

Mental health prediction EDA(Exploratory Data Analysis) and ML Models

1. Import Libraries

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy import stats
6 from scipy.stats import randint
7 from sklearn.model_selection import train_test_split
8 from sklearn import preprocessing
9 from sklearn.datasets import make_classification
10 from sklearn.preprocessing import binarize, LabelEncoder, MinMaxScaler
```

2. Data Preprocessing

```
1 #Print the dataframe
2 #Dataset link : "https://www.kaggle.com/datasets/ron2112/mental-health-data"
3 url = r'/content/drive/MyDrive/ML Innovative/Mental Health Data.csv'
4 data=pd.read_csv(url)
5 data.head(10)
```

Saving...



| | Are you self-employed? | How many employees does your company or organization have? | Is your employer primarily a tech company/organization? | Is your primary role within your company related to tech/IT? | Does your employer provide mental health benefits as part of healthcare coverage? | Do you know the options for mental health care available under your employer-provided coverage? | |
|---|------------------------|--|---|--|---|---|--|
| 0 | 0 | 1 to 5 | 1.0 | NaN | Yes | Yes | |
| 1 | 0 | 1 to 5 | 1.0 | NaN | No | No | |
| 2 | 0 | 1 to 5 | 1.0 | NaN | Yes | Yes | |
| 3 | 0 | 1 to 5 | 1.0 | NaN | No | No | |
| 4 | 0 | 1 to 5 | 0.0 | 1.0 | I don't know | No | |
| 5 | 0 | 1 to 5 | 1.0 | NaN | No | Yes | |

Not eligible

Saving...

×

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1433 entries, 0 to 1432
Data columns (total 31 columns):
#   Column
---  ---
0   Are you self-employed?
1   How many employees does your company or organization have?
2   Is your employer primarily a tech company/organization?
3   Is your primary role within your company related to tech/IT?
4   Does your employer provide mental health benefits as part of healthcare coverage?
5   Do you know the options for mental health care available under your employer-provic
6   Has your employer ever formally discussed mental health (for example, as part of a
7   Does your employer offer resources to learn more about mental health concerns and c
```

```

8   If a mental health issue prompted you to request a medical leave from work, asking
9   Would you feel comfortable discussing a mental health disorder with your coworkers?
10  Do you feel that your employer takes mental health as seriously as physical health?
11  Do you know local or online resources to seek help for a mental health disorder?
12  Do you believe your productivity is ever affected by a mental health issue?
13  If yes, what percentage of your work time (time performing primary or secondary job)?
14  How willing would you be to share with friends and family that you have a mental illness?
15  Do you have a family history of mental illness?
16  Have you had a mental health disorder in the past?
17  Do you currently have a mental health disorder?
18  Have you been diagnosed with a mental health condition by a medical professional?
19  If so, what condition(s) were you diagnosed with?
20  Have you ever sought treatment for a mental health issue from a mental health professional?
21  If you have a mental health issue, do you feel that it interferes with your work?
22  If you have a mental health issue, do you feel that it interferes with your work?
23  What is your age?
24  What is your gender?
25  What country do you live in?
26  What US state or territory do you live in?
27  What country do you work in?
28  What US state or territory do you work in?
29  Which of the following best describes your work position?
30  Do you work remotely?
dtypes: float64(2), int64(3), object(26)
memory usage: 347.2+ KB

```

```

1 #Check the Shape of dataset
2 print(data.shape)

```

```
(1433, 31)
```

```

1 #Make the list of columns
2 a=list(data.columns)
3 print(a)
4 # New name of the all columns

```

Saving...

```

7   'tech_company', 'role_IT',
8   'mental_healthcare_coverage',
9   'knowledge_about_mental_healthcare_options_workplace',
10  'employer_discussed_mental_health ',
11  'employer_offer_resources_to_learn_about_mental_health',
12  'medical_leave_from_work ',
13  'comfortable_discussing_with_coworkers',
14  'employer_take_mental_health_seriously',
15  'knowledge_of_local_online_resources ',
16  'productivity_affected_by_mental_health ',
17  'percentage_work_time_affected_mental_health',
18  'openness_of_family_friends',
19  'family_history_mental_illness',
20  'mental_health_disorder_past',

```

```

21 'currently_mental_health_disorder',
22 'diagnosed_mental_health_condition',
23 'type_of_disorder',
24 'treatment_from_professional',
25 'while_effective_treatment_mental_health_issue_interferes_work',
26 'while_not_effective_treatment_interferes_work ',
27 'age',
28 'gender',
29 'country',
30 'US state',
31 'country work ',
32 'US state work',
33 'role_in_company',
34 'work_remotely','']
35 for i,j in zip(a,b):
36     data.rename(columns={i:j},inplace=True)

```

['Are you self-employed?', 'How many employees does your company or organization have?'],

```

1 # Information of dataframe after the rename
2 data.info()

```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 1433 entries, 0 to 1432
```


```
Data columns (total 31 columns):
```

| # | Column | Non-Null Count | Dtype |
|----|---|----------------|---------|
| 0 | self_employed | 1433 non-null | int64 |
| 1 | no_of_employees | 1146 non-null | object |
| 2 | tech_company | 1146 non-null | float64 |
| 3 | role_IT | 263 non-null | float64 |
| 4 | mental_healthcare_coverage | 1146 non-null | object |
| 5 | knowledge_about_mental_healthcare_options_workplace | 1013 non-null | object |
| 6 | employer_discussed_mental_health | 1146 non-null | object |
| 7 | do_you_learn_about_mental_health | 1146 non-null | object |
| 8 | h_coworkers | 1146 non-null | object |
| 9 | employer_take_mental_health_seriously | 1146 non-null | object |
| 10 | knowledge_of_local_online_resources | 287 non-null | object |
| 11 | productivity_affected_by_mental_health | 287 non-null | object |
| 12 | percentage_work_time_affected_mental_health | 204 non-null | object |
| 13 | openess_of_family_friends | 1433 non-null | object |
| 14 | family_history_mental_illness | 1433 non-null | object |
| 15 | mental_health_disorder_past | 1433 non-null | object |
| 16 | currently_mental_health_disorder | 1433 non-null | object |
| 17 | diagnosed_mental_health_condition | 1433 non-null | object |
| 18 | type_of_disorder | 711 non-null | object |
| 19 | treatment_from_professional | 1433 non-null | int64 |
| 20 | while_effective_treatment_mental_health_issue_interferes_work | 1433 non-null | object |
| 21 | while_not_effective_treatment_interferes_work | 1433 non-null | object |
| 22 | age | 1433 non-null | int64 |
| 23 | gender | 1430 non-null | object |

Saving...



```
25 country 1433 non-null obje
26 US state 840 non-null obje
27 country work 1433 non-null obje
28 US state work 851 non-null obje
29 role_in_company 1433 non-null obje
30 work_remotely 1433 non-null obje
dtypes: float64(2), int64(3), object(26)
memory usage: 347.2+ KB
```



```
1 ## Now We Find the Missing values in different Columns
2 columns=data.columns
3 pd.DataFrame({'no of missing values':data.isnull().sum()})
```

Saving...



| | |
|---|------|
| self_employed | 0 |
| no_of_employees | 287 |
| tech_company | 287 |
| role_IT | 1170 |
| mental_healthcare_coverage | 287 |
| knowledge_about_mental_healthcare_options_workplace | 420 |
| employer_discussed_mental_health | 287 |
| employer_offer_resources_to_learn_about_mental_health | 287 |
| medical_leave_from_work | 287 |
| comfortable_discussing_with_coworkers | 287 |

```
1 # Now we copy the dataset in data1
2 data1=data.copy()
3 data1
```

Saving...

×

| self_employed | no_of_employees | tech_company | role_IT | mental_healthcare_coverage |
|---------------|-----------------|--------------|---------|----------------------------|
|---------------|-----------------|--------------|---------|----------------------------|

| 0 | 0 | 1 to 5 | 1.0 | NaN | Yes |
|---|---|--------|-----|-----|-----|
|---|---|--------|-----|-----|-----|

```

1 # Now there are sum columns which has so many tuple are not have any value so it is unnece
2 remove_columns = ['role_IT',
3                   'knowledge_of_local_online_resources ',
4                   'productivity_affected_by_mental_health ',
5                   'percentage_work_time_affected_mental_health']
6 data2=data1.drop(remove_columns,axis=1)
7 data2.shape

```

```
(1433, 27)
```

▸ Cleaning Different Columns

| 1429 | 1 | NaN | NaN | NaN | NaN |
|------|---|-----|-----|-----|-----|
|------|---|-----|-----|-----|-----|

```

1 # No of employee column
2 print(data2.no_of_employees.unique())
3 data2.no_of_employees.unique()

['1 to 5' '6 to 25' '26-99' '100-500' '26-100' '500-1000' 'More than 1000'
 nan]
array(['1 to 5', '6 to 25', '26-99', '100-500', '26-100', '500-1000',
       'More than 1000', nan], dtype=object)

1 # change the value format
2 data2.no_of_employees.replace(to_replace=['1 to 5', '6 to 25', 'More than 1000', '26-99'],
3                               value=['1-5', '6-25', '>1000', '26-100'], inplace=True)
4
5 print(data2.no_of_employees.value_counts())

```

Saving...

```

>1000      256
100-500     248
6-25       176
500-1000     80
1-5         60
Name: no_of_employees, dtype: int64

```

1

```

1 # Cleaning Mental Health Care coverage column
2 data2.mental_healthcare_coverage.unique()

```

```
array(['Yes', 'No', "I don't know", 'Not eligible for coverage / N/A',
```

```
nan], dtype=object)
```

```
1 data2.mental_healthcare_coverage.replace(to_replace=['Not eligible for coverage / N/A'],
2                                           value='No',inplace=True)
3 print(data2.mental_healthcare_coverage.unique())
4 print(data2.mental_healthcare_coverage.value_counts())
```

```
['Yes' 'No' "I don't know" nan]
Yes      531
I don't know  319
No       296
Name: mental_healthcare_coverage, dtype: int64
```

```
1 # openess_of_family_friends column
2 data2.openess_of_family_friends.unique()
```

```
array(['Somewhat open', 'Very open', 'Somewhat not open', 'Neutral',
      'Not applicable to me (I do not have a mental illness)',
      'Not open at all'], dtype=object)
```

```
1 data2.openess_of_family_friends.replace(to_replace=['Not applicable to me (I do not have a
2                                           value="I don't know",inplace=True)
3 data2.openess_of_family_friends.unique()
```

```
array(['Somewhat open', 'Very open', 'Somewhat not open', 'Neutral',
      "I don't know", 'Not open at all'], dtype=object)
```

```
1 print(data2.openess_of_family_friends.value_counts())
```

```
Somewhat open      640
Very open          251
Somewhat not open  214
Neutral            141
I don't know       112
```

Saving...



