																																					



Multivariate Analysis	-
-----------------------	---

Outliers and Anomalies	-
------------------------	---

Data Preprocessing Code Screenshots

Loading Data

#1.Loading Data Set

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sb
data = pd.read_csv(r"C:\Users\jaind\OneDrive\Desktop\ML_Project\data\survey.csv")
|
data
```

	Timestamp	Age	Gender	Country	state	self_employed	family_history	treatment	work_interfere	no_employees	...	leave	mental_health_consequence
0	2014-08-27 11:29:31	37	Female	United States	IL	NaN	No	Yes	Often	6-25	...	Somewhat easy	No
1	2014-08-27 11:29:37	44	M	United States	IN	NaN	No	No	Rarely	More than 1000	...	Don't know	Maybe
2	2014-08-27 11:29:44	32	Male	Canada	NaN	NaN	No	No	Rarely	6-25	...	Somewhat difficult	No
3	2014-08-27 11:29:46	31	Male	United Kingdom	NaN	NaN	Yes	Yes	Often	26-100	...	Somewhat difficult	Yes
4	2014-08-27 11:30:22	31	Male	United States	TX	NaN	No	No	Never	100-500	...	Don't know	No

Handling Missing Data	<pre>data.isnull().sum()</pre> <p><i>#we Have 4 Columns Which Have Null Values</i> <i>#Change The Value not available to anything Efficient To Reduce overfitting And Increase Accuracy</i></p> <pre>data['self_employed'].value_counts()</pre> <pre>self_employed No 1095 Yes 146 Name: count, dtype: int64</pre> <p><i>#As majority self employed Come with No Replace Null values To No</i></p> <pre>data['self_employed'] = data['self_employed'].fillna('no')</pre>
Data Transformation	<pre># Convert 'Age' to numeric and remove rows with invalid age data['Age'] = pd.to_numeric(data['Age'], errors='coerce') data.dropna(subset=['Age'], inplace=True)</pre> <p><i># Drop unnecessary columns (if they exist)</i></p> <pre>data.drop(columns=['Country', 'Timestamp', 'state', 'comments', 'leave'], errors='ignore', inplace=True)</pre> <p><i># Split features and Label</i></p> <pre>x = data.drop('treatment', axis=1) y = data['treatment']</pre> <p><i># Encode Label (target variable)</i></p> <pre>le = LabelEncoder() y = le.fit_transform(y)</pre> <p><i># Identify numerical and categorical columns</i></p> <pre>num_cols = ['Age'] cat_cols = [col for col in x.columns if col not in num_cols]</pre> <p><i># Define column transformer</i></p>
Feature Engineering	Attached the codes in final submission.
Save Processed Data	-