

Configuration Manual

MSc Research Project MSc In Data Analytics

Sonal Suryakant Puradkar Student ID: X22130250

School of Computing National College of Ireland

Supervisor:

Arjun Chikkankod

National College of Ireland



MSc Project Submission Sheet

School of Computing

Student Name:	Sonal Suryakant Puradkar					
Student ID:	X22130250					
Programme:	MSc In Data Analytics MSc Research Project	Year:	2023-2024.			
Module:						
Lecturer: Submission Due	Arjun Chikkankod					
Date:	14 Dec 2023					

Project Title:Forecasting Global Mental Health Disorders : A Machine Learning
Approach using Socioeconomic Indicators

Word Count:1137

Page Count:19

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

 Signature:
 Sonal Suryakant

 Puradkar.....

 Date:
 14 Dec

2023.....

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple	
copies)	
Attach a Moodle submission receipt of the online project	
submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project, both	
for your own reference and in case a project is lost or mislaid. It is not	
sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only

Office use Offiy	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Sonal Suryakant Puradkar Student ID:X22130250

1 Introduction

The objective of the project is to find out if the countries where mental health disorders namely Anxiety Disorders, Bipolar Disorders, Eating Disorders, Schizophrenia and Depressive disorders will be highest can be predicted based on socio economic indicators and to see the forecasted values for 3 years. For this study we are taking country wise year wise data of the below socio economic indicators for the years 1960-2022:

Indicator name	Indicator Code
Adjusted net national income per capita (current US\$)	NY.ADJ.NNTY.PC.CD
Consumer price index (2010 = 100)	FP.CPI.TOTL
Inflation, consumer prices (annual %)	FP.CPI.TOTL.ZG
Employment in industry (% of total employment) (modeled ILO estimate)	SL.IND.EMPL.ZS
Employment in services (% of total employment) (modeled ILO estimate)	SL.SRV.EMPL.ZS
Self-employed, total (% of total employment) (modeled ILO estimate)	SL.EMP.SELF.ZS
Proportion of people living below 50 percent of median income (%)	SI.DST.50MD
Unemployment with advanced education (% of total labor force with	
advanced education)	SL.UEM.ADVN.ZS
Unemployment, total (% of total labor force) (modeled ILO estimate)	SL.UEM.TOTL.ZS
New businesses registered (number)	IC.BUS.NREG
Multidimensional poverty index (scale 0-1)	SI.POV.MDIM.XQ

For mental health disorders we are taking country wise year wise Disability Adjusted Life Years (DALYs) data from 1990 -2019.

To further find out the workplace mental health affecting factors this study creates visualizations for the factors affecting workplace mental health.

Python is used for data pre-processing and data modelling and Power BI is used for visualization in this study.

2 Data Gathering

Socio economic indicators data has been downloaded from https://databank.worldbank.org/source/world-development-indicators and mental health DALY data has been downloaded from https://ourworldindata.org/mental-health . Workplace mental health data has been collected from Kaggle survey on

https://www.kaggle.com/code/aditimulye/mental-health-at-workplace/input. All the data has been gathered and kept into the root folder of the project.

3 Data Pre Processing

The data pre processing is done in python. It starts with reading data files from downloaded CSV files. In section 1 of the notebook CrispDM_Mental_Health_forecasting.ipynb all the libraries are imported and files are read into dataframes as shown below