```
, dtype: int64
```

```
1 # Cleaning the age column remove outliers
2 med_age = data2[(data2['age'] >= 18) | (data2['age'] <= 75)]['age'].median()
3 print(med_age)
4 data2['age'].replace(to_replace = data2[(data2['age'] < 18) | (data2['age'] > 75)]['age'].
5                      value = med_age, inplace = True)
6 data2.age.unique()
```

```
33.0
array([33., 40., 21., 36., 42., 26., 29., 30., 56., 35., 51., 24., 38.,
      44., 27., 55., 22., 25., 28., 23., 32., 31., 43., 37., 39., 45.,
      46., 20., 54., 34., 61., 41., 48., 66., 19., 52., 50., 49., 47.,
      57., 74., 53., 58., 70., 59., 62., 63., 65.]
```



```

1 # gender column
2 data2.gender.unique()

array(['Male', 'male', 'F', 'Transitioned, M2F', 'Other/Transfeminine',
      'M', 'female', 'm', 'Female', 'f', 'non-binary', 'woman', 'male ',
      'Male ', 'Bigender', 'Genderfluid (born female)',
      'male 9:1 female, roughly', 'Male (cis)', 'Other', 'Sex is male',
      'genderqueer', 'Human', 'mail', 'Cis-woman',
      'female-bodied; no feelings about gender', 'Transgender woman',
      'Genderfluid', 'female ', 'Male/genderqueer', 'fem', 'Nonbinary',
      'Female', 'Female ', 'Genderqueer', nan, 'I identify as female.',
      'fm', 'Cis female ', 'female/woman', 'Androgynous', 'man',
      'nb masculine', 'Cisgender Female', 'Woman', 'Cis Male',
      'Female or Multi-Gender Femme', 'Male.', 'Enby', 'Agender',
      'Female (props for making this a freeform field, though)',
      'cis man', 'Female assigned at birth ', 'Cis male', 'Man',
      'none of your business', 'cis male', 'genderqueer woman', 'Queer',
      'Dude', 'Male (trans, FtM)', 'cisdude', 'Genderflux demi-girl',
      'Malr', 'mtf', 'Fluid',
      "I'm a man why didn't you make this a drop down question. You should of asked
      sex? And I would of answered yes please. Seriously how much text can this take? ",
      'M|', 'human', 'Unicorn', 'AFAB', 'MALE'], dtype=object)

1 data2['gender'].replace(to_replace = ['Male', 'male', 'Male ', 'M', 'm',
2   'man', 'Cis male', 'Male.', 'male 9:1 female, roughly', 'Male (cis)', 'Man', 'Sex i
3   'cis male', 'Malr', 'Dude', "I'm a man why didn't you make this a drop down questio
4   'mail', 'M|', 'Male/genderqueer', 'male ',
5   'Cis Male', 'Male (trans, FtM)',
6   'cisdude', 'cis man', 'MALE'], value = 'male', inplace = True)
7 data2['gender'].replace(to_replace = ['Female', 'female', 'I identify as female.', 'female
8   'Female assigned at birth ', 'F', 'Woman', 'fm', 'f', 'Cis female ', 'Transitioned,
9   'Genderfluid (born female)', 'Female or Multi-Gender Femme', 'Female ', 'woman', 'f
10  'Cisgender Female', 'fem', 'Female (props for making this a freeform field, though)
11  'Female', 'Cis-woman', 'female-bodied; no feelings about gender',
12  'AFAB'], value = 'female', inplace = True)
13 data2['gender'].replace(to_replace = ['Bigender', 'non-binary', 'Other/Transfeminine',
    'nb masculine',
    'genderqueer', 'Human', 'Genderfluid',
16  'Enby', 'genderqueer woman', 'mtf', 'Queer', 'Agender', 'Fluid',
17  'Nonbinary', 'human', 'Unicorn', 'Genderqueer',
18  'Genderflux demi-girl', 'Transgender woman'], value = 'other', inplace = True)

```

```

1 data2.gender.unique()

array(['male', 'female', 'other', nan], dtype=object)

```

```
1 data2.gender.value_counts()
```

```

male      1060
female    343

```

```

other          27
Name: gender, dtype: int64

1 ## Cleaning the role_in_company
2 tech_list = []
3 tech_list.append(data2[data2['role_in_company'].str.contains('Back-end')]['role_in_company'])
4 tech_list.append(data2[data2['role_in_company'].str.contains('Front-end')]['role_in_company'])
5 tech_list.append(data2[data2['role_in_company'].str.contains('Dev')]['role_in_company'])
6 tech_list.append(data2[data2['role_in_company'].str.contains('DevOps')]['role_in_company'])
7 flat_list = [item for sublist in tech_list for item in sublist]
8 flat_list = list(dict.fromkeys(flat_list))

1 ## Replace tech role=1 and other=0 in a new tech role operation
2 data2['tech_role']=data2['role_in_company']
3 data2['tech_role'].replace(to_replace=flat_list,value=1,inplace=True)
4 remain_list=data2['tech_role'].unique()[1:]
5 data2['tech_role'].replace(to_replace=remain_list,value=0,inplace=True)

1 data2.tech_role.value_counts()

1    1045
0     388
Name: tech_role, dtype: int64

1 data2=data2.drop(['role_in_company'],axis=1)

```

▼ Handling Missing values

```

1 data3=pd.concat([data2['type_of_disorder'],data2['US state'],data2['US state work']],axis=
2 print(data3.info())

```

Saving...

```

RangeIndex: 1433 entries, 0 to 1432
Data columns (total 3 columns):
#   Column                Non-Null Count  Dtype
---  -
0   type_of_disorder      711 non-null   object
1   US state               840 non-null   object
2   US state work         851 non-null   object
dtypes: object(3)
memory usage: 33.7+ KB
None

```

```
1 data2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 1433 entries, 0 to 1432
Data columns (total 24 columns):
#   Column                                                                 Non-Null Count  Dtype
---  -
0   self_employed                                                         1433 non-null   int64
1   no_of_employees                                                       1146 non-null   object
2   tech_company                                                          1146 non-null   float64
3   mental_healthcare_coverage                                           1146 non-null   object
4   knowledge_about_mental_healthcare_options_workplace                 1013 non-null   object
5   employer_discussed_mental_health                                     1146 non-null   object
6   employer_offer_resources_to_learn_about_mental_health               1146 non-null   object
7   medical_leave_from_work                                              1146 non-null   object
8   comfortable_discussing_with_coworkers                                1146 non-null   object
9   employer_take_mental_health_seriously                               1146 non-null   object
10  openness_of_family_friends                                           1433 non-null   object
11  family_history_mental_illness                                         1433 non-null   object
12  mental_health_disorder_past                                          1433 non-null   object
13  currently_mental_health_disorder                                     1433 non-null   object
14  diagnosed_mental_health_condition                                    1433 non-null   object
15  treatment_from_professional                                          1433 non-null   int64
16  while_effective_treatment_mental_health_issue_interferes_work       1433 non-null   object
17  while_not_effective_treatment_interferes_work                       1433 non-null   object
18  age                                                                    1433 non-null   float64
19  gender                                                                1430 non-null   object
20  country                                                                1433 non-null   object
21  country work                                                          1433 non-null   object
22  work_remotely                                                         1433 non-null   object
23  tech_role                                                             1433 non-null   int64
dtypes: float64(2), int64(3), object(19)
memory usage: 268.8+ KB
```




```
1 from sklearn.impute import SimpleImputer
2 imp = SimpleImputer(missing_values=np.nan, strategy='most_frequent')
3 imp.fit(data2)
4 imp_data=pd.DataFrame(data=imp.transform(data2),columns=data2.columns)
```

Saving...

✕

```
2 data4.isnull().sum().to_frame()
axis=1)
```

| | | |
|---|---|---|
| | 0 |  |
| self_employed | 0 | |
| no_of_employees | 0 | |
| tech_company | 0 | |
| mental_healthcare_coverage | 0 | |
| knowledge_about_mental_healthcare_options_workplace | 0 | |
| employer_discussed_mental_health | 0 | |
| employer_offer_resources_to_learn_about_mental_health | 0 | |
| medical_leave_from_work | 0 | |
| comfortable_discussing_with_coworkers | 0 | |
| employer_take_mental_health_seriously | 0 | |
| openess_of_family_friends | 0 | |
| family_history_mental_illness | 0 | |
| mental_health_disorder_past | 0 | |
| currently_mental_health_disorder | 0 | |
| diagnosed_mental_health_condition | 0 | |
| treatment_from_professional | 0 | |
| while_effective_treatment_mental_health_issue_interferes_work | 0 | |
| while_not_effective_treatment_interferes_work | 0 | |
| age | 0 | |
| gender | 0 | |

Saving...

✕

Data Analysis)

work_remotely0

Questions with regard to the Target:

type of disorder722

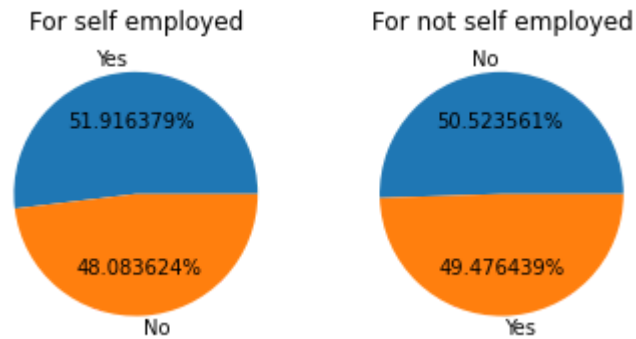
```
1 # 1. If a person is self employed then is there a higher chance of mental disorder?
2
3 plt.subplot(1,2,1)
4 plt.title("For self employed")
5 plt.pie(data4[data4.self_employed==1]['diagnosed_mental_health_condition'].value_counts(),
6
7 plt.subplot(1,2,2)
```

```

8 plt.title("For not self employed")
9 plt.pie(data4[data4.self_employed==0]['diagnosed_mental_health_condition'].value_counts(),

([<matplotlib.patches.Wedge at 0x7fd81da7e3d0>,
  <matplotlib.patches.Wedge at 0x7fd81da7e9d0>],
 [Text(-0.018092161764598828, 1.0998512052467295, 'No'),
  Text(0.018092161764598828, -1.0998512052467295, 'Yes')],
 [Text(-0.00986845187159936, 0.5999188392254888, '50.523561%'),
  Text(0.00986845187159936, -0.5999188392254888, '49.476439%')])

```

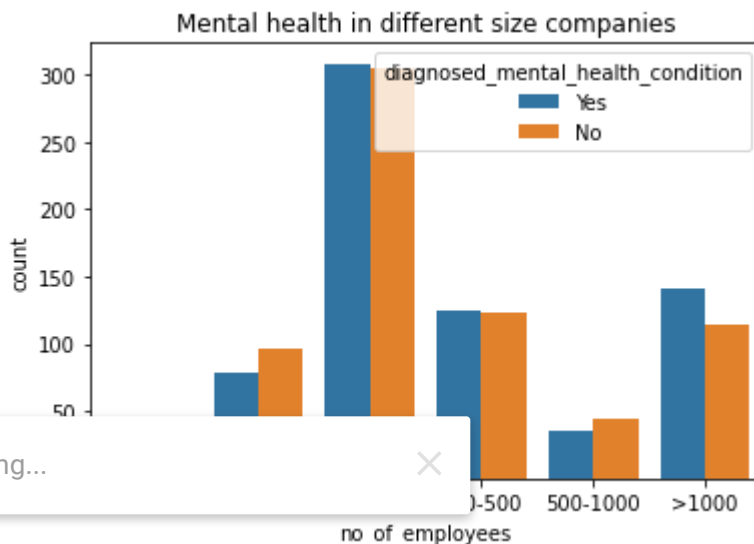


```

1 # 2. Does big size of the company affect your mental health condition adversely?
2 import seaborn as sns
3 sns.countplot(data=data4,x='no_of_employees',hue='diagnosed_mental_health_condition')
4 plt.title('Mental health in different size companies')

```

```
Text(0.5, 1.0, 'Mental health in different size companies')
```

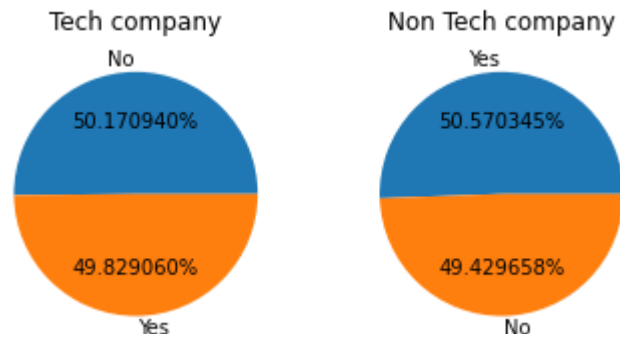


```

1 # 3. Does working in a tech company affect adversely to your mental well being?
2
3 plt.subplot(1,2,1)
4 plt.title("Tech company")
5 plt.pie(data4[data4.tech_company==1]['diagnosed_mental_health_condition'].value_counts(),a
6
7 plt.subplot(1,2,2)
8 plt.title("Non Tech company")
9 plt.pie(data4[data4.tech_company==0]['diagnosed_mental_health_condition'].value_counts(),a

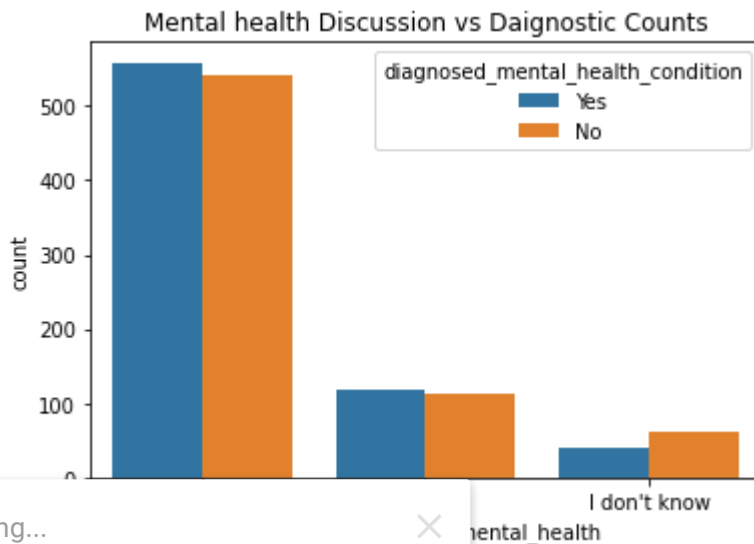
```

```
([<matplotlib.patches.Wedge at 0x7fd81d927890>,
  <matplotlib.patches.Wedge at 0x7fd81d933110>],
 [Text(-0.019708651065785975, 1.0998234263158642, 'Yes'),
  Text(0.019708754038686934, -1.0998234244706022, 'No')],
 [Text(-0.01075017330861053, 0.5999036870813804, '50.570345%'),
  Text(0.010750229475647417, -0.5999036860748739, '49.429658%')])
```



```
1 # 4. Does the employers discussion on mental health reduces the chance of getting postive
2
3 sns.countplot(data=data4,x='employer_discussed_mental_health ',hue='diagnosed_mental_healt
4 plt.title('Mental health Discussion vs Daignostic Counts')
5
```

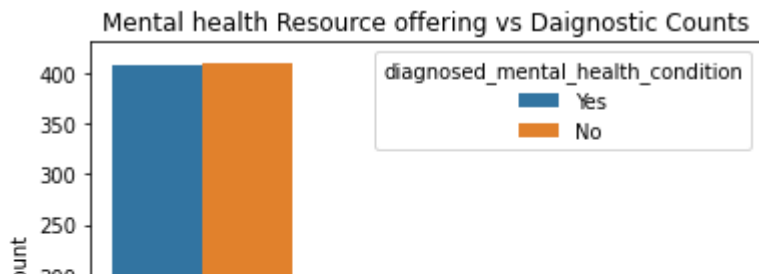
```
Text(0.5, 1.0, 'Mental health Discussion vs Daignostic Counts')
```



Saving...

