

PUBLIC HEALTH AWARENESS CAMPAIGN

– Phase3

TEAM MEMBERS:

KEERTHI VARSHINI S	- 2021115053
KONETI SASANK	- 2021115054
KOWSHIK T	- 2021115055
LEKHA S	- 2021115056
SANTHOSH	-2021115311

Dataset Preprocessing and Cleaning :

In this phase, the primary task to complete is to analyze the dataset and make sure the dataset is clean. In the cleaning process, we may have to remove null values and make sure that the dataset finally contains all non-null values

▼ Import necessary libraries

```
#imports necessary libraries to do basic things on the dataset
import pandas as pd
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt

print('Successfully imported')

Successfully imported
```

▼ Read Dataset

```
#Reading data
data = pd.read_csv('survey.csv')
data.head()
```

	Timestamp	Age	Gender	Country	state	self_employed	family_history	treatment
0	2014-08-27 11:29:31	37	Female	United States	IL	NaN	No	Yes
1	2014-08-27 11:29:37	44	M	United States	IN	NaN	No	No
2	2014-08-27 11:29:44	32	Male	Canada	NaN	NaN	No	No
3	2014-08-27 11:29:46	31	Male	United Kingdom	NaN	NaN	Yes	Yes
4	2014-08-27 11:30:22	31	Male	United States	TX	NaN	No	No

5 rows × 27 columns

▼ Preprocessing and Cleaning dataset

```
#Check the dataset for missing data
if data.isnull().sum().sum() == 0 :
    print ('There is no missing data in our dataset')
else:
    print('There is {} missing data in our dataset '.format(data.isnull().sum().sum()))
```

There is 1892 missing data in our dataset

```
#Check our missing data from which columns and how many unique features they have.
frame = pd.concat([data.isnull().sum(), data.nunique(), data.dtypes], axis = 1, sort= False)
frame
```

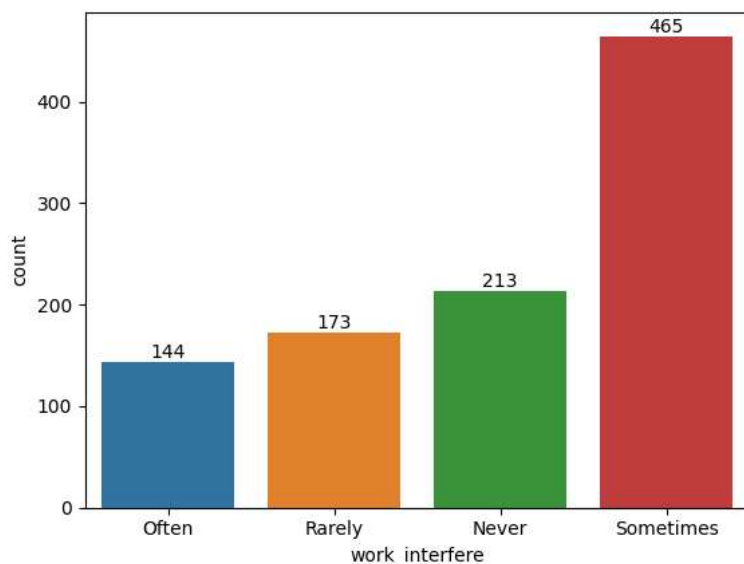
	0	1	2
Timestamp	0	1246	object
Age	0	53	int64
Gender	0	49	object
Country	0	48	object
state	515	45	object
self_employed	18	2	object
family_history	0	2	object
treatment	0	2	object
work_interfere	264	4	object
no_employees	0	6	object
remote_work	0	2	object
tech_company	0	2	object
benefits	0	3	object
care_options	0	3	object
wellness_program	0	3	object
seek_help	0	3	object
anonymity	0	3	object

```
#Look at what is in the 'Work_interfere' column to choose a suitable method to fill nan values.
data['work_interfere'].unique()
```

```
array(['Often', 'Rarely', 'Never', 'Sometimes', nan], dtype=object)

coworkers      0      3  object
```

```
#Plot **work_interfere**
ax = sns.countplot(data = data , x = 'work_interfere');
#Add the value of each parameter on the Plot
ax.bar_label(ax.containers[0]);
```



```
from sklearn.impute import SimpleImputer
import numpy as np
columns_to_drop = ['state', 'comments', 'Timestamp']
for column in columns_to_drop:
    if column in data.columns:
        data = data.drop(columns=[column])

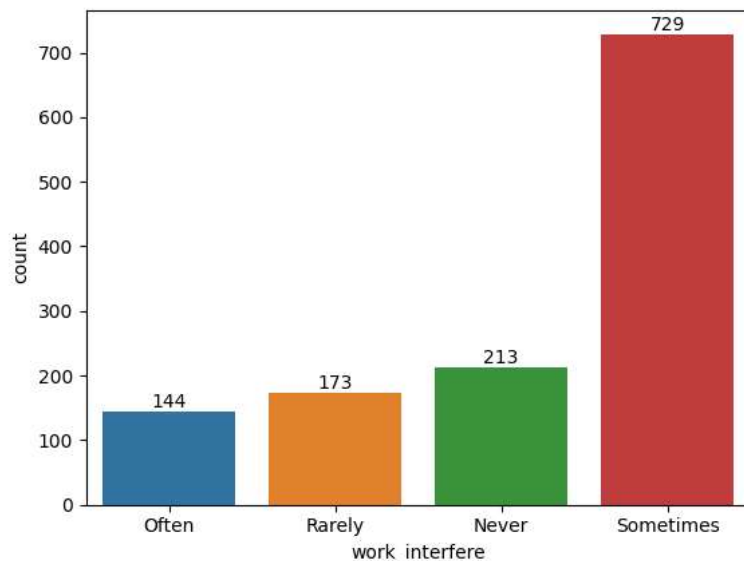
# Fill in missing values in work_interfere column
data['work_interfere'] = np.ravel(SimpleImputer(strategy = 'most_frequent').fit_transform(data['work_interfere'].values.reshape(-1,1)))
data['self_employed'] = np.ravel(SimpleImputer(strategy = 'most_frequent').fit_transform(data['self_employed'].values.reshape(-1,1)))

data.head()
```