	•													
1.	1 Importii	ng Libr	aries											
1 1 1 1 1 1 W 1 W	1 import requests 2 import pprint 3 import pandes as pd 4 import pandes as osl 6 import numpy as np 7 import numpy as np 7 import matplotilo.pyrot as plt													
C: fo	\ProgramDa r this ver warnings.wa	ta∖Anacc sion of arn(f"A	nda3\lib\sit SciPy (detec NumPy versio	e-packages\scipy\ ted version 1.23. n >={np_minversio	init 5 on} and <	.py:146: {np_maxv	UserWarn ersion}"	ning: A≬	lumPy ver	sion >=1	.16	.5 and <1.2	3.0 is requ	uired
1.:	2 Importii	ng Data	Files and o	checking first 1	0 rows									
3	1 metadata 2 Indicato 3 Mental_d	_df = po or_data_o lata_df	d.read_csv('I df = pd.read_ = pd.read_csv	ndicators_Metada csv('Indicators_ ('Mental_Health_	ta.csv') Data.csv' Data.csv'	3								
1 Indicator_data_df.head(10)														
1	1 Indicato	r_data_d	df.head(10)											
1	Country Name	Country Code	Series Name	Series Code	1960 [YR1960]	1961 [YR1961]	1962 [YR1962]	1963 [YR1963]	1964 [YR1964]	1965 [YR1965]		2013 [YR2013]	2014 [YR2014]	C
0	Afghanistan	Country Code	Adjusted net Adjusted net national income per capita (courre	Series Code	1960 [YR1960] 	1961 [YR1961]	1962 [YR1962]	1963 [YR1963]	1964 [YR1964]	1965 [YR1965]		2013 [YR2013] 586.6136889	2014 [YR2014] 575.6807964	6 26.1
1 0 1	Country Name Afghanistan	Country Code AFG AFG	Adjusted net national income per capita (curre Consumer price index (2010 = 100)	Series Code NY:ADJ.NNTY:PC.CD FP:CPI.TOTL	1960 [YR1960] 	1961 [YR1961] 	1962 [YR1962] 	1963 [YR1963] 	1964 [YR1964] 	1965 [YR1965] 		2013 [VR2013] 588.8138889 127.7952231	2014 [YR2014] 575.6807984 133.7683667	528.1 132.2
0 1 2	I Indicato Country Name Afghanistan Afghanistan	Country Code AFG AFG AFG	Adjusted net national income per capita (curre Consumer price index (2010 = 100) Inflation, consumer prices (annual %)	Series Code NY.ADJ.NNTY.PC.CD FP.CPI.TOTL FP.CPI.TOTL 2G	1960 [YR1960] 	1961 [YR1961] 	1962 [YR1962] 	1963 [YR1963] 	1964 [YR1964] 	1965 [YR1965] 		2013 [VR2013] 588.6138880 127.7952231 7.385771784	2014 [YR2014] 575.6807064 133.7683667 4.673996035	528.1 132.3 -0.00
1 2 3	I Indicato Country Name Afghanistan Afghanistan Afghanistan	Country Code AFG AFG AFG AFG	Adjusted net Adjusted net ational income per capita (corre Consume price index (collo = 100) inflation, consume prices (annea biologymett) Employmet)	Series Code NY.ADJ.NNTY.PC.CD FR.CPI.TOTL FR.CPI.TOTL ZG SL.IND.EMPL.ZS	1960 [YR1960] 	1961 [YR1961] 	1962 [YR1962] 	1963 [YR1963] 	1964 [YR1964] 	1965 [YR1965] 		2013 [YR2013] 588.8138889 127.7952231 7.385771784 19.53853	2014 [YR2014] 575.8807084 133.7683667 4.673996035 22.04671	526.1 132.1 -0.00
1 0 1 2 3 4	Indicato Country Name Afghanistan Afghanistan Afghanistan Afghanistan	Country Code AFG AFG AFG AFG AFG	Series Name Adjusted net national income per capita (curre Consumer proce index (2010 - 100) Inflation, industry (% of employment) Employment in services (% of employment)	Series Code NYADJ.NNTY,PC.CD FP.CPI.TOTL PP.CPI.TOTL2G SL.IND.EMPL2S SL.SRV.EMPL2S	[YR1960] [YR1960] 	1961 [YR1961] 	1962 [YR1962] 	1963 [YR1963] 	1964 [YR1964] 	1965 [YR1965] 	 	2013 [YR2013] 586.6136680 127.7052231 7.385771784 19.53653 32.82066	2014 [YR2014] 575.6807964 133.7683667 4.673990035 22.04671 33.15474	528.4 132.4 -0.68 2 3
0 1 2 3 4 5	I Indicato Country Name Afghanistan Afghanistan Afghanistan Afghanistan Afghanistan	Country Code AFG AFG AFG AFG AFG AFG AFG	Series Name Adjusted net national income per cacita consumer proce index (2010 = 100) Inflation, consumer prices (anny) Employment) Employment) Employment) Employment) Employment) Employment) Beirtemployed. total (% of total employmer)	Series Code NYADJ.NNTY.PC.OD FP.CPI.TOTL FP.CPI.TOTLZG SL.IND.EMPLZS SL.SRV.EMPLZS SL.EMP.SELF.ZS	1960 [YR1960] 	1961 [YR1961] 	1962 [YR1962]	1963 [YR1963] 	1964 [YR1964] 	1965 [YR1965] 		(vR2013) 598.6136889 127.7082231 7.385771784 19.53853 32.82088 87.09895	(YR2014) 575.8807084 133.7683667 4.673996035 22.04671 33.15474 86.01729	526.1 132.1 -0.06 2 3 8

In [4]: 1 metadata_df.head(10)

	Series Name	Series Code
0	Adjusted net national income per capita (curre	NY.ADJ.NNTY.PC.CD
1	Consumer price index (2010 = 100)	FP.CPI.TOTL
2	Inflation, consumer prices (annual %)	FP.CPI.TOTL.ZG
3	Employment in industry (% of total employment)	SL.IND.EMPL.ZS
4	Employment in services (% of total employment)	SL.SRV.EMPL.ZS
5	Self-employed, total (% of total employment) (SL.EMP.SELF.ZS
6	Proportion of people living below 50 percent o	SI.DST.50MD
7	Unemployment with advanced education (% of tot	SL.UEM.ADVN.ZS
8	Unemployment, total (% of total labor force) (SL.UEM.TOTL.ZS
9	New businesses registered (number)	IC.BUS.NREG