```
1 # 5. Will offering more options to learn about mental health reduces the chance of getting
2
3 sns.countplot(data=data4,x='employer_offer_resources_to_learn_about_mental_health',hue='di
4 plt.title('Mental health Resource offering vs Daignostic Counts')
```

Text(0.5, 1.0, 'Mental health Resource offering vs Daignostic Counts')



1 # 6. Does providing no leaves increases the less reporting of mental health issues?

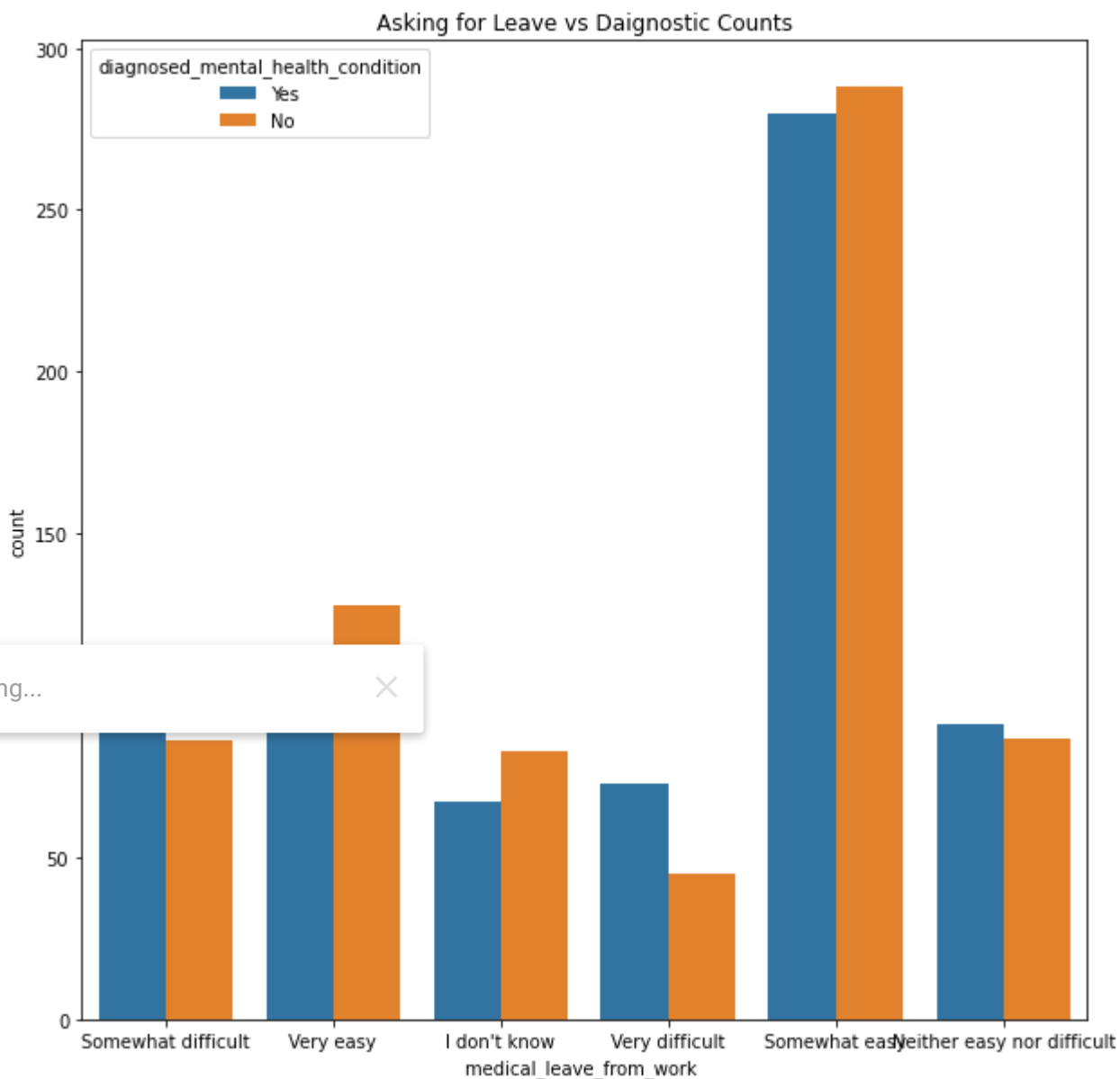
2

3 plt.figure(figsize=(10,10))

4 sns.countplot(data=data4,x='medical_leave_from_work ',hue='diagnosed_mental_health_conditi

5 plt.title('Asking for Leave vs Daignostic Counts')

Text(0.5, 1.0, 'Asking for Leave vs Daignostic Counts')

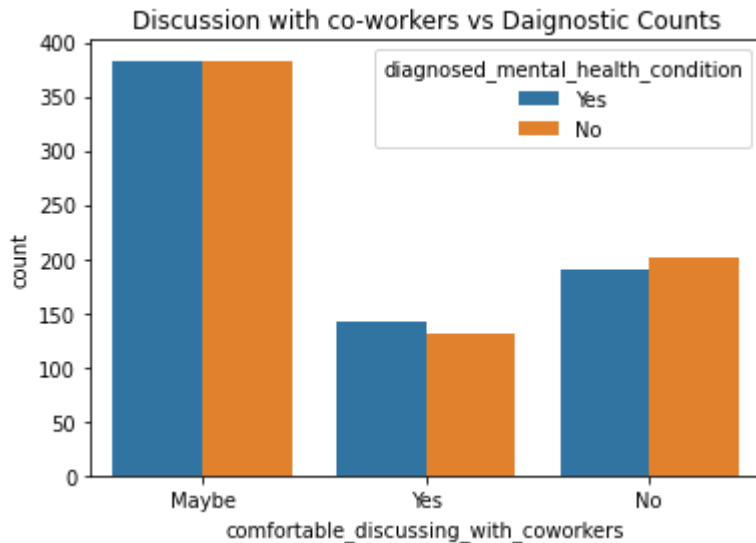


1 # 7. Does discussion with coworkers about mental health care reduces the chance of positiv

2 sns.countplot(data=data4,x='comfortable_discussing_with_coworkers',hue='diagnosed_mental_h

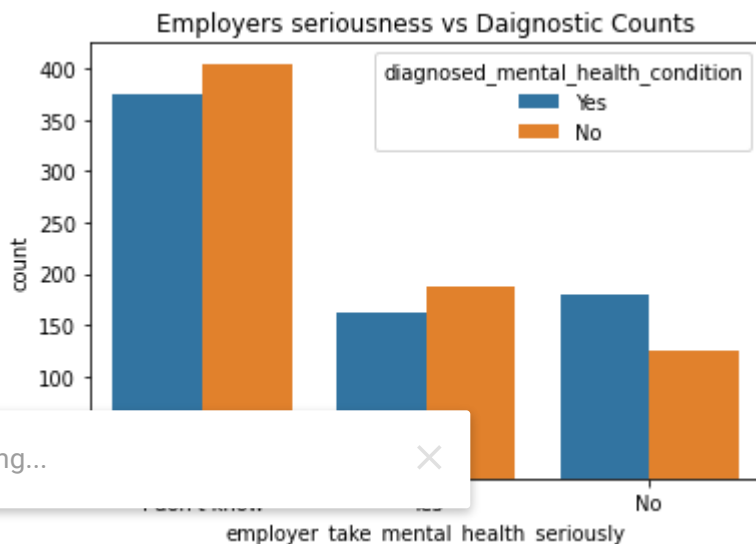
```
3 plt.title('Discussion with co-workers vs Daignostic Counts')
```

```
Text(0.5, 1.0, 'Discussion with co-workers vs Daignostic Counts')
```



```
1 # 8. If Employer takes mental health seriously, then will it reduce the chance of positive
2 sns.countplot(data=data4,x='employer_take_mental_health_seriously',hue='diagnosed_mental_h
3 plt.title('Employers seriousness vs Daignostic Counts')
```

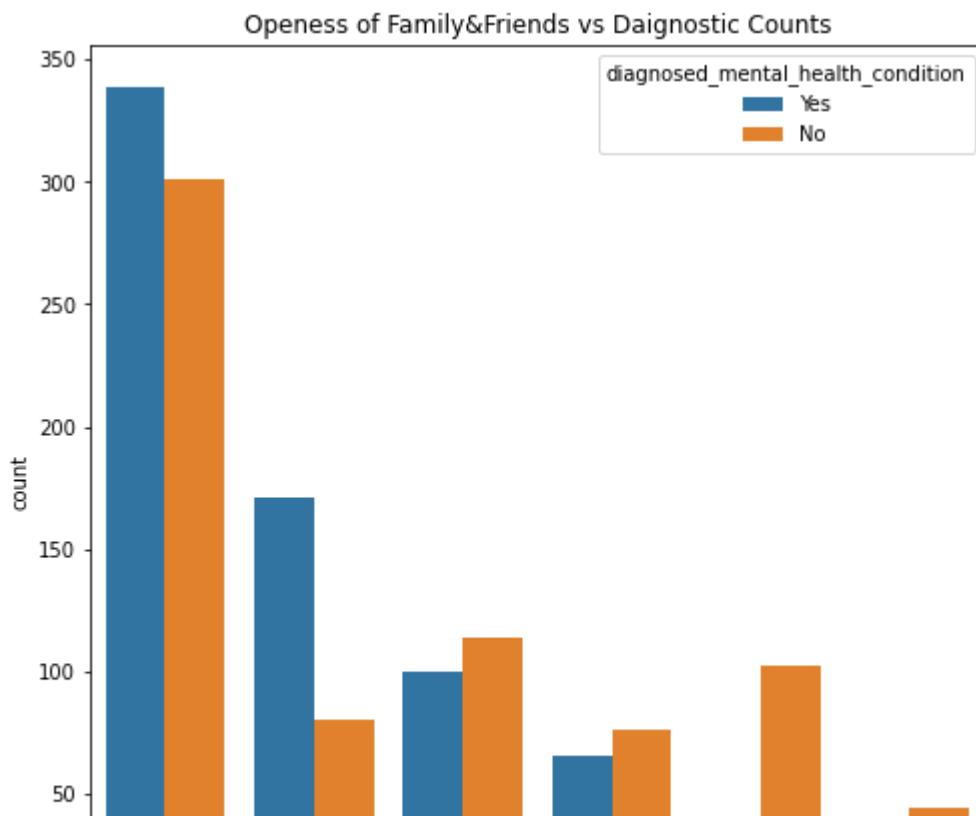
```
Text(0.5, 1.0, 'Employers seriousness vs Daignostic Counts')
```



Saving...

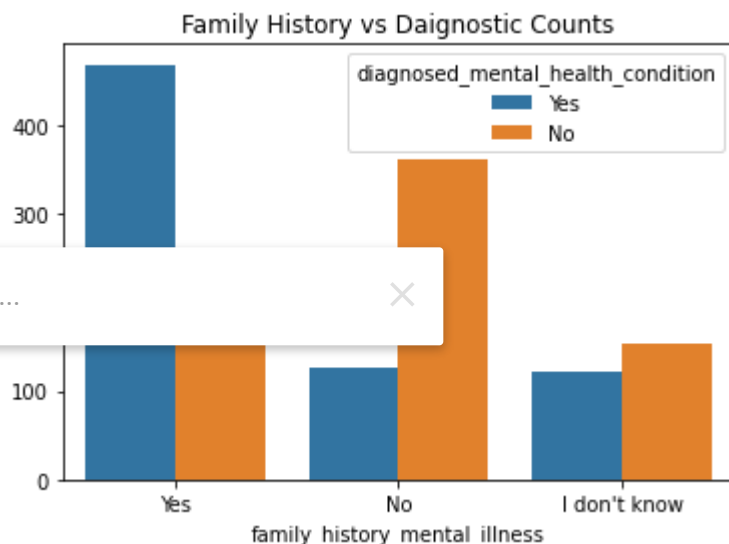
```
1 # 9. If family friends are open about the mental health then will it reduce the positive d
2
3 plt.figure(figsize=(8,8))
4 sns.countplot(data=data4,x='openess_of_family_friends',hue='diagnosed_mental_health_condit
5 plt.title('Openess of Family&Friends vs Daignostic Counts')
```


Text(0.5, 1.0, 'Openess of Family&Friends vs Daignostic Counts')