	Age	Gender	Country	self_employed	family_history	treatment	work_interfere	no_employees	remote_work	tech_company	...	anor
0	37	Female	United States	No	No	Yes	Often	6-25	No	Yes	...	
1	44	M	United States	No	No	No	Rarely	More than 1000	No	No	...	Don
2	32	Male	Canada	No	No	No	Rarely	6-25	No	Yes	...	Don
3	31	Male	United Kingdom	No	Yes	Yes	Often	26-100	No	Yes	...	
4	31	Male	United States	No	No	No	Never	100-500	Yes	Yes	...	Don

5 rows × 24 columns

```
ax = sns.countplot(data=data, x='work_interfere');
ax.bar_label(ax.containers[0]);
```



```
#Check unique data in gender columns
print(data['Gender'].unique())
print('')
print('-'*75)
print('')
#Check number of unique data too.
print('number of unique Gender in our dataset is :', data['Gender'].nunique())
```

```
['Female' 'M' 'Male' 'male' 'female' 'm' 'Male-ish' 'maile' 'Trans-female'
'Cis Female' 'F' 'something kinda male?' 'Cis Male' 'Woman' 'f' 'Mal'
'Male (CIS)' 'queer/she/they' 'non-binary' 'Femake' 'woman' 'Make' 'Nah'
'All' 'Enby' 'fluid' 'Genderqueer' 'Female ' 'Androgyne' 'Agender'
'cis-female/femme' 'Guy (-ish) ^_^' 'male leaning androgynous' 'Male '
'Man' 'Trans woman' 'msle' 'Neuter' 'Female (trans)' 'queer'
'Female (cis)' 'Mail' 'cis male' 'A little about you' 'Malr' 'p' 'femail'
'Cis Man' 'ostensibly male, unsure what that really means']
```

number of unique Gender in our dataset is : 49

```
#Gender data contains dictation problems, nonsense answers, and too unique Genders.
#_So Let's clean it and organize it into Male, Female, and other categories
```

```
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',
'femail', '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',], 'Other', inplace = True)

print(data['Gender'].unique())

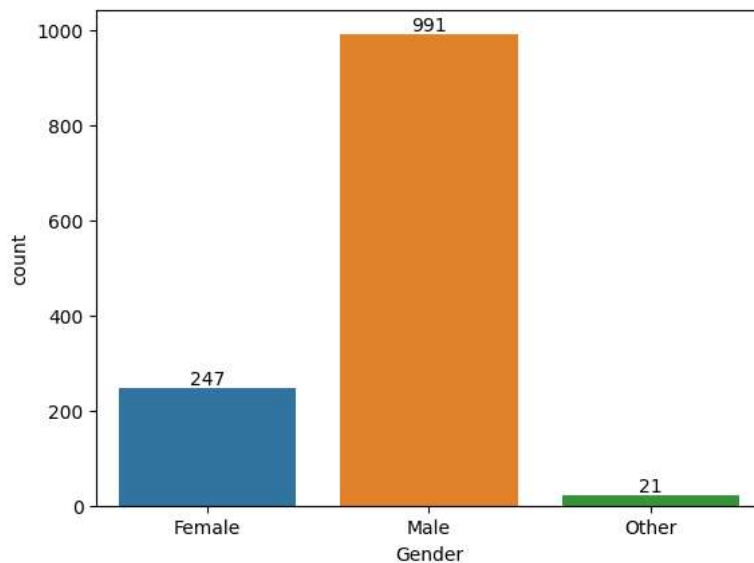
['Female' 'Male' 'Other']