In [5]: 1 Mental_data_df.head(10)

Out[5]:

:

Out[4]

	Entity	Code	Year	DALYs from depressive disorders per 100,000 people in, both sexes aged age-standardized	DALYs from schizophrenia per 100,000 people in, both sexes aged age-standardized	DALYs from bipolar disorder per 100,000 people in, both sexes aged age-standardized	DALYs from eating disorders per 100,000 people in, both sexes aged age-standardized	DALYs from anxiety disorders per 100,000 people in, both sexes aged age-standardized
0	Afghanistan	AFG	1990	895.22565	138.24825	147.64412	26.471115	440.33000
1	Afghanistan	AFG	1991	893.88434	137.76122	147.56696	25.548681	439.47202
2	Afghanistan	AFG	1992	892.34973	137.08030	147.13088	24.637949	437.60718
3	Afghanistan	AFG	1993	891.51587	136.48602	146.78812	23.863169	436.69104
4	Afghanistan	AFG	1994	891.39160	136.18323	146.58481	23.189074	436.76800
5	Afghanistan	AFG	1995	891.21344	135.65398	146.63217	22.503244	436.69098
6	Afghanistan	AFG	1996	891.92096	135.16324	146.67957	21.827528	438.52713
7	Afghanistan	AFG	1997	893.02045	134.20360	146.54947	21.124723	436.42557
8	Afghanistan	AFG	1998	894.73680	133.56137	146.69704	20.412657	436.77020
9	Afghanistan	AFG	1999	896.13806	132.82457	146.76817	19.931143	437.12753

Data Pre Processing starts in section 2 in the notebook. The null values are represented as "..." in the downloaded CSVs, these values are replaced by nulls in the section 2.1. Also the data has years in columns and indicators in rows, so the data is unpivoted about years in section 2.2. Following this the year column is formatted to show numerical values in section 2.3 :



2.2 Unpivot Indicators data about years and Pivot about values

In [9]:	1	#Getting List of Columns
	2	
	З	Indicator_data_df_columns = Indicator_data_df.columns
	4	<pre>Indicator_data_columns_list = Indicator_data_df_columns.tolist()</pre>
	5	pprint.pprint(Indicator_data_columns_list)
	6	
	7	

2.3 Format year column , convert value to numeric column



The column Year and Country code are renamed in section 2.4 to make it consistent in mental health disorders dataset and indicators dataset. This is followed up by joining the two dataframes in section 2.5

```
2.4 Renaming columnnames to make it consistent across all dataframes
             1 Indicators_df = Indicators_df.rename(columns={'year': 'Year', 'Country Code' : 'Code'})
In [13]:
             3 Indicators_df['value'].replace('', np.nan, inplace=True)
3 Indicators_df.head()
Out[13]:
                                             Series Code Year value
                Country Name Code
            0 Afghanistan AFG NY.ADJ.NNTY.PC.CD 1980 NaN
            1
                  Afghanistan AFG
                                           EPICPLITOTI 1980 NaN
                                      FP.CPI.TOTL.ZG 1980 NaN
            2
                 Afghanistan AFG
                  Afghanistan AFG
                                        SLIND.EMPL.ZS 1960 NaN
            3
            4 Afghanistan AFG SL.SRV.EMPL.ZS 1980 NaN
In [14]: 1 Indicators_df['value'] = Indicators_df['value'].apply(lambda x: float(x) if x else float('nan'))
2 column_types = Indicators_df.dtypes
             4 print(column_types)
             5
            Country Name
                                object
                                object
object
            Code
            Series Code
            Year
                                object
            value
                               float64
            dtype: object
            2.5 Join Indicators and mental health dataframes based on country code and year
             1 Indicators_df['Year'] = Indicators_df['Year'].astype(str)
2 Mental_data_df['Year'] = Mental_data_df['Year'].astype(str)
In [15]:
             3 data_df = pd.merge(Indicators_df, Mental_data_df, on=['Code', 'Year'], how='inner')
In [16]: 1 data_df.head()
Out[16]:
                                                                                     DALYs from
                                                                                                                    DALYs from
bipolar disorder
per 100,000
people in, both
sexes aged age-
standardized
                                                                                                                                                              DALYs from
                                                                                   DALYs from
depressive DALYs from
disorders per
schizophrenia per
100,000 people
in, both sexes
aged age-
tandrecified
                                                                                                                                            DALYs from
                                                                                                                                                           anxiety
disorders per
100,000 people
in, both sexes
                                                                                                                                      DALYs from
eating disorders
per 100,000
people in, both
sexes aged age-
standardized
                  Country Code
                                          Series Code Year value
                                                                         Entity
                                                                                    aged age-
standardized
                                                                                                                                                             aged age-
standardized
            0 Afghanistan AFG NYADJ.NNTY.PC.CD 1990 NaN Afghanistan 895.22565 138.24825 147.64412 26.471115
                                                                                                                                                                  440.33
             1 Afohanistan AFG
                                      FP.CPI.TOTL 1990 NaN Afghanistan
                                                                                       895.22565
                                                                                                          138.24825
                                                                                                                            147.64412
                                                                                                                                              26.471115
                                                                                                                                                                   440.33
            2 Afghanistan AFG
                                     FP.CPI.TOTL.ZG 1990 NaN Afghanistan
                                                                                      895.22565
                                                                                                         138.24825
                                                                                                                           147.64412
                                                                                                                                             26.471115
                                                                                                                                                                   440.33
             3 Afghanistan AFG
                                     SL.IND.EMPL.ZS 1990 NaN Afghanistan
                                                                                       895.22565
                                                                                                          138.24825
                                                                                                                            147.64412
                                                                                                                                              26.471115
                                                                                                                                                                   440.33
            4 Afghanistan AFG SL.SRV.EMPL.ZS 1990 NaN Afghanistan
                                                                                       895.22565
                                                                                                          138.24825
                                                                                                                            147.64412
                                                                                                                                              26.471115
                                                                                                                                                                   440.33
In [17]: 1 data_df.info()
```