```
1 # 10. What are the chances that if a person having family history of mental illness then h
2 sns.countplot(data=data4,x='family_history_mental_illness',hue='diagnosed_mental_health_co
3 plt.title('Family History vs Daignostic Counts')
```

Text(0.5, 1.0, 'Family History vs Daignostic Counts')



```
1 # 11. Does having mental illness of the past affect the diagonosis?
2
3 plt.figure(figsize=(10,10))
4 plt.subplot(1,2,1)
5 plt.title("Had past mental illness")
6 plt.pie(data4[data4.mental_health_disorder_past=='Yes']['diagnosed_mental_health_condition
```

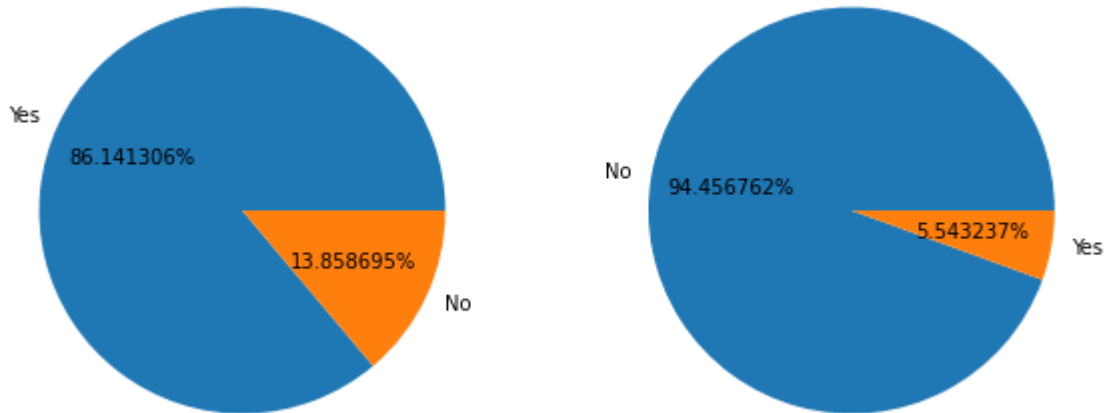
```

7
8 plt.subplot(1,2,2)
9 plt.title("Did Not had mental illness of past")
10 plt.pie(data4[data4.mental_health_disorder_past=='No']['diagnosed_mental_health_condition']

([<matplotlib.patches.Wedge at 0x7fd81d4ebb10>,
 <matplotlib.patches.Wedge at 0x7fd81d4f83d0>],
 [Text(-1.0833623634614518, 0.19059378120814285, 'No'),
  Text(1.083362359000282, -0.19059380656606933, 'Yes')],
 [Text(-0.5909249255244282, 0.10396024429535064, '94.456762%'),
  Text(0.5909249230910628, -0.10396025812694691, '5.543237%')])

```

Had past mental illness Did Not had mental illness of past



```

1 # 12. Is self proclaimed mental health disorders increases the chances of being diagonised
2 sns.countplot(data=data4,x='currently_mental_health_disorder',hue='diagnosed_mental_health
3 plt.title('Self proclaimed disorder vs Daignostic Counts')

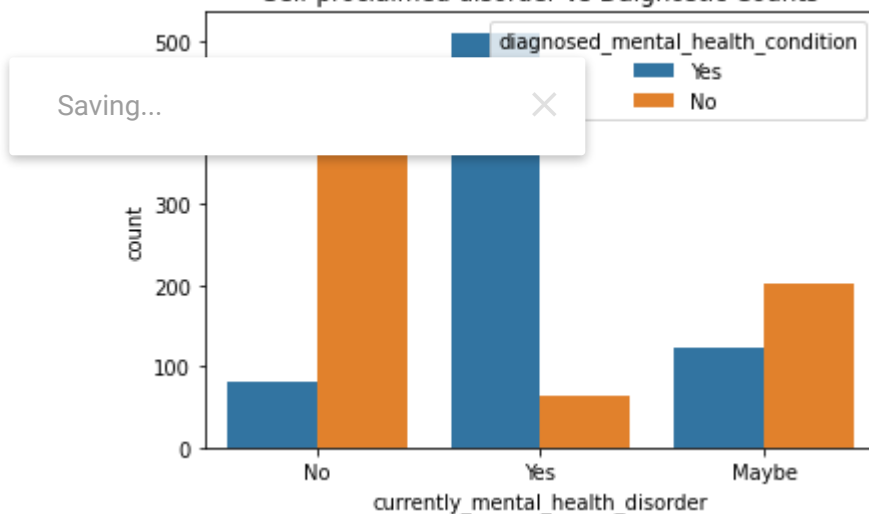
```

```

Text(0.5, 1.0, 'Self proclaimed disorder vs Daignostic Counts')

```

Self proclaimed disorder vs Daignostic Counts



```

1 # 13. How many of those who has diagonised positively will seek help of professional?
2

```

```

3 plt.figure(figsize=(10,10))
4 plt.subplot(1,2,1)
5 plt.title("Taking Help from professional")
6 plt.pie(data4[data4.treatment_from_professional==1]['diagnosed_mental_health_condition'].v
7
8 plt.subplot(1,2,2)
9 plt.title("Not taking help from professional")
10 plt.pie(data4[data4.treatment_from_professional==0]['diagnosed_mental_health_condition'].v

```

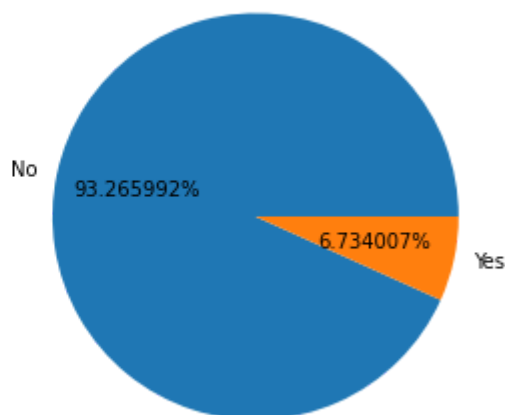
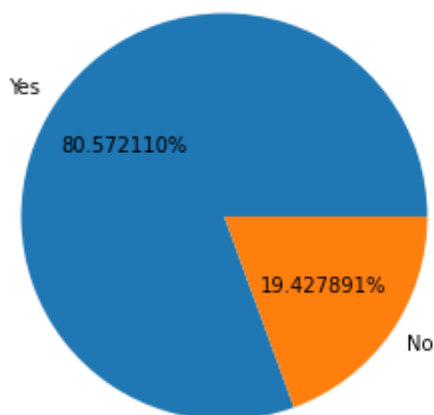
```

([<matplotlib.patches.Wedge at 0x7fd81d3fc810>,
 <matplotlib.patches.Wedge at 0x7fd81d409090>],
 [Text(-1.0754761146685294, 0.2309786283998678, 'No'),
  Text(1.0754761092620841, -0.23097865357320332, 'Yes')],
 [Text(-0.5866233352737433, 0.12598834276356424, '93.265992%'),
  Text(0.5866233323247732, -0.1259883564944745, '6.734007%')])

```

Taking Help from professional

Not taking help from professional



1 # 14. If one is diagnosed positive how effective and not effective medication affecting t
2

```

3 plt.figure(figsize=(10,10))
4 plt.subplot(1,2,1)

```

Saving...



```

5 plt.title("How effective medication")

```

```

6 plt.pie(data4[data4.diagnosed_mental_health_condition=='Yes']['while_effective_treatment_m

```

```

8 plt.subplot(1,2,2)
9 plt.title("Positive diagnosed under not-effective medication")
10 plt.pie(data4[data4.diagnosed_mental_health_condition=='Yes']['while_not_effective_treatme

```

```
([<matplotlib.patches.Wedge at 0x7fd81d3a3e90>,
 <matplotlib.patches.Wedge at 0x7fd81d333650>,
 <matplotlib.patches.Wedge at 0x7fd81d33c090>,
 <matplotlib.patches.Wedge at 0x7fd81d33c950>,
 <matplotlib.patches.Wedge at 0x7fd81d3483d0>],
 [Text(-0.387286831783129, 1.0295673411328599, 'Often'),
 Text(0.11082048010455826, -1.0944034087983259, 'Sometimes'),
 Text(1.0153415504745185, -0.4231802640483144, 'Rarely'),
 Text(1.0863057567337249, -0.1730312194003421, 'Not applicable to me'),
 Text(1.0997352976415675, -0.02413037756051461, 'Never')],
 [Text(-0.21124736279079762, 0.5615821860724689, '61.452514%'),
 Text(0.06044753460248632, -0.5969473138899958, '30.307263%'),
 Text(0.5538226638951919, -0.2308255985718078, '3.910615%'),
 Text(0.5925304127638499, -0.09438066512745932, '3.631285%'),
 Text(0.5998556168954003, -0.013162024123917057, '0.698324%')])
```

Positive diagnosed under effective medication Positive diagnosed under not-effective medication



```
1 # 15. Is the chances of getting positively diagnosed increases with age?
2
3 plt.figure(figsize=(7,7))
4 plt.subplot(1,2,1)
5 plt.hist(data4[data4.diagnosed_mental_health_condition=='Yes']['age'],bins=5)
6 plt.title("Positive Diagonosis with age")
7
8 plt.subplot(1,2,2)
9 plt.hist(data4[data4.diagnosed_mental_health_condition=='No']['age'],bins=5)
10 plt.title("Negetive Diagonosis with age")
```

Saving...

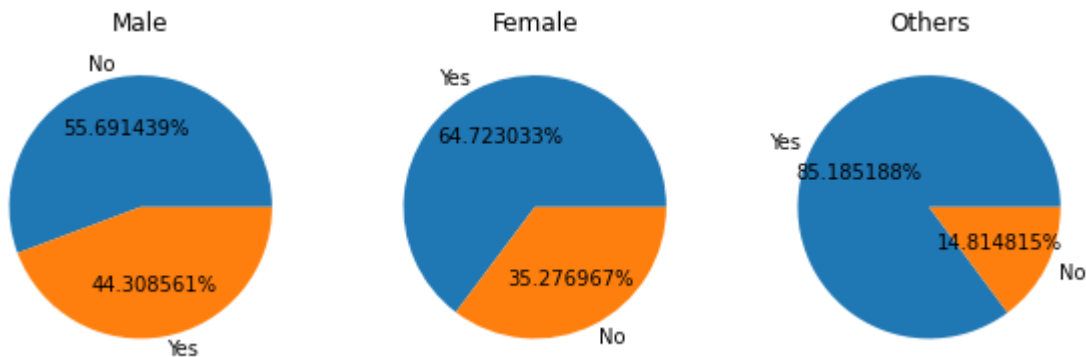


Text(0.5, 1.0, 'Negative Diagonosis with age')



```
1 # 16. Is there any chance that a other gender category is more positively diagnosed than
2
3 plt.figure(figsize=(10,10))
4 plt.subplot(1,3,1)
5 plt.title("Male")
6 plt.pie(data4[data4.gender=='male']['diagnosed_mental_health_condition'].value_counts(), a
7
8 plt.subplot(1,3,2)
9 plt.title("Female")
10 plt.pie(data4[data4.gender=='female']['diagnosed_mental_health_condition'].value_counts(),
11
12 plt.subplot(1,3,3)
13 plt.title("Others")
14 plt.pie(data4[data4.gender=='other']['diagnosed_mental_health_condition'].value_counts(),
```

```
([<matplotlib.patches.Wedge at 0x7fd81d183250>,
  <matplotlib.patches.Wedge at 0x7fd81d183bd0>],
 [Text(-0.9829959488654123, 0.49367900959448124, 'Yes'),
  Text(0.9829959950869729, -0.4936789175597559, 'No')],
 [Text(-0.536179608472043, 0.2692794597788079, '85.185188%'),
  Text(0.5361796336838033, -0.26927940957804863, '14.814815%')])
```



Saving...