```

```

#Plot Genders column after cleaning and new categorizing
ax = sns.countplot(data=data, x='Gender');
ax.bar_label(ax.containers[0]);

```



▼ CHECKING AFTER CLEANING THE DATASET

```

#Our data is clean now ? let's see.
if data.isnull().sum().sum() == 0:
    print('There is no missing data')
else:
    print('There is {} missing data'.format(data.isnull().sum().sum()))

```

There is no missing data

```

#Let's check duplicated data.
if data.duplicated().sum() == 0:
    print('There is no duplicated data:')
else:
    print('There is {} duplicated data:'.format(data.duplicated().sum()))
    #If there is duplicated data drop it.
    data.drop_duplicates(inplace=True)

print('-'*50)
print(data.duplicated().sum())

```

There is no duplicated data:

```

-----
0

```

```

#Look unique data in Age column
data['Age'].unique()

```

```

array([
    37,    44,    32,    31,    33,
    35,    39,    42,    23,    29,
    36,    27,    46,    41,    34,
    30,    40,    38,    50,    24,
    18,    28,    26,    22,    19,
    25,    45,    21,   -29,    43,
    56,    60,    54,   329,    55,
    9999999999,    48,    20,    57,    58,
    47,    62,    51,    65,    49,
   -1726,     5,    53,    61,     8,
    11,    -1,    72])

```

```

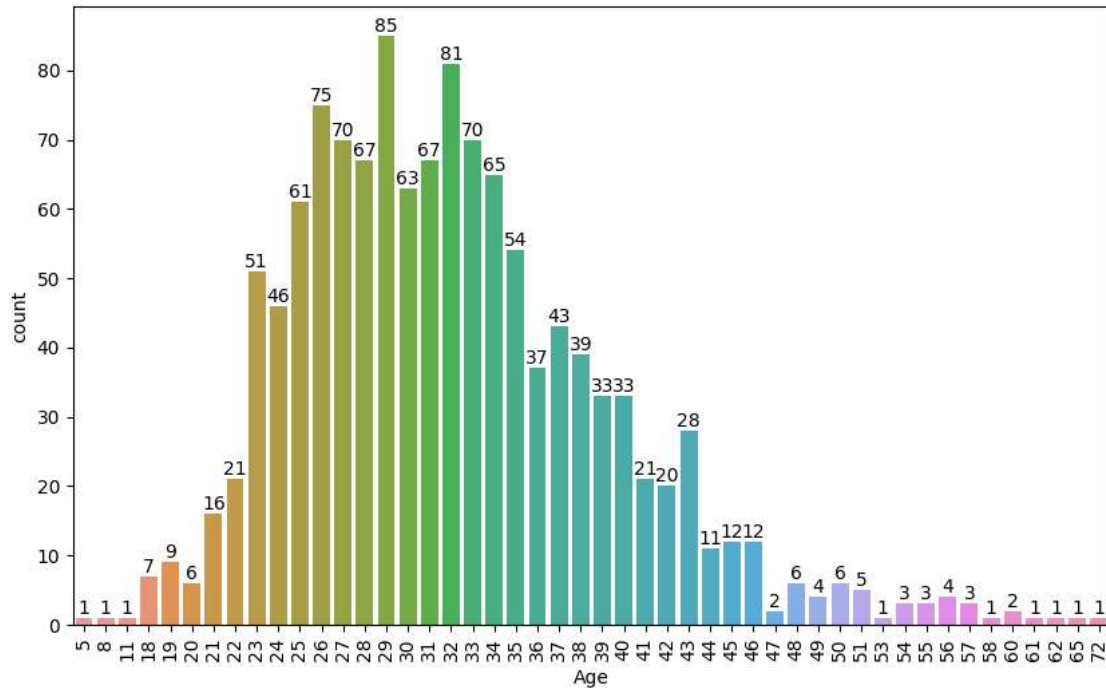
#We had a lot of nonsense answers in the Age column too
#This filtering will drop entries exceeding 100 years and those indicating negative values.
data.drop(data[data['Age']<0].index, inplace = True)
data.drop(data[data['Age']>99].index, inplace = True)