In section 2.6 the disorders DALYs are renamed to a shorter name for ease of understanding :

2.6 Renaming the disorder features to shorter names In [18]: 1 data_df.columns Out[18]: Index(['Country Name', 'Code', 'Series Code', 'Year', 'value', 'Entity', 'DALYS from depressive disorders per 100,000 people in, both sexes aged age-standardized', 'DALYS from schlzophrenia per 100,000 people in, both sexes aged age-standardized', 'DALYS from bipolar disorder per 100,000 people in, both sexes aged age-standardized', 'DALYS from eating disorders per 100,000 people in, both sexes aged age-standardized', 'DALYS from eating disorders per 100,000 people in, both sexes aged age-standardized', 'DALYS from exitety disorders per 100,000 people in, both sexes aged age-standardized', 'dtype='object') })
data_df = data_df.drop(columns=['Entity']) In [20]: 1 data_df.head()
2 #data_df.to_csv('data.csv') out[20]:
 value
 Depressive_Disorders
 Schize

 NaN
 895.22585
 13
 Series Code Year Country Name Code nia Bipolar_Disorders Eating_Disorders Diso Afghanistan AFG NY.ADJ.NNTY.PC.CD 1990 138.24825 147.64412 26.471115 440.33 0 Afghanistan AFG FP.CPI.TOTL 1990 NaN 895.22565 138.24825 147.64412 26.471115 440.33 2 Afghanistan AFG FP.CPI.TOTL.ZG 1990 NaN 895.22565 138.24825 147.64412 26.471115 440.33 AFG SL.IND.EMPL.ZS 1990 147.64412 440.33 Afghanistan NaN 895.22565 138.24825 26.471115 4 Afghanistan AFG SL.SRV.EMPL.ZS 1990 NaN 895.22505 138.24825 147.64412 26.471115 440.33

In section 2.7 the data is unpivoted about indicators so they show up in row and can be used in prediction. These indicator codes in the column series code are then renamed for ease of understanding

	2.7	Piv	oting ab	out i	ndic	ators and renam	ing indicate	or columns					
In [21]:	1	piv ∢	voted_df =	data	_df.p	ivot(index=['Count	ry Name','Co	de','Year','Depi	ressive_Disorde	rs','Schizophre	nia','Bipol	ar_Disorders'	, 'E ▶
In [22]:	1 pivoted_df.reset_index(inplace=True)												
In [23]:	1	piv	voted_df.c	olumn	s								
Out[23]:	<pre>Index(['Country Name', 'Code', 'Year', 'Depressive_Disorders', 'Schizophrenia',</pre>												
In [24]:	1	piv	voted_df.h	nead()									
Out[24]:	Seri Co	ies de	Country Name	Code	Year	Depressive_Disorders	Schizophrenia	Bipolar_Disorders	Eating_Disorders	Anxiety_Disorders	FP.CPI.TOTL	FP.CPI.TOTL.ZG	IC.
		0	Afghanistan	AFG	1990	895.22565	138.24825	147.64412	28.471115	440.33000	NaN	NaN	I T
		1	Afghanistan	AFG	1991	893.88434	137.76122	147.56696	25.548681	439.47202	NaN	NaN	J
		2	Afghanistan	AFG	1992	892.34973	137.08030	147.13086	24.637949	437.60718	NaN	NaN	1
		3	Afghanistan	AFG	1993	891.51587	136.48602	146.78812	23.863169	436.69104	NaN	NaN	1
		4	Afghanistan	AFG	1994	891.39160	136.18323	146.58481	23.189074	436.76800	NaN	NaN	(
	4												•
In [25]:	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18	#re #a #b #c #d #e #f #f #i #i #k	enaming cc = FP.CPI. = FC.RUS. = SI.UEM. = SL.UEM. = SL.UEM. = SL.UEM. = SL.UEM. = SL.UEM. = SL.VEM. = SL.VEM. = SL.SRV. = SL.IND. voted_df =	TOTL TOTL. NREG SOMD ADVN. SELF. NTY. MDIM. EMPL. EMPL.	indic ZG ZS ZS ZS ZS ZS ZS ZS ted_d	ator columns f.rename(columns={	'FP.CPI.TOTL	': 'A', 'FP.CPI.	.TOTL.ZG' : 'B'	, 'IC.BUS.NREG	': 'C', 'S	I.DST.50MD' :	: 'D
	15	4											•

There are missing values for indicators and indicators with more than 60 percent data missing are dropped from this analysis. Missing values are handled in section 2.8.