```
!pip install chart-studio
```

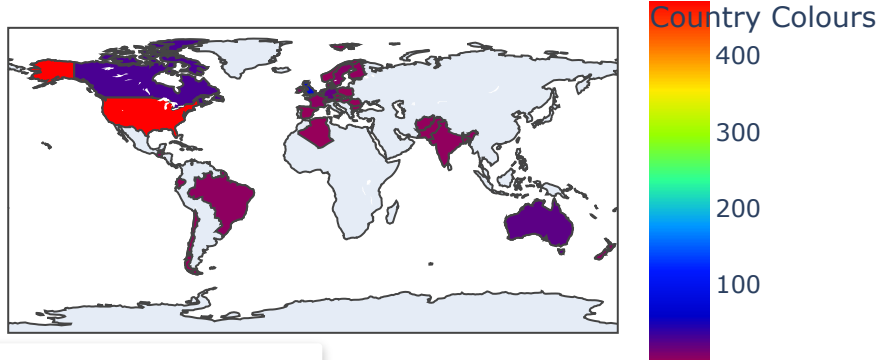
Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public>
 Requirement already satisfied: chart-studio in /usr/local/lib/python3.7/dist-packages (1
 Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from chart
 Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from
 Requirement already satisfied: plotly in /usr/local/lib/python3.7/dist-packages (from ch
 Requirement already satisfied: retrying>=1.3.3 in /usr/local/lib/python3.7/dist-packages
 Requirement already satisfied: tenacity>=6.2.0 in /usr/local/lib/python3.7/dist-packages
 Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packag
 Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packa
 Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (f
 Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr/local/lib

```

1 # 17. Country wise positive disorder cases?
2 import chart_studio.plotly as py
3 import plotly.graph_objs as gobj
4 from plotly.offline import download_plotlyjs,init_notebook_mode,plot,iplot

1 data =dict( type = 'choropleth',
2             locations = list(data4[data4.diagnosed_mental_health_condition=='Yes']['country']
3             locationmode = 'country names',
4             colorscale= 'Rainbow',
5             z=list(data4[data4.diagnosed_mental_health_condition=='Yes']['country'].value_
6             colorbar = {'title':'Country Colours', 'len':200,'lenmode':'pixels' })
7 layout = dict(geo = {'scope':'world'})
8 col_map=gobj.Figure(data = [data],layout = layout)
9 iplot(col_map)

```



Saving...

```

1 # 18. Does being involved in tech role increases chances of diagnosed positive?
2
3 plt.figure(figsize=(10,10))
4 plt.subplot(1,2,1)
5 plt.title("Tech Role")
6 plt.pie(data4[data4.tech_role==1]['diagnosed_mental_health_condition'].value_counts(), aut

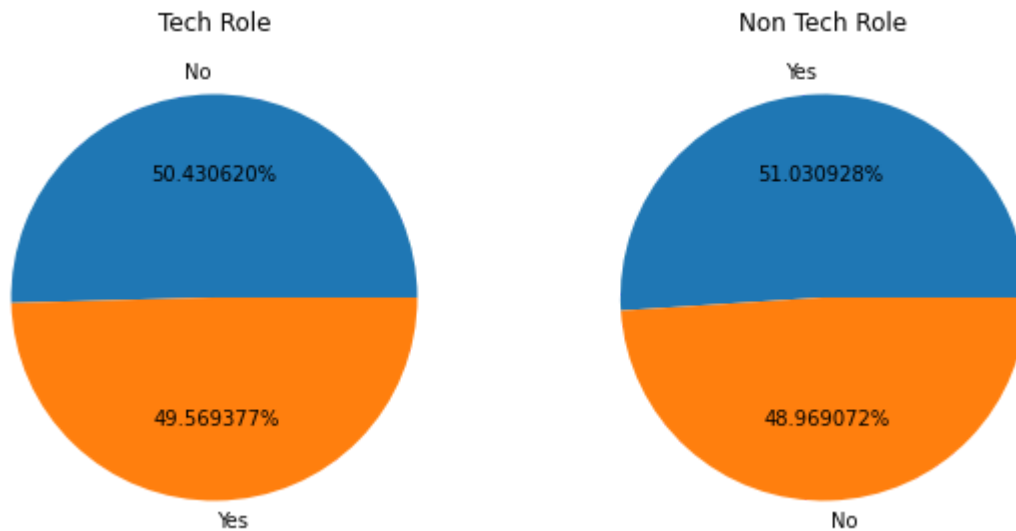
```

```

7
8 plt.subplot(1,2,2)
9 plt.title("Non Tech Role")
10 plt.pie(data4[data4.tech_role==0]['diagnosed_mental_health_condition'].value_counts(), aut

([<matplotlib.patches.Wedge at 0x7fd81d0d8490>,
 <matplotlib.patches.Wedge at 0x7fd81d0d8e10>],
 [Text(-0.03562008250107137, 1.0994231258813036, 'Yes'),
  Text(0.03562008250107123, -1.0994231258813039, 'No')],
 [Text(-0.01942913590967529, 0.5996853413898019, '51.030928%'),
  Text(0.019429135909675214, -0.599685341389802, '48.969072%')])

```



```

1 # 19. Will working remotely helps to better the mental health condition?
2
3 plt.figure(figsize=(10,10))
4 plt.subplot(1,3,1)
5 plt.title("Always Remote work")
6 plt.pie(data4[data4.work_remotely=='Always']['diagnosed_mental_health_condition'].value_co
7

```

Saving...

```

11 plt.pie(data4[data4.work_remotely=='Sometimes']['diagnosed_mental_health_condition'].value
12 plt.subplot(1,3,3)
13 plt.title("Never Remote work")
14 plt.pie(data4[data4.work_remotely=='Never']['diagnosed_mental_health_condition'].value_cou

```

```
([<matplotlib.patches.Wedge at 0x7fd81d00b6d0>,
  <matplotlib.patches.Wedge at 0x7fd81d016050>],
 [Text(-0.15000672147010105, 1.089723810657449, 'No'),
  Text(0.15000661944279384, -1.089723824702087, 'Yes')],
 [Text(-0.08182184807460056, 0.5943948058131538, '54.354352%'),
  Text(0.08182179242334207, -0.5943948134738656, '45.645645%')])
```

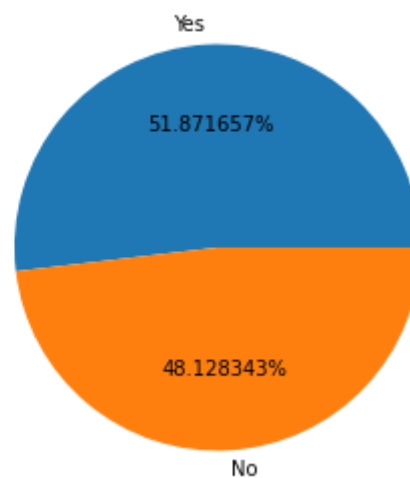
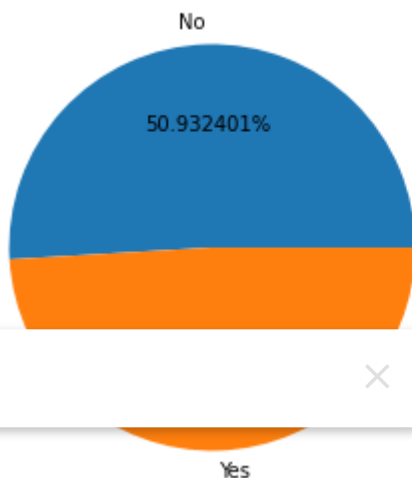
Always Remote work Sometime Remote work Never Remote work

No Yes No

```
1 # 20. Does a person in tech role in tech company has higher chance of diagnosis than a te
2
3 plt.figure(figsize=(10,10))
4 plt.subplot(1,2,1)
5 plt.title("Tech Role in tech company")
6 plt.pie(data4[(data4.tech_role==1) & (data4.tech_company==1)]['diagnosed_mental_health_con
7
8 plt.subplot(1,2,2)
9 plt.title("Tech Role in non tech company")
10 plt.pie(data4[(data4.tech_role==1) & (data4.tech_company==0)]['diagnosed_mental_health_con
```

```
([<matplotlib.patches.Wedge at 0x7fd81cfa0c90>,
  <matplotlib.patches.Wedge at 0x7fd81cf2d310>],
 [Text(-0.06464257105613563, 1.0980989654886542, 'Yes'),
  Text(0.06464257105613574, -1.0980989654886542, 'No')],
 [Text(-0.03525958421243761, 0.5989630720847204, '51.871657%'),
  Text(0.035259584212437675, -0.5989630720847204, '48.128343%')])
```

Tech Role in tech company Tech Role in non tech company



Saving...



Questions other than Target

```
1 # 1. For self employed does the past mental disorder more than those who are not self empl
2
3 plt.subplot(1,2,1)
4 plt.title("Self Employed")
5 plt.pie(data4[data4.self_employed==1]['mental_health_disorder_past'].value_counts(), autop
6
7 plt.subplot(1,2,2)
```

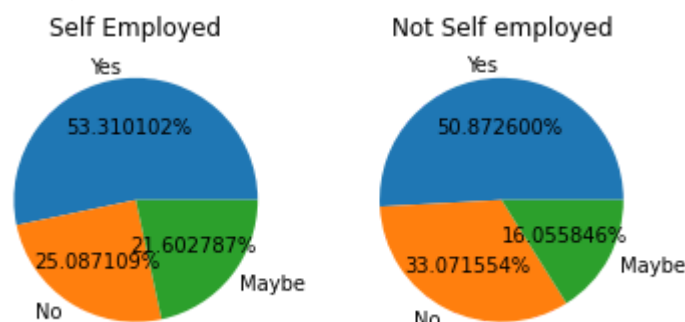


```

8 plt.title("Not Self employed")
9 plt.pie(data4[data4.self_employed==0]['mental_health_disorder_past'].value_counts(), autop

([<matplotlib.patches.Wedge at 0x7fd81cebead0>,
 <matplotlib.patches.Wedge at 0x7fd81cecc350>,
 <matplotlib.patches.Wedge at 0x7fd81cecc4d0>],
 [Text(-0.030151117190896082, 1.0995866996886334, 'Yes'),
 Text(-0.5050235705581543, -0.9772160422243861, 'No'),
 Text(0.9630061253565793, -0.5316194151135837, 'Maybe')],
 [Text(-0.01644606392230695, 0.5997745634665272, '50.872600%'),
 Text(-0.27546740212262955, -0.5330269321223924, '33.071554%'),
 Text(0.5252760683763159, -0.28997422642559106, '16.055846%')])

```



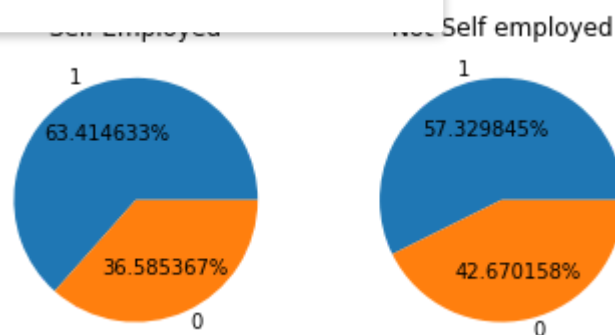
```

1 # 2. Does self employed people shy away to seek help?
2 plt.subplot(1,2,1)
3 plt.title("Self Employed")
4 plt.pie(data4[data4.self_employed==1]['treatment_from_professional'].value_counts(), autop
5
6 plt.subplot(1,2,2)
7 plt.title("Not Self employed")
8 plt.pie(data4[data4.self_employed==0]['treatment_from_professional'].value_counts(), autop

