

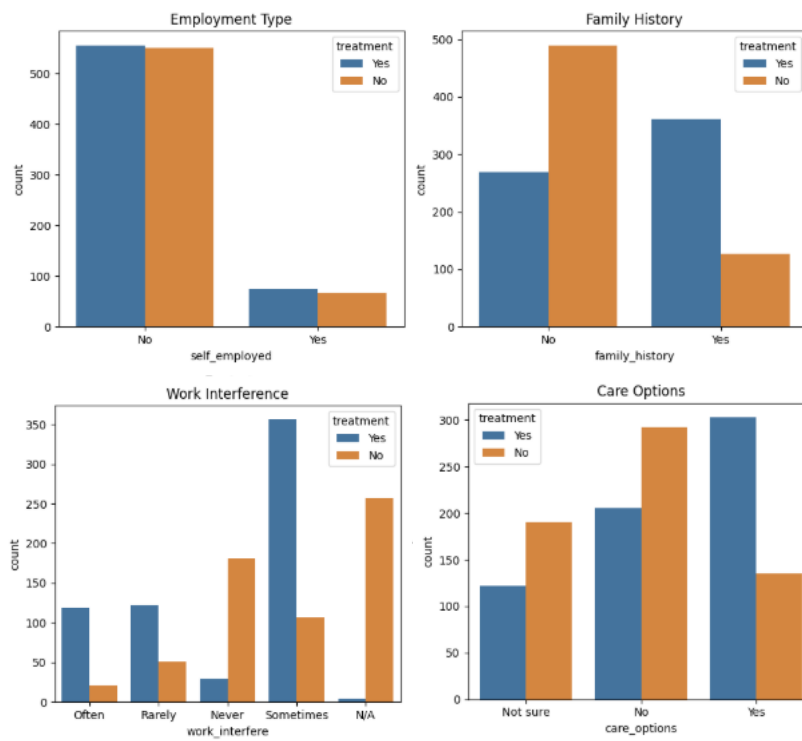
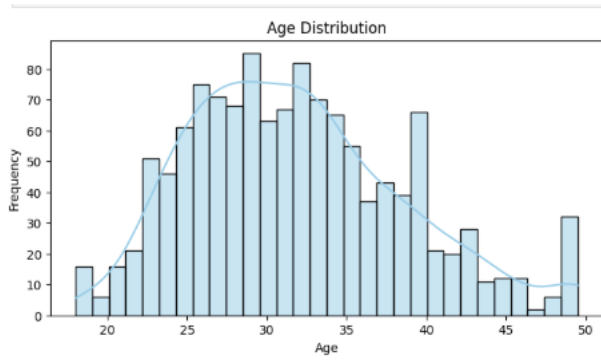
## Data Collection and Pre-processing Phase

Date	15 March 2024
Team ID	SWTID1749622322
Project Title	Mental Health Prediction
Maximum Marks	6 Marks