2.8 Handling Missing values

	2.8.1	Checki	ng dis	tributi	on of missing/null va	lues and visu	alizing them						
In [26]:	1 2 3 4 5 6 7 8	# Count blank_c blank_ #bLank_	counts	bLank 5 = pi 5.head ts.to_	values in all indi voted_df.groupby([() csv('blanks.csv')	cator column	ns ne', 'Code','Yea	n', 'Depressive	e_Disorders','So	hizophr:	enia',	'Bipo	blar_Disorders'
Out[26]:									Series Code	Country Name	Code	Year	Depressive_Disorc
		Name	Code	Year	Depressive_Disorders	Schizophrenia	Bipolar_Disorders	Eating_Disorders	Anxiety_Disorders				
	Afgh	anistan	AFG	1990	895.22565	138.24825	147.64412	26.471115	440.33000	0	0	0	
				1991	893.88434	137,76122	147.56696	25.548681	439.47202	0	0	0	
				1992	892.34973	137.08030	147.13086	24.637949	437.60718	0	0	0	
				1993	891.51587	136.48602	146.78812	23.863169	436.69104	0	0	0	
				1994	891.39160	136.18323	146.58481	23.189074	436.76800	0	0	0	
	4												F.
In [27]:	1 2 3 4 5 6	<pre># Count cols = blanks_ print()</pre>	['A', agg_colanks	total 'B',' If = b Lagg_	blank values in ea C','D','E','F','G' lank_counts[cols]. df)	sum()	,'К']						
	Seri A B C D E F G H I J K dtvp	es code 984 1044 4214 4293 4141 722 1036 5917 722 e: int6	- -										

4 5 print (percentage_distribution)

Series Code

6

- A 16.400000 17.400000 в
- č
- D 71.550000 69.016667 Е
- 12.033333
- G 12.033333 17.266667
- н
- I 98.616667 Э
- 12.033333 12.033333 к

dtype: float64

[29]: 1 series = pd.Series(blanks_agg_df)
2 # Calculating percentage out of 6000 which is total number of records
3 percentages = series / 6000 * 100 4
5 # PLot the percentages
6 plt.bar(percentages.index, percentages.values, color='skyblue')
7 plt.xlabel('categories')
8 plt.ylabel('Percentage of 6000')
9 plt.title('Values as Percentage of 6000')
10 plt.ylim(0, 100) # Set the y-axis Limit to 0-100 for percentage representation
11 11 12 # Show the plot 13 plt.show() Values as Percentage of 6000



2.8.2 Dropping indicators where more than 60 percent data is null



1
2 pivoted_df = pivoted_df[['Country Name', 'Code','Year', 'Depressive_Disorders','Schizophrenia','Bipolar_Disorders','Eating_D 3 4 ⊧

For the columns remaining in the analysis, missing values are replaced by mean values in section 2.8.2.

	2.8.	3 Replacing missing values with mean()
In [32]:	1 2	<pre>numerical_columns = ['A', 'B', 'F', 'G', 'H', 'J', 'K'] pivoted_df[numerical_columns] = pivoted_df[numerical_columns].fillna(pivoted_df[numerical_columns].mean())</pre>
		~ ^ ~

This is followed up with exploratory data anlaysis using seaborn and feature scaling. In EDA distribution of all the columns is explored and correlation matrix is prepared for indicators and mental health disorders to find out any association.



In [40]: 1 #correlation matrix

- correlation_matrix = target_v_df.corr()
- plt.figure(figsize=(10, 8))
 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm') 4



plt.title('Correlation Heatmap')





Depressive DisordersSchizophrenia Bipolar Disorders Eating Disorders Anxiety Disorders Series Code



It is found that the disorders bipolar disorders and anxiety disorders are correlated according to available data, hence bipolar disorder are dropped. Similarly indicators J , K and F are correlated so the indicators J and K are dropped.



```
In [44]: 1
2 # dropping J and K as F , J and K are correlated and creating target dataframes for all the disorders
3
4
5 df_ad = df[['Country Name','Code','Year','A','B','F','G','H','Anxiety_Disorders']]
6
6 df_dd = df[['Country Name','Code','Year','A','B','F','G','H','Depressive_Disorders']]
8 df_s = df[['Country Name','Code','Year','A','B','F','G','H','Schizophrenia']]
9 df_ed = df[['Country Name','Code','Year','A','B','F','G','H','Eating_Disorders']]
```

This is followed up by creating separate dataframes for model implementation for all the disorders as seen in the picture above.