```

```
print(data['Age'].unique())
```

```
[37 44 32 31 33 35 39 42 23 29 36 27 46 41 34 30 40 38 50 24 18 28 26 22
 19 25 45 21 43 56 60 54 55 48 20 57 58 47 62 51 65 49  5 53 61  8 11 72]
```

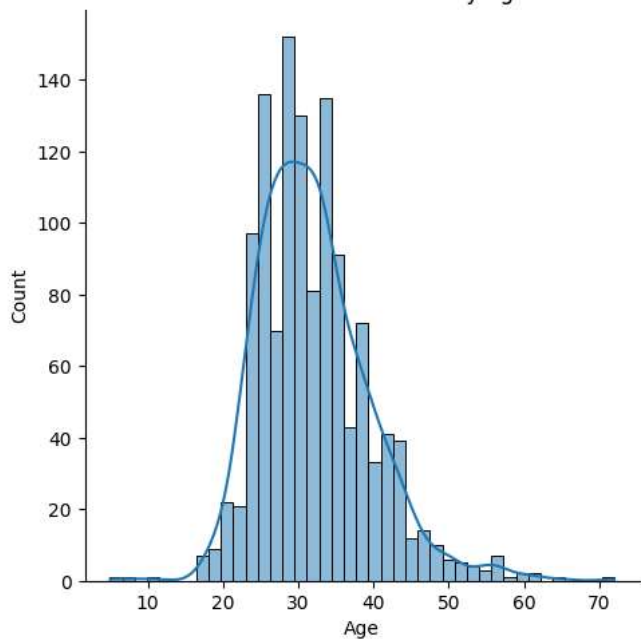
```
#Let's see the Age distribution in this dataset.
plt.figure(figsize = (10,6))
age_range_plot = sns.countplot(data = data, x = 'Age');
age_range_plot.bar_label(age_range_plot.containers[0]);
plt.xticks(rotation=90);
```



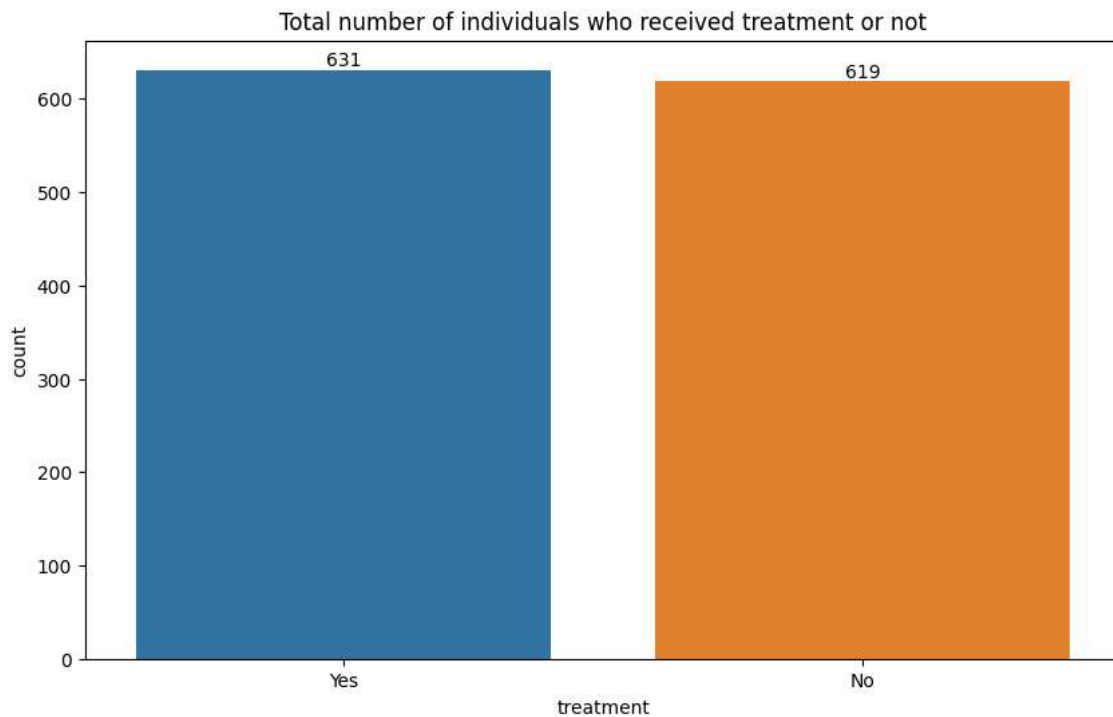
```
#In this plot moreover on Age distribution we can see treatment distribution by age
plt.figure(figsize=(10, 6));
sns.displot(data['Age'], kde = 'treatment');
plt.title('Distribution treatment by age');
```