([<matplotlib.patches.Wedge at 0x7fd81cde70d0>,
 <matplotlib.patches.Wedge at 0x7fd81cde7ad0>],
 [Text(-0.25106860689769006, 1.0709643106240532, '1'),
 Text(0.2510687071686046, -1.0709642871173088, '0')],
 [Text(0.841623512494835, 0.57329845, '57.329845%'),
 Text(0.584162338427623, 0.42670158, '42.670158%')])

```

Saving...

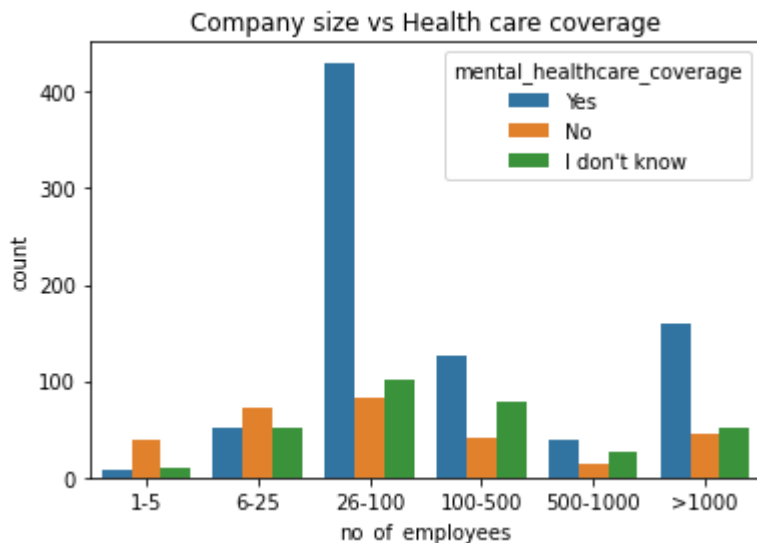


```

1 # 3. Is the large company is more serious about mental health than small companies?
2 sns.countplot(data=data4,x='no_of_employees',hue='mental_healthcare_coverage')
3 plt.title('Company size vs Health care coverage')

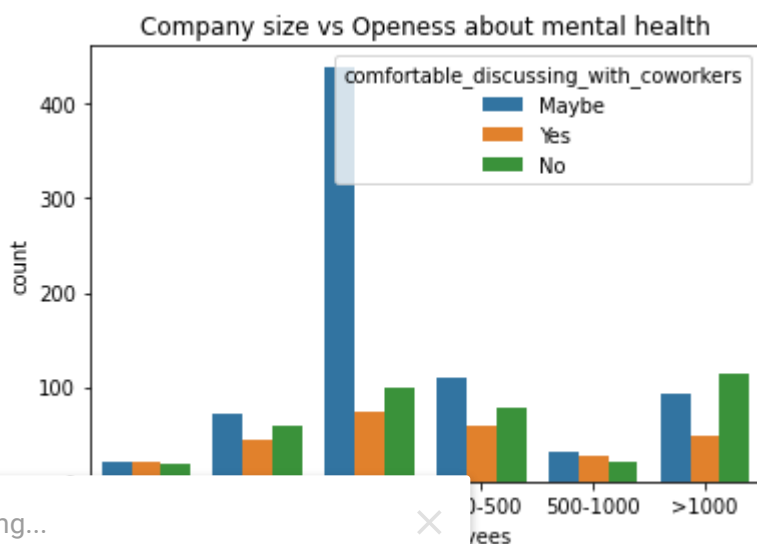
```

```
Text(0.5, 1.0, 'Company size vs Health care coverage')
```



```
1 # 4. Does openness about the mental health varies with size of the companies?
2 sns.countplot(data=data4,x='no_of_employees',hue='comfortable_discussing_with_coworkers')
3 plt.title('Company size vs Openess about mental health')
```

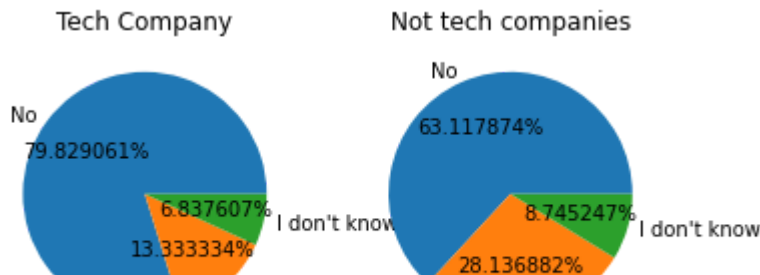
```
Text(0.5, 1.0, 'Company size vs Openess about mental health')
```



Saving...

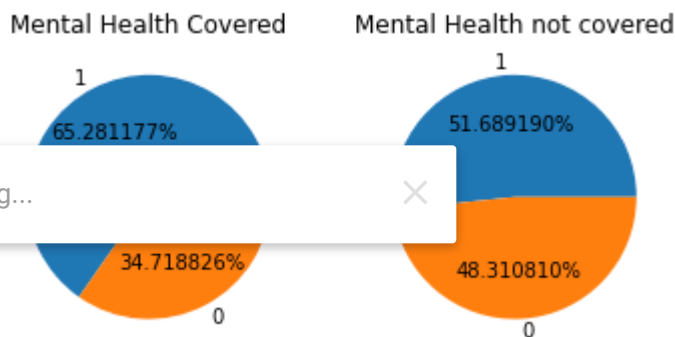
```
1 # 5. Does the tech companies take mental health seriously than other non tech companies?
2 plt.subplot(1,2,1)
3 plt.title("Tech Company")
4 plt.pie(data4[data4.tech_company==1]['employer_discussed_mental_health '].value_counts(),
5
6 plt.subplot(1,2,2)
7 plt.title("Not tech companies")
8 plt.pie(data4[data4.tech_company==0]['employer_discussed_mental_health '].value_counts(),
9
```

```
([<matplotlib.patches.Wedge at 0x7fd81cc30ad0>,
 <matplotlib.patches.Wedge at 0x7fd81cc3e350>,
 <matplotlib.patches.Wedge at 0x7fd81cc3e4d0>],
 [Text(-0.44059807728808054, 1.0079054193177288, 'No'),
 Text(0.1506324301763201, -1.0896374952153474, 'Yes'),
 Text(1.058745387606376, -0.2984262123578033, "I don't know")],
 [Text(-0.24032622397531664, 0.5497665923551248, '63.117874%'),
 Text(0.08216314373253823, -0.5943477246629167, '28.136882%'),
 Text(0.5774974841489323, -0.16277793401334723, '8.745247%')])
```



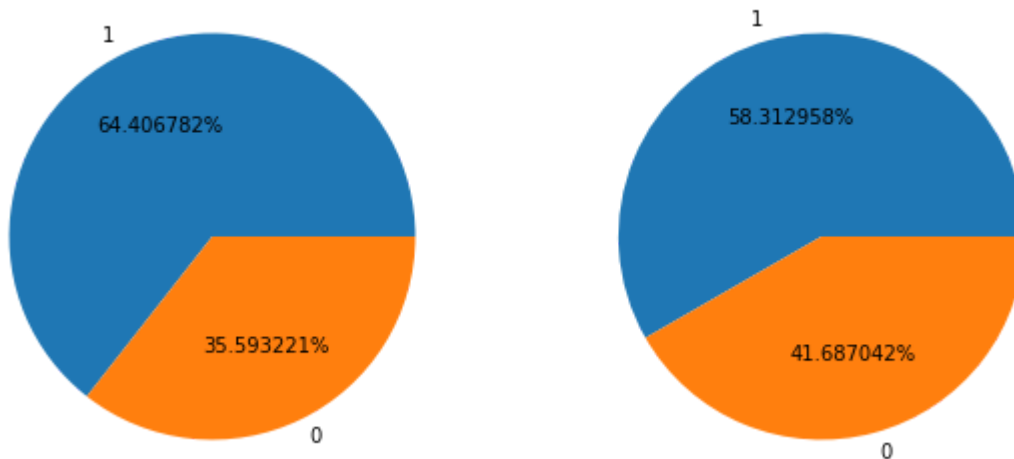
```
1 # 6. Does providing more health care benefits provide seeking for professional health?
2 plt.subplot(1,2,1)
3 plt.title("Mental Health Covered")
4 plt.pie(data4[data4.mental_healthcare_coverage=='Yes']['treatment_from_professional'].valu
5
6 plt.subplot(1,2,2)
7 plt.title("Mental Health not covered")
8 plt.pie(data4[data4.mental_healthcare_coverage=='No']['treatment_from_professional'].value
```

```
([<matplotlib.patches.Wedge at 0x7fd81cbd6250>,
 <matplotlib.patches.Wedge at 0x7fd81cbd6c50>],
 [Text(-0.05834680990910332, 1.0984514781151833, '1'),
 Text(0.058346809909102695, -1.0984514781151833, '0')],
 [Text(-0.03182553267769272, 0.5991553516991909, '51.689190%'),
 Text(0.03182553267769238, -0.5991553516991909, '48.310810%')])
```



```
1 # 7. Does providing more information about mental health increase help seeking behaviour?
2 plt.figure(figsize=(10,10))
3 plt.subplot(1,2,1)
4 plt.title("Mental Healthcare options Knowledge")
5 plt.pie(data4[data4.employer_offer_resources_to_learn_about_mental_health=='Yes']['treatme
6
7 plt.subplot(1,2,2)
8 plt.title("Mental Healthcare options no knowledge")
9 plt.pie(data4[data4.employer_offer_resources_to_learn_about_mental_health=='No']['treatmen
```

```
([<matplotlib.patches.Wedge at 0x7fd81caf1350>,
<matplotlib.patches.Wedge at 0x7fd81caf1d50>],
[Text(-0.2840207786894458, 1.062700426871393, '1'),
Text(0.2840207786894455, -1.062700426871393, '0')],
[Text(-0.1549204247396977, 0.579654778293487, '58.312958%'),
Text(0.1549204247396975, -0.5796547782934871, '41.687042%')])
Mental Healthcare options Knowledge      Mental Healthcare options no knowledge
```



```
1 # 8. The family where there are history mental health issuses are they open about discussi
2 plt.figure(figsize=(10,10))
3 plt.subplot(1,2,1)
4 plt.title("Having family history mental illness")
5 plt.pie(data4[data4.family_history_mental_illness=='Yes']['openess_of_family_friends'].valu
6
7 plt.subplot(1,2,2)
8 plt.title("No family history of mental illness")
9 plt.pie(data4[data4.family_history_mental_illness=='No']['openess_of_family_friends'].valu
```

Saving...



```
([<matplotlib.patches.Wedge at 0x7fd81d7b9610>,
 <matplotlib.patches.Wedge at 0x7fd81d979510>,
 <matplotlib.patches.Wedge at 0x7fd81d76f390>,
 <matplotlib.patches.Wedge at 0x7fd81d742e50>,
 <matplotlib.patches.Wedge at 0x7fd81d6fee50>,
 <matplotlib.patches.Wedge at 0x7fd81d706890>],
 [Text(0.21807097813164583, 1.078167449191779, 'Somewhat open'),
 Text(-1.09854150500076, -0.05662651137643169, 'Very open'),
 Text(-0.6373633449424054, -0.8965310739309755, "I don't know"),
 Text(0.266419245583799, -1.0672491675248847, 'Somewhat not open'),
 Text(0.9165788254046539, -0.6081802831560928, 'Neutral'),
 Text(1.0908949119861056, -0.14123841900427447, 'Not open at all')],
 [Text(0.11894780625362497, 0.5880913359227885, '43.647540%'),
 Text(-0.5992044572731418, -0.03088718802350819, '14.344262%'),
 Text(0.2476577272604084, 0.48001604041680566, '14.344262%')])
```

1 # 9. Does willing ness among family memebers increases the chance of seeking more professi

2

3 plt.figure(figsize=(10,10))

4 plt.subplot(1,3,1)

5 plt.title("Very open family")

6 plt.pie(data4[data4.openess_of_family_friends=='Very open']['treatment_from_professional'])

7

8 plt.subplot(1,3,2)

9 plt.title("Somewhat Open family")

10 plt.pie(data4[data4.openess_of_family_friends=='Somewhat open']['treatment_from_profession

11

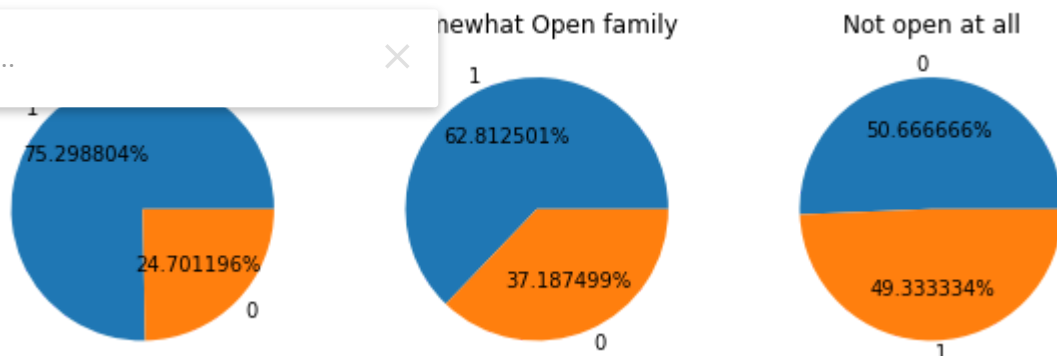
12 plt.subplot(1,3,3)

13 plt.title("Not open at all")

14 plt.pie(data4[data4.openess_of_family_friends=='Not open at all']['treatment_from_professi