4 Modelling Evaluation and Visualizations

All the disorders are modelled using naiive forecasting and exponential smoothing model and evaluated against mean squared error, mean absolute error, root mean squared error, R squared and Adjusted R squared values. Following this forecasted values are generated in separate CSVs to be used for data visualizations using power bi. Section 3.1 has modelling and forecasting for anxiety disorders.

3.1 Modelling for Anxiety Disorders

```
In [53]: 1 data = df_ad
In [54]: 1 #renaming anxiety disorder to t and exporting to csv
            2 dataad = data.rename(columns={'Anxiety_Disorders': 'T'})
           4 dataad.to csv('addata.csv')
In [55]: 1 dataad.head()
Out[55]:
           Series Code Country Name Code Year
                                                     Α
                                                                 в
                                                                             F
                                                                                         G
                                                                                                     н
                                                                                                              т
                 0 Afghanistan AFG 1990 6.410890e-16 1.504054e-16 2.443852e-15 1.882180e-15 7.886385e-17 440.33000
                      Afghanistan AFG 1991 6.410690e-16 1.504054e-16 1.854854e+00 -8.010559e-03 7.866365e-17 439.47202
                   1
               2 Afghanistan AFG 1992 6.410890e-18 1.504054e-16 1.853048e+00 2.012347e-04 7.886385e-17 437.80718
                  3 Afghanistan AFG 1993 6.410890e-16 1.504054e-16 1.850257e+00 -7.681121e-03 7.886385e-17 436.69104
                  4 Afghanistan AFG 1994 6.410690e-16 1.504054e-16 1.848748e+00 -9.757749e-03 7.866365e-17 436.76800
In [56]: 1 #pip install pmdarima
         3.1.1 Naiive Forecasting
```

```
In [57]: 1 import pandas as pd
               2 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
                 # Load data
              5 df = pd.read_csv('addata.csv')
                     Prepare data
              8 df = df.set_index('Year')
             10 # Split data into train and test
             11 train = df[:-12]
12 test = df[-12:]
              14 # Naive forecast
              15 predictions = train['T'][-12:].values
             16
              17 # Calculate MSE
             19 mse = mean_squared_error(test['T'].values, predictions)
19 mae = mean_absolute_error(test['T'].values, predictions)
20 r2 = r2_score(test['T'].values, predictions)
             24 print('Naive MSE:', mse)
25 print('Naive MAE:', mae)
26 print('R-squared (R<sup>2</sup>):', r2)
             28 n = len(test)
             20 4 - 1
                                    ,
her of predictors (in this case, 1 for the naive forecast)
```

3.1.2 Exponential Smoothing Model

In [58]:	1	import matplotlib.pyplot as plt
	3	from sklearn.metrics import mean squared error
	4	from math import sqrt
	5	
	6	<pre>df = pd.read_csv('addata.csv')</pre>
	8	# Specify the number of steps to forecast
	9	forecast_steps = 3
	10	
	11	# Fit the Exponential Smoothing model
	12	<pre>model = ExponentialSmoothing(df['T'], seasonal='add', seasonal_periods=4, trend='add', damped=True)</pre>
	13	model_tit = model.tit()
	14	# Mabe forerasts
	16	forecast values = model fit.forecast(stens=forecast stens)
	17	
	18	# Evaluate the model using Mean Squared Error (MSE)
	19	<pre>y_test = df['T'][-forecast_steps:]</pre>
	20	<pre>mse = mean_squared_error(y_test, forecast_values)</pre>
	21	rmse = sqrt(mse)
	22	print(f"Mean Squared Error: {mse}")
	23	print(f"Root Mean Squared Error: {rmse}")
	24	# Calculate B caused
	25	w mean = df('t') mean()
	27	y_mean = or(r j,mean)**2).sum()
	28	ssr = ((v test - forecast values)**2).sum()
	29	$r^{2} = 1 - (ssr / sst)$
	30	<pre>print(f"R-squared (R²): {r2}")</pre>
	31	
	32	# Calculate Adjusted R-squared
	33	n = len(y_test)
	34	K = 3 # number of predictors (seasonal, trend, and damping terms)
	35	$aojusted_{2} = 1 - ((1 - r2) * (n - 1) / (n - k - 1))$
	30	princ(+ Aujusieu k-squareu; {aujusieu_r2; }
	C:\U 'dan mc	Jsers\KOMALP~1\AppData\Local\Temp/ipykernel_11200/1313939485.py:12: FutureWarning: the 'damped'' keyword is deprecated, use uped_trend' instead del = ExponentialSmoothing(df['T'], seasonal='add', seasonal_periods=4, trend='add', damped=True)
	ust wa	rngrambata\knaconda>\ilb\site-packages\statsmodeis\tsa\noitwinters\modei.py:42/: Futurewarning: After 0.13 initialization m be handled at model creation arnings.warn(
	Mear Root R-so Adju	n Squared Error: 16.58262528429255 : Mean Squared Error: 4.07217697114118 quared (R ²): 1.0 isted R-squared: 1.0