<Figure size 1000x600 with 0 Axes>

Distribution treatment by age



```
#In this plot We can see Total number of individuals who received treatment or not.
plt.figure(figsize = (10,6));
treat = sns.countplot(data = data, x = 'treatment');
treat.bar_label(treat.containers[0]);
plt.title('Total number of individuals who received treatment or not');
```



```
#Check Dtypes
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1250 entries, 0 to 1258
Data columns (total 24 columns):
#   Column                Non-Null Count  Dtype
---  -
0    Age                   1250 non-null  int64
1    Gender                1250 non-null  object
2    Country               1250 non-null  object
3    self_employed         1250 non-null  object
4    family_history        1250 non-null  object
5    treatment             1250 non-null  object
6    work_interfere        1250 non-null  object
7    no_employees          1250 non-null  object
8    remote_work           1250 non-null  object
9    tech_company          1250 non-null  object
10   benefits              1250 non-null  object
11   care_options          1250 non-null  object
12   wellness_program      1250 non-null  object
13   seek_help             1250 non-null  object
14   anonymity             1250 non-null  object
15   leave                 1250 non-null  object
16   mental_health_consequence 1250 non-null  object
17   phys_health_consequence 1250 non-null  object
18   coworkers             1250 non-null  object
19   supervisor            1250 non-null  object
20   mental_health_interview 1250 non-null  object
21   phys_health_interview  1250 non-null  object
22   mental_vs_physical     1250 non-null  object
23   obs_consequence       1250 non-null  object
dtypes: int64(1), object(23)
memory usage: 244.1+ KB
```

```
#Use LabelEncoder to change the Dtypes to 'int'
from sklearn.preprocessing import LabelEncoder
```

```
le = LabelEncoder()
#Make the dataset include all the columns we need to change their dtypes
columns_to_encode = ['Gender', 'Country', 'self_employed', 'family_history', 'treatment', 'work_interfere', 'no_employees',
                    'remote_work', 'tech_company', 'benefits', 'care_options', 'wellness_program',
                    'seek_help', 'anonymity', 'leave', 'mental_health_consequence', 'phys_health_consequence',
                    'coworkers', 'supervisor', 'mental_health_interview', 'phys_health_interview',
                    'mental_vs_physical', 'obs_consequence']
#Write a Loop for fitting LabelEncoder on columns_to_encode
for columns in columns_to_encode:
```

```
data[columns] = le.fit_transform(data[columns])

data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1250 entries, 0 to 1258
Data columns (total 24 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Age                                   1250 non-null   int64
1   Gender                               1250 non-null   int64
2   Country                             1250 non-null   int64
3   self_employed                       1250 non-null   int64
4   family_history                      1250 non-null   int64
5   treatment                           1250 non-null   int64
6   work_interfere                      1250 non-null   int64
7   no_employees                       1250 non-null   int64
8   remote_work                         1250 non-null   int64
9   tech_company                       1250 non-null   int64
10  benefits                            1250 non-null   int64
11  care_options                       1250 non-null   int64
12  wellness_program                   1250 non-null   int64
13  seek_help                         1250 non-null   int64
14  anonymity                          1250 non-null   int64
15  leave                             1250 non-null   int64
16  mental_health_consequence          1250 non-null   int64
17  phys_health_consequence             1250 non-null   int64
18  coworkers                          1250 non-null   int64
19  supervisor                         1250 non-null   int64
20  mental_health_interview            1250 non-null   int64
21  phys_health_interview              1250 non-null   int64
22  mental_vs_physical                 1250 non-null   int64
23  obs_consequence                    1250 non-null   int64
dtypes: int64(24)
memory usage: 244.1 KB
```