```
([<matplotlib.patches.Wedge at 0x7fd81ca1e190>,
 <matplotlib.patches.Wedge at 0x7fd81ca1eb90>],
 [Text(-0.02303663990543217, 1.0997587522824575, '0'),
 Text(0.02303663990543179, -1.0997587522824575, '1')],
 [Text(-0.012565439948417547, 0.5998684103358859, '50.666666%'),
 Text(0.012565439948417339, -0.5998684103358859, '49.333334%')])
```

Saving...



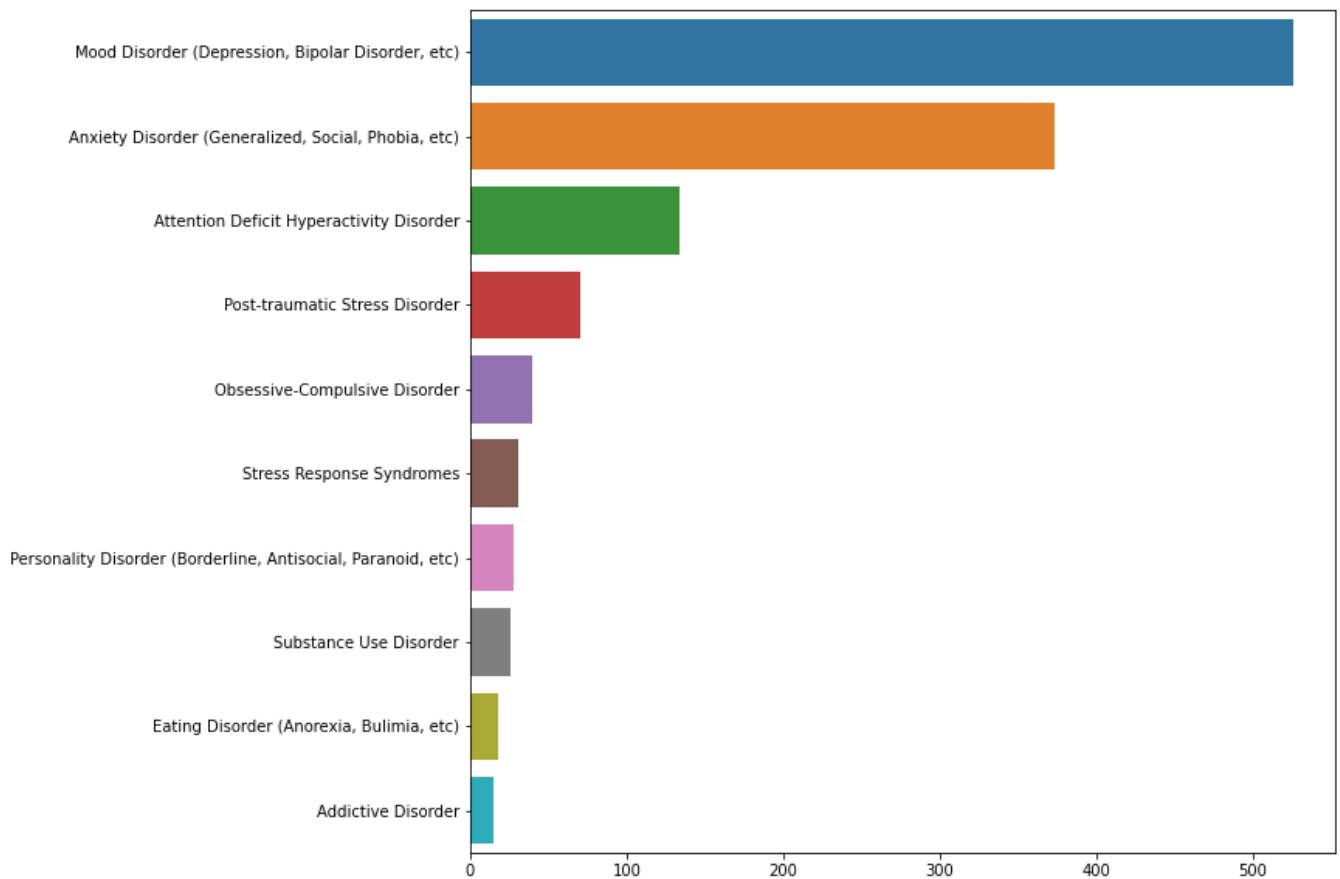
1 # 10. Which kind of disorder occur most?

2 disorder_type=pd.DataFrame(data4[data4.type_of_disorder.isnull() != True]['type_of_disorde

3 plt.figure(figsize=(10,10))

4 sns.barplot(x=disorder_type.value_counts()[0:10],y=disorder_type.value_counts().index[0:10

<matplotlib.axes._subplots.AxesSubplot at 0x7fd81d6d2810>



▼ 4. ML Models

Saving...

Target Variable Column: "diagnosed_mental_health_condition"

Aim : Here our main task is that knowing certain parameters of the respondent's background we have to predict if one will be diagnosed positive or negative.

▼ Stop Data Leakage:

- So, now comes the important part where we have to stop data leakage. To stop Data leakage we have to drop certain columns which we can't have while making prediction. Like Treatment from professional column we might not know when we are making prediction because if one

is diagnosed then one takes help from professional. These type of columns are like false target if we include it might train on that and then make prediction.

Stop Train and Test contamination :

- To stop this issue we have split the data then done all the preprocessing separately.

```
1 print(data4.shape)
```

```
(1433, 27)
```

```
1 data4.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1433 entries, 0 to 1432
Data columns (total 27 columns):
#   Column                                                                 Non-Null Count  Dtype
---  -
0   self_employed                                                         1433 non-null   object
1   no_of_employees                                                       1433 non-null   object
2   tech_company                                                          1433 non-null   object
3   mental_healthcare_coverage                                           1433 non-null   object
4   knowledge_about_mental_healthcare_options_workplace                 1433 non-null   object
5   employer_discussed_mental_health                                     1433 non-null   object
6   employer_offer_resources_to_learn_about_mental_health               1433 non-null   object
7   medical_leave_from_work                                              1433 non-null   object
8   comfortable_discussing_with_coworkers                               1433 non-null   object
9   employer_take_mental_health_seriously                               1433 non-null   object
10  openness_of_family_friends                                           1433 non-null   object
11  family_history_mental_illness                                         1433 non-null   object
12  mental_health_disorder_past                                           1433 non-null   object
13  currently_mental_health_disorder                                      1433 non-null   object
14  diagnosed_mental_health_condition                                    1433 non-null   object
15  mental_health_issue_interferes_work                                  1433 non-null   object
16  while_not_effective_treatment_interferes_work                       1433 non-null   object
17  age                                                                    1433 non-null   object
18  gender                                                                1433 non-null   object
19  country                                                                1433 non-null   object
20  country_work                                                          1433 non-null   object
21  work_remotely                                                         1433 non-null   object
22  tech_role                                                             1433 non-null   object
23  type_of_disorder                                                       711 non-null    object
24  US_state                                                              840 non-null    object
25  US_state_work                                                         851 non-null    object
dtypes: object(27)
memory usage: 302.4+ KB
```

▼ Data Preperation

```
1 data4.shape
```

```
(1433, 27)
```

```
1 # Here We Dropping unnecessary columns
```

```
2 y=data4.diagnosed_mental_health_condition
```

```
3 x=data4.drop(['diagnosed_mental_health_condition','treatment_from_professional','while_eff
```

```
1 print(x.shape)
```

```
2 print(y.shape)
```

```
(1433, 20)
```

```
(1433,)
```

```
1 # Splitting the data
```

```
2 x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.8,test_size=0.2,random_sta
```

```
3 print(x_train.shape)
```

```
4 print(x_test.shape)
```

```
5 print(y_train.shape)
```

```
6 print(y_test.shape)
```

```
(1146, 20)
```

```
(287, 20)
```

```
(1146,)
```

```
(287,)
```

```
1 cat_columns=['self_employed',
```

```
2             'no_of_employees',
```

```
3             'tech_company',
```

```
4             'mental_health_care_coverage',
```

```
5             'mental_healthcare_options_workplace',
```

```
6             'mental_health ',
```

```
7             'employer_offer_resources_to_learn_about_mental_health',
```

```
8             'medical_leave_from_work ',
```

```
9             'comfortable_discussing_with_coworkers',
```

```
10            'employer_take_mental_health_seriously',
```

```
11            'openess_of_family_friends',
```

```
12            'family_history_mental_illness',
```

```
13            'mental_health_disorder_past',
```

```
14            'currently_mental_health_disorder',
```

```
15            'age',
```

```
16            'gender',
```

```
17            'country',
```

```
18            'country work ',
```

```
19            'work_remotely',
```

```
20            'tech_role']
```

Saving...




```
1 print(data4['diagnosed_mental_health_condition'].unique())

['Yes' 'No']
```

```
1 for col in cat_columns:
2     print('The Unique value',col,'is')
3     print(data4[col].unique())
4     print()
```

```
The Unique value self_employed is
[0 1]
```

```
The Unique value no_of_employees is
['1-5' '6-25' '26-100' '100-500' '500-1000' '>1000']
```

```
The Unique value tech_company is
[1.0 0.0]
```

```
The Unique value mental_healthcare_coverage is
['Yes' 'No' "I don't know"]
```

```
The Unique value knowledge_about_mental_healthcare_options_workplace is
['Yes' 'No' 'I am not sure']
```

```
The Unique value employer_discussed_mental_health is
['No' 'Yes' "I don't know"]
```

```
The Unique value employer_offer_resources_to_learn_about_mental_health is
['No' "I don't know" 'Yes']
```

```
The Unique value medical_leave_from_work is
['Somewhat difficult' 'Very easy' "I don't know" 'Very difficult'
 'Somewhat easy' 'Neither easy nor difficult']
```

```
The Unique value comfortable_discussing_with_coworkers is
```

Saving...

```
mental_health_seriously is
["I don't know" 'Yes' 'No']
```

```
The Unique value openness_of_family_friends is
['Somewhat open' 'Very open' 'Somewhat not open' 'Neutral' "I don't know"
 'Not open at all']
```

```
The Unique value family_history_mental_illness is
['Yes' 'No' "I don't know"]
```

```
The Unique value mental_health_disorder_past is
['Yes' 'No' 'Maybe']
```

```
The Unique value currently_mental_health_disorder is
['No' 'Yes' 'Maybe']
```

The Unique value age is

```
[33.0 40.0 21.0 36.0 42.0 26.0 29.0 30.0 56.0 35.0 51.0 24.0 38.0 44.0
 27.0 55.0 22.0 25.0 28.0 23.0 32.0 31.0 43.0 37.0 39.0 45.0 46.0 20.0
 54.0 34.0 61.0 41.0 48.0 66.0 19.0 52.0 50.0 49.0 47.0 57.0 74.0 53.0
 58.0 70.0 59.0 62.0 63.0 65.0]
```

The Unique value gender is

```
['male' 'female' 'other']
```

The Unique value country is

```
['Canada' 'Netherlands' 'United Kingdom' 'Brazil'
 'United States of America' 'Denmark' 'Mexico' 'Australia' 'India' 'Iran'
 'Switzerland' 'Finland' 'Austria' 'Romania' 'Spain' 'Germany' 'Ireland'
 'Vietnam' 'South Africa' 'Slovakia' 'Norwav' 'France' 'Sweden']
```

```
1 from sklearn.preprocessing import LabelEncoder
2 import numpy as np
3
4
5 class LabelEncoderExt(object):
6     def __init__(self):
7         """
8         It differs from LabelEncoder by handling new classes and providing a value for it
9         Unknown will be added in fit and transform will take care of new item. It gives un
10        """
11        self.label_encoder = LabelEncoder()
12        # self.classes_ = self.label_encoder.classes_
13
14    def fit(self, data_list):
15        """
16        This will fit the encoder for all the unique values and introduce unknown value
17        :param data_list: A list of string
18        :return: self
19        """
20        self.label_encoder = self.label_encoder.fit(list(data_list) + ['Unknown'])
21        self.classes_ = self.label_encoder.classes_
```

Saving...