3.1.3 Forecasted Values



Section 3.2 has modelling and forecasting for depressive disorders.

```
3.2 Modelling for Depressive disorders
```

```
3.2.1 Naive Forecasting
```

```
[60]: 1 data = df_dd
[61]: 1 #renaming anxiety disorder to t and exporting to csv
2 datadd = data.rename(columns={'Depressive_Disorders': 'T'})
               4 datadd.to csv('dddata.csv')
[62]: 1 import pandas as pd
2 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
3 from sklearn.model_selection import train_test_split
                   # Load data
              6 df = pd.read_csv('dddata.csv')
                                 re data
                 df = df.set_index('Year')
             10
             11 # split data into train and test (70% train, 30% test)
12 train, test = train_test_split(df, test_size=0.3, shuffle=False)
             14 # Naive forecast
15 predictions = train['T'][-len(test):].values
             17 # Calculate MSE
             17 w Culculute Nos
18 mse = mean_squared_error(test['T'].values, predictions)
19 mae = mean_absolute_error(test['T'].values, predictions)
20 r2 = r2_score(test['T'].values, predictions)
             21
22 print('Naive MSE:', mse)
23 print('Naive MAE:', mae)
24 print('R-squared (R<sup>2</sup>):', r2)
             25
            26 n = len(test)
27 k = 1 # number
             26 n = len(test)
27 k = 1 # number of predictors (in this case, 1 for the naive forecast)
adjusted_r2 = 1 - ((1 - r2) * (n - 1) / (n - k - 1))
             30 print('Adjusted R-squared:', adjusted_r2)
            Naive MSE: 66881.41718150298
Naive MAE: 200.21860311666666
R-squared (R<sup>2</sup>): -0.6468914186754229
Adjusted R-squared: -0.6418040390417608
```

3.2.2 Exponential Smoothing Model

```
In [63]:
1 import matplotlib.pyplot as plt
2 from Statsmodel.tisa.Noltwinters import ExponentialSmoothing
3 from Statsmodel.tisa.Noltwinters import ExponentialSmoothing
6 forecast_steps = 3
1 # fit the Exponential Smoothing model
1 model_fit = model.fit()
1 model_fit = model_fit()
1 model_fit()
1
```

3.2.3 Forecasted Values

```
In [64]:
           1 import pandas as pd
               from statsmodels.tsa.holtwinters import ExponentialSmoothing
            3 from sklearn.metrics import mean_squared_error
              from math import sqrt
            6 # Print column names to identify the correct ones
             7 print(df.columns)
           9 # Specify the column name for the year
10 year_column = 'Year'
           12 # Specifying the number of steps to forecast (3 years)
13 forecast_steps = 3
           14
           15 # Create an empty dataframe to store the forecasted values
           16 forecast_df = pd.DataFrame()
           18 # Loop through each unique country in the dataframe
           19 for country in df['Country Name'].unique():
20 # Filter data for the current country
           21
                   country_df = df[df['Country Name'] == country]
           23
                   # Fit the Exponential Smoothing model
           24
25
                   model = ExponentialSmoothing(country_df['T'], seasonal='add', seasonal_periods=4, trend='add', damped=True)
                    model_fit = model.fit()
           26
27
                   # Make forecasts for the next 3 years
           28
                   forecast_values = model_fit.forecast(steps=forecast_steps)
           29
                  30
31
32
           33
           34
35
                                                             Forecast': forecast_values})
                   # Concatenate the current country's forecast dataframe to the overall forecast dataframe
forecast_df = pd.concat([forecast_df, country_forecast_df], ignore_index=True)
           36
           38
           39 # Print the forecasted values dataframe
40 print(forecast_df)
           41
           42 forecast df.to csv('forecasted values dipressive disorders.csv', index=False)
           43
          Index(['Unnamed: 0', 'Country Name', 'Code', 'Year', 'A', 'B', 'F', 'G',
                                                                                           'H',
                                                                                                                                                      'T'],
                 dtype='object')
          C:\Users\KOMALP~1\AppData\Local\Temp/ipykernel_11200/923382177.py:24: FutureWarning: the 'damped'' keyword is deprecated, us
             'damped_trend' instead
          e
          model = ExponentialSmoothing(country_df['T'], seasonal='add', seasonal_periods=4, trend='add', damped=True)
C:\ProgramData\Anaconda3\lib\site-packages\statsmodels\tsa\holtwinters\model.py:427: FutureWarning: After 0.13 initializatio
```

Section 3.3 has modelling evaluation and forecasting for Schizophrenia :