### Data Exploration and Preprocessing Report

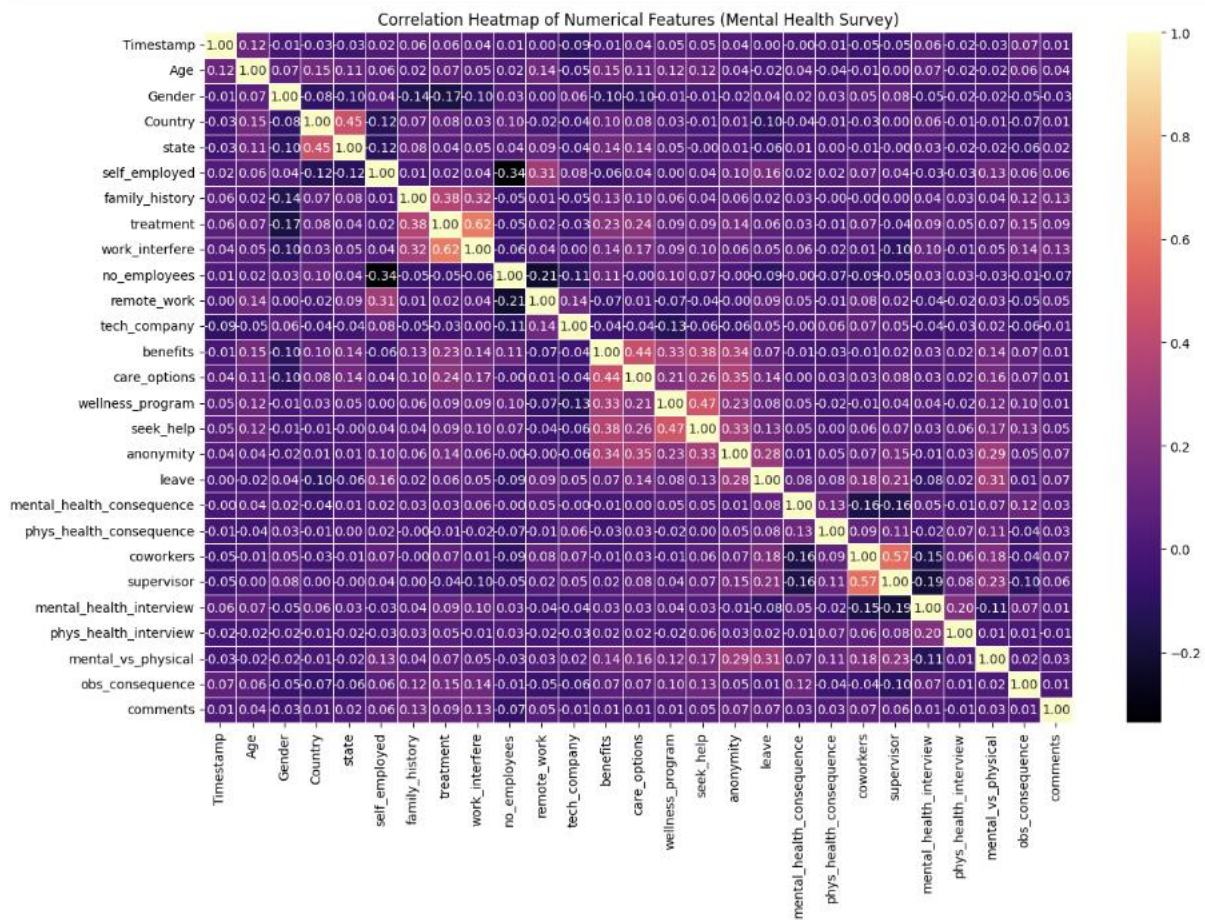
Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong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

Section	Description																																																																																																																																																									
Data Overview	<table><thead><tr><th></th><th>Timestamp</th><th>Age</th><th>Gender</th><th>Country</th><th>state</th><th>self_employed</th><th>family_history</th><th>treatment</th><th>work_interfere</th><th>no_employees</th></tr></thead><tbody><tr><td>count</td><td>1247</td><td>1247.000000</td><td>1247</td><td>1247</td><td>735</td><td>1247</td><td>1247</td><td>1247</td><td>1247</td><td>1247</td></tr><tr><td>unique</td><td>1235</td><td>NaN</td><td>5</td><td>46</td><td>45</td><td>2</td><td>2</td><td>2</td><td>5</td><td>5</td></tr><tr><td>top</td><td>2014-08-27 12:43:28</td><td>NaN</td><td>Male</td><td>United States</td><td>CA</td><td>No</td><td>No</td><td>Yes</td><td>Sometimes</td><td>6-7</td></tr><tr><td>freq</td><td>2</td><td>NaN</td><td>979</td><td>743</td><td>137</td><td>1107</td><td>759</td><td>630</td><td>463</td><td>28</td></tr><tr><td>mean</td><td>NaN</td><td>31.873296</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>NaN</td><td>6.756424</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>NaN</td><td>18.000000</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></tr><tr><td>25%</td><td>NaN</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></tr><tr><td>50%</td><td>NaN</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></tr><tr><td>75%</td><td>NaN</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></tr><tr><td>max</td><td>NaN</td><td>49.500000</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 colspan="11">11 rows x 27 columns</td></tr></tbody></table>												Timestamp	Age	Gender	Country	state	self_employed	family_history	treatment	work_interfere	no_employees	count	1247	1247.000000	1247	1247	735	1247	1247	1247	1247	1247	unique	1235	NaN	5	46	45	2	2	2	5	5	top	2014-08-27 12:43:28	NaN	Male	United States	CA	No	No	Yes	Sometimes	6-7	freq	2	NaN	979	743	137	1107	759	630	463	28	mean	NaN	31.873296	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	std	NaN	6.756424	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	min	NaN	18.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	25%	NaN	27.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	50%	NaN	31.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	75%	NaN	36.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	max	NaN	49.500000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	11 rows x 27 columns										
		Timestamp	Age	Gender	Country	state	self_employed	family_history	treatment	work_interfere	no_employees																																																																																																																																															
	count	1247	1247.000000	1247	1247	735	1247	1247	1247	1247	1247																																																																																																																																															
	unique	1235	NaN	5	46	45	2	2	2	5	5																																																																																																																																															
	top	2014-08-27 12:43:28	NaN	Male	United States	CA	No	No	Yes	Sometimes	6-7																																																																																																																																															
	freq	2	NaN	979	743	137	1107	759	630	463	28																																																																																																																																															
	mean	NaN	31.873296	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																															
	std	NaN	6.756424	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																															
	min	NaN	18.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																															
	25%	NaN	27.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																															
	50%	NaN	31.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																															
	75%	NaN	36.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																															
max	NaN	49.500000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN																																																																																																																																																
11 rows x 27 columns																																																																																																																																																										
Univariate Analysis																																																																																																																																																										