```
25    def transform(self, data_list):
26        """
27        This will transform the data_list to id list where the new values get assigned to
28        :param data_list:
29        :return:
30        """
31        new_data_list = list(data_list)
32        for unique_item in np.unique(data_list):
33            if unique_item not in self.label_encoder.classes_:
34                new_data_list = ['Unknown' if x==unique_item else x for x in new_data_list]
35
36        return self.label_encoder.transform(new_data_list)
```

```
1 label_encode=LabelEncoderExt()  
2  
3 label_x_train=x_train.copy()  
4 label_x_test=x_test.copy()  
5  
6 for col in cat_columns:  
7     label_x_train[col]=label_encode.fit(x_train[col])  
8     label_encode.classes_  
9     label_x_train[col]=label_encode.transform(x_train[col])  
10    label_x_test[col] = label_encode.transform(label_x_test[col])
```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:33: FutureWarning:
elementwise comparison failed; returning scalar instead, but in the future will perform



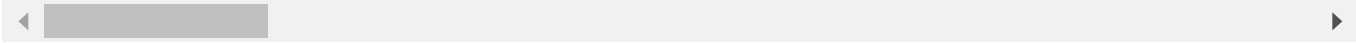
```
1 label_x_train
```

| | self_employed | no_of_employees | tech_company | mental_healthcare_coverage | knowled |
|------|---------------|-----------------|--------------|----------------------------|---------|
| 867 | 2 | 3 | 2 | | 1 |
| 608 | 2 | 2 | 2 | | 3 |
| 511 | 2 | 1 | 2 | | 3 |
| 1214 | 2 | 2 | 2 | | 3 |
| 1244 | 2 | 2 | 2 | | 3 |
| ... | ... | ... | ... | | ... |
| 763 | 2 | 2 | 2 | | 0 |
| 835 | 2 | 3 | 2 | | 3 |
| | | 2 | 2 | | 3 |
| 559 | 2 | 2 | 2 | | 1 |
| 684 | 2 | 2 | 2 | | 0 |

Saving...

✕

1146 rows × 20 columns



```
1 label_x_test
```

| | self_employed | no_of_employees | tech_company | mental_healthcare_coverage | knowled |
|-------------|---------------|-----------------|--------------|----------------------------|---------|
| 1059 | 2 | 5 | 2 | | 3 |
| 411 | 2 | 1 | 2 | | 3 |
| 342 | 2 | 1 | 2 | | 3 |
| 1295 | 2 | 2 | 2 | | 3 |
| 483 | 2 | 1 | 2 | | 3 |
| ... | ... | ... | ... | | ... |
| 1045 | 2 | 5 | 2 | | 3 |
| 1309 | 2 | 2 | 2 | | 3 |
| 520 | 2 | 2 | 2 | | 0 |
| 993 | 2 | 5 | 2 | | 3 |
| 333 | 2 | 1 | 2 | | 3 |

287 rows × 20 columns



```
1 df = pd.DataFrame(label_x_test)
2
3 for col in cat_columns:
4     print('The Unique value',col,'is')
5     print(df[col].unique())
6     #print(type(df["Subjects"].unique()))
7
8 type(label_x_test)
```

The Unique value self_employed is
[2]

Saving...



is is

The Unique value tech_company is
[2]

The Unique value mental_healthcare_coverage is
[3 0 1]

The Unique value knowledge_about_mental_healthcare_options_workplace is
[0 1 3]

The Unique value employer_discussed_mental_health is
[1 3 0]

The Unique value employer_offer_resources_to_learn_about_mental_health is
[0 1 3]

The Unique value medical_leave_from_work is
[5 1 0 3 6 2]

The Unique value comfortable_discussing_with_coworkers is
[0 1 3]

The Unique value employer_take_mental_health_seriously is

```

[1 0 3]
The Unique value openness_of_family_friends is
[3 6 1 4 0 2]
The Unique value family_history_mental_illness is
[3 1 0]
The Unique value mental_health_disorder_past is
[1 3 0]
The Unique value currently_mental_health_disorder is
[3 0 1]
The Unique value age is
[48]
The Unique value gender is
[2 1 3]
The Unique value country is
[49 48 33 31 44 21 45 7 40 10 15 50 2 25 11 28 37 3 18 27 46 19 20 30
32]
The Unique value country_work is
[48 47 32 43 20 44 6 39 9 49 14 1 24 10 27 36 2 17 26 45 18 19 31]
The Unique value work_remotely is
[0 2 1]
The Unique value tech_role is
[2]
pandas.core.frame.DataFrame

```

```

1 # For Y label Encode
2 label_encode_1=LabelEncoder()
3 label_y_train_1=label_encode_1.fit_transform(y_train)
4 label_y_test_1=label_encode_1.transform(y_test)

```

```

1 st=pd.DataFrame(label_y_train_1)
2 print(st)

```

```

      0
0      0
1      0
2      0

```

Saving...

```

1141  0
1142  1
1143  0
1144  1
1145  0

```

```
[1146 rows x 1 columns]
```

```

1 st=pd.DataFrame(label_y_test_1)
2 print(st)

```

```

      0
0      1
1      1

```

```

2      1
3      0
4      1
..    ..
282    0
283    1
284    1
285    1
286    1

```

```
[287 rows x 1 columns]
```

1. Logistic Regression

```

1 import sklearn
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score
4 logistic=LogisticRegression(C=1,penalty='l1',solver='liblinear',random_state=0)
5
6 logistic.fit(label_x_train,label_y_train_1)
7 preds3=logistic.predict(label_x_test)
8 accuracy_score(label_y_test_1,preds3)

```

```
0.89198606271777
```

```

1 from sklearn.metrics import confusion_matrix
2 from sklearn.metrics import accuracy_score
3 from sklearn.metrics import classification_report
4 from sklearn.metrics import roc_auc_score
5 from sklearn.metrics import log_loss
6
7 results = confusion_matrix(label_y_test_1,preds3)

```

Saving...



```

... accuracy_score(label_y_test_1,preds3))
11 print ('Classification Report : ')
12 print (classification_report(label_y_test_1,preds3))
13 print('AUC-ROC:',roc_auc_score(label_y_test_1,preds3))
14 print('LOGLOSS Value is',log_loss(label_y_test_1,preds3))

```

Confusion Matrix :

```
[[127  13]
 [ 18 129]]
```

Accuracy Score is 0.89198606271777

Classification Report :

| | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| 0 | 0.88 | 0.91 | 0.89 | 140 |
| 1 | 0.91 | 0.88 | 0.89 | 147 |

| | | | | |
|--------------|------|------|------|-----|
| accuracy | | | 0.89 | 287 |
| macro avg | 0.89 | 0.89 | 0.89 | 287 |
| weighted avg | 0.89 | 0.89 | 0.89 | 287 |

AUC-ROC: 0.8923469387755102

LOGLOSS Value is 3.730705446023776

2. Decision Tree

```

1 from sklearn.tree import DecisionTreeClassifier
2 clf = DecisionTreeClassifier()
3 clf = clf.fit(label_x_train,label_y_train_1)
4 y_pred = clf.predict(label_x_test)
5 accuracy_score(label_y_test_1,y_pred)

```

0.7909407665505227

```

1 from sklearn.metrics import confusion_matrix
2 from sklearn.metrics import accuracy_score
3 from sklearn.metrics import classification_report
4 from sklearn.metrics import roc_auc_score
5 from sklearn.metrics import log_loss
6
7 results = confusion_matrix(label_y_test_1,y_pred)
8 print ('Confusion Matrix :')
9 print(results)
10 print ('Accuracy Score is',accuracy_score(label_y_test_1,y_pred))
11 print ('Classification Report : ')
12 print (classification_report(label_y_test_1,y_pred))
13 print('AUC-ROC:',roc_auc_score(label_y_test_1,y_pred))
14 print('LOGLOSS Value is',log_loss(label_y_test_1,y_pred))

```

Saving...

```

[[ 31 116]]
Accuracy Score is 0.7909407665505227
Classification Report :

```

| | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| 0 | 0.78 | 0.79 | 0.79 | 140 |
| 1 | 0.80 | 0.79 | 0.79 | 147 |

| | | | | |
|--------------|------|------|------|-----|
| accuracy | | | 0.79 | 287 |
| macro avg | 0.79 | 0.79 | 0.79 | 287 |
| weighted avg | 0.79 | 0.79 | 0.79 | 287 |

AUC-ROC: 0.7909863945578232

LOGLOSS Value is 7.220730912962079

3. Random Forest

```

1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score
3 model=RandomForestClassifier(n_estimators=1000, max_depth=10, random_state=0)
4 model.fit(label_x_train,label_y_train_1)
5 preds=model.predict(label_x_test)
6 accuracy_score(label_y_test_1,preds)

```

0.9337979094076655

```

1 from sklearn.metrics import confusion_matrix
2 from sklearn.metrics import accuracy_score
3 from sklearn.metrics import classification_report
4 from sklearn.metrics import roc_auc_score
5 from sklearn.metrics import log_loss
6
7 results = confusion_matrix(label_y_test_1,preds)
8 print ('Confusion Matrix :')
9 print(results)
10 print ('Accuracy Score is',accuracy_score(label_y_test_1,preds))
11 print ('Classification Report : ')
12 print (classification_report(label_y_test_1,preds))
13 print('AUC-ROC:',roc_auc_score(label_y_test_1,preds))
14 print('LOGLOSS Value is',log_loss(label_y_test_1,preds))

```

Confusion Matrix :

[[126 14]

[5 142]]

Accuracy Score is 0.9337979094076655

Classification Report :

| | 1 | 0 | precision | recall | f1-score | support |
|--------------|------|------|-----------|--------|----------|---------|
| 0 | 126 | 14 | 0.90 | 0.90 | 0.93 | 140 |
| 1 | 5 | 142 | 0.91 | 0.97 | 0.94 | 147 |
| accuracy | | | | | 0.93 | 287 |
| macro avg | 0.94 | 0.93 | 0.93 | | 0.93 | 287 |
| weighted avg | 0.94 | 0.93 | 0.93 | | 0.93 | 287 |

AUC-ROC: 0.9329931972789115

LOGLOSS Value is 2.2865782085969544

4. KNN


```

1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 scaler.fit(label_x_train)
4 label_x_train = scaler.transform(label_x_train)
5 label_x_test = scaler.transform(label_x_test)
6 from sklearn.neighbors import KNeighborsClassifier
7 classifier = KNeighborsClassifier(n_neighbors=8)
8 classifier.fit(label_x_train, label_y_train_1)

```

```

KNeighborsClassifier(n_neighbors=8)

```

```

1 y_pred1 = classifier.predict(label_x_test)

```

```

1 from sklearn.metrics import confusion_matrix
2 from sklearn.metrics import accuracy_score
3 from sklearn.metrics import classification_report
4 from sklearn.metrics import roc_auc_score
5 from sklearn.metrics import log_loss
6
7 results = confusion_matrix(label_y_test_1,y_pred1)
8 print ('Confusion Matrix :')
9 print(results)
10 print ('Accuracy Score is',accuracy_score(label_y_test_1,y_pred1))
11 print ('Classification Report : ')
12 print (classification_report(label_y_test_1,y_pred1))
13 print('AUC-ROC:',roc_auc_score(label_y_test_1,y_pred1))
14 print('LOGLOSS Value is',log_loss(label_y_test_1,y_pred1))

```

Confusion Matrix :

```

[[129  11]
 [ 30 117]]

```

Accuracy Score is 0.8571428571428571

Classification Report :

| | precision | recall | f1-score | support |
|--|-----------|--------|----------|---------|
|--|-----------|--------|----------|---------|

Saving...

| | | |
|----|------|-----|
| 92 | 0.86 | 140 |
| 80 | 0.85 | 147 |

| | | | |
|--------------|------|------|-----|
| accuracy | | 0.86 | 287 |
| macro avg | 0.86 | 0.86 | 287 |
| weighted avg | 0.86 | 0.86 | 287 |

AUC-ROC: 0.8586734693877552

LOGLOSS Value is 4.9341415601500715

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 4:51 PM



Saving...