3.3.2 Exponential Smoothing Model



3.3.3 Forecasted Values



Section 3.4 has modelling evaluation and forecasting for eating disorders :

			•										
n [71]:	1	1 #Running standard scaler for all indicator columns with target variable as anxiety disorders											
	<pre>3 from sklearn.preprocessing import StandardScaler 4</pre>												
	5	data = df	f_ed										
	6	columns t	to scale = ['	A'. '	8'.'F	'.'g'.'н'	1						
	8				1		-						
	10	scaler =	StandardScal	ler()									
	11												
	12	# Fit and	d transform t	the nu	meric	al column	s						
n [72]:	1	<pre>#renaming depressive disorder to t and exporting to csv dataad = data papame(c)umes_(lating Disorders', 'IT'))</pre>											
	3	<pre>2 datadu = data.rename(columns={ Eating_Disorders : 1 }) 3</pre>											
	4 dataed.to_csv('eddata.csv')												
	3.4.1 Naiive Forecasting												
. (72).		innent av	adag ag ad										
n [/s]:	2	I Import pandas as po 2 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score											
	3	# Logd de	+-										
	5	df = pd.r	read_csv('edd	lata.c	sv')								
	6	df.head()											
ut[73]:				- .			_	_	-		_		
	_	Unnamed: 0	Country Name	Code	Year	A	в	F	G	н	1		
	0	0	Afghanistan	AFG	1990	91.735775	27.1965	44.252894	8.166848	7903.252158	26.471115		
							07 408E	02 702270	0 404000	7003 252158	25 542821		
	1	1	Afghanistan	AFG	1991	91.735775	27.1900	82.783270	0.121000	1000.202100	20.040001		
	1 2	1	Afghanistan Afghanistan	AFG AFG	1991	91.735775 91.735775	27.1965	92.745990	8.168000	7903.252158	24.637949		
	1 2 3	1 2 3	Afghanistan Afghanistan Afghanistan	AFG AFG AFG	1991 1992 1993	91.735775 91.735775 91.735775	27.1965 27.1965 27.1965	92.745990 92.672950	8.168000 8.123000	7903.252158 7903.252158	24.637949 23.863169		

n [74]:
1
2 # Prepare data
3 df = df.set_index('Year')
4
4
5 # Split data into train and test
6 train = df[-12]
7 test = df[-12:]
8
9 # Naive forecast
10 predictions = train['T'][-12:].values
11
12 # Calculate MSE
13 mse = mean_squared_error(test['T'].values, predictions)
14 mae = mean_absolute_error(test['T'].values, predictions)
15 r2 = r2_score(test['T'].values, predictions)
16

3.4.2 Exponential Smoothing Model

3.4 Modelling for Eating Disorders

```
In [75]:
1 import matplotlib.pyplot as plt
2 from stlamodels.tss.holtwinters import ExponentialSmoothing
7 from stlamodels.tss.holtwinters import men_squared_error
4 from math import sqrt
4 from math import of steps to forecast
7 forecast_steps = 3
9
9 df = pd.read_sw('eddata.csv')
14 # fit the Exponential Smoothing model
12 model = Exponential Smoothing model
13 model = Exponential Smoothing model
14 model = Exponential Smoothing model
15 model_fit = model.fit.forecast_steps
16 # Nuke forecasts
17 forecast_values = model_fit.forecast_steps)
17 forecast_values = model_fit.forecast_values)
18 # Kuthe the model using Nean Squared Error (NSE)
19 y_test = df[''][_forecast_stepsi]
10 me = mean_squared_Error: (mse)''
11 me = mean_squared_Error: (mse)''
12 for if(''Not Nean Squared Error: (mse)''
13 # Evolute the model by calculating Acquared
14 y_mean = df['']['.exest_values)
15 # Evolute the model by calculating Acquared
16 y_mean = df[''.forest_values,'+2).sum()
17 # Evolute the model by calculating Acquared
17 y_test = (f_1''.forest_values,'+2).sum()
18 st = ((d_1''.forest_values,'+2).sum()
19 print(f'Most_Mean Squared Error: (mse)'')
19 print(f'Most_Mean Squared Error:
```

Mean Squared Error: 0.026787367780580437 Root Mean Squared Error: 0.16366846910929556 R-squared (R³): 1.0 Adjusted R-squared: 1.0 3.4.3 Forecasted Values



Forecasted CSVs are used to create visualizations in power bi. Below reports show the top 5 countries forecasted to have highest number of cases of these mental health disorders. For this analysis average value of forecast for 2020, 2021 and 2022 is taken and plotted :





Below dashboard shows countries where percentage increase over these three years will be highest. For this visualization a new metric for percentage increase is calculated and maximum of the value is taken , the dashboard is then filtered for showing top 5 values in power bi.



To further drill down to workplace related mental health factors, the downloaded dataset is explored against all factors and pie charts are created in power bi to show the distribution of mental health disorders according to these factors.



Percentage distribution of various factors affecting workplace mental well being