Bivariate Analysis

## Multivariate Analysis



## Outliers and Anomalies

## Data Pre-processing Code Screenshots

## Loading Data

```
data = pd.read_csv(r'D:\Anshveer - Git - Mental Health\Mental_Health_FINAL\data\survey.csv')
```

```
data.head()
```

	Timestamp	Age	Gender	Country	state	self_employed	family_history	treatment	work_interfere	no_employees	...	lr
0	2014-08-27 11:29:31	37	Female	United States	IL	NaN	No	Yes	Often	6-25	...	Some
1	2014-08-27 11:29:37	44	M	United States	IN	NaN	No	No	Rarely	More than 1000	...	C
2	2014-08-27 11:29:44	32	Male	Canada	NaN	NaN	No	No	Rarely	6-25	...	Some
3	2014-08-27 11:29:46	31	Male	United Kingdom	NaN	NaN	Yes	Yes	Often	26-100	...	Some
4	2014-08-27 11:30:22	31	Male	United States	TX	NaN	No	No	Never	100-500	...	C

5 rows x 27 columns

## Handling Missing Data

```
data['self_employed'].value_counts()
```

```
self_employed
No      1095
Yes      146
Name: count, dtype: int64
```

```
data['self_employed'].fillna('No',inplace=True)
```

```
data['work_interfere'].value_counts()
```

```
work_interfere
Sometimes  465
Never      213
Rarely     173
Often      144
Name: count, dtype: int64
```

```
data['work_interfere'].fillna('N/A',inplace=True)
```

```
data.drop(data[(data['Age']>60) | (data['Age']<18)].index,inplace=True)
```

## Data Transformation

```
data['Gender'].replace(['Male ', 'male', 'M', 'm', 'Male', 'Cis Male',
                        'Man', 'cis male', 'Mail', 'Male-ish', 'Male (CIS)',
                        'Cis Man', 'msle', 'Malr', 'Mal', 'maile', 'Make'],
                        'Male', inplace=True)

data['Gender'].replace(['Female ', 'female', 'F', 'f', 'Woman', 'Female',
                        'femal', 'Cis Female', 'cis-female/femme', 'Femake', 'Female (cis)',
                        'woman'],
                        'Female', inplace=True)

data["Gender"].replace(['Female (trans)', 'queer/she/they', 'non-binary',
                        'fluid', 'queer', 'Androgyne', 'Trans-female', 'male leaning androgynous',
                        'Agender', 'A little about you', 'Nah', 'All',
                        'Ostensibly male, unsure what that really means',
                        'Genderqueer', 'Enby', 'p', 'Neuter', 'something kinda male?',
                        'Guy (-ish) ^_^', 'Trans woman', ],
                        'Non-Binary', inplace=True)
```

Feature Engineering	Attached the codes in final submission.
Save Processed Data	-