

Mental Fitness Tracker :

The Mental Fitness Tracker is an AI-powered project aimed at monitoring and supporting mental well-being using advanced algorithms. It leverages regression models to analyze and predict mental fitness based on various factors. The project utilizes a dataset containing information about mental and substance use disorders, prevalence rates, and other relevant factors.

About the Data:

- **Data Overview:** This is a "Mental Health".csv data

Import Required Libraries:

```
1 # Ignore warning messages to prevent them from being displayed during code execution
2 import warnings
3 warnings.filterwarnings('ignore')

1 import numpy as np    # Importing the NumPy library for linear algebra operations
2 import pandas as pd   # Importing the Pandas library for data processing and CSV file handling

1 from google.colab import drive  # Importing the necessary library to mount Google Drive
2 drive.mount('/content/drive')    # Mounting Google Drive to the specified directory '/content/drive'

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

1 import seaborn as sns          # Importing the Seaborn library for statistical data visualization
2 import matplotlib.pyplot as plt # Importing the Matplotlib library for creating plots and visualizations
3 import plotly.express as px     # Importing the Plotly Express library for interactive visualizations
4
```

Exploratory Data Analysis

Load and Prepare Data:

```
1 # Reading the CSV file containing prevalence data for mental and substance use disorder
2 df1 = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/prevalence-by-mental-and-substance-use-disorder.csv")
3
4 # Reading the CSV file containing data on mental and substance use as a share of disease
5 df2 = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/mental-and-substance-use-as-share-of-disease.csv")

1 df1.head()
2 # Displaying the first few rows of the DataFrame df1
```

Entity	Code	Year	Prevalence - Schizophrenia	Prevalence - Bipolar disorder	Prevalence - Eating disorders	Prevalence - Anxiety disorders	Prevalence - Drug use disorders	Prevalence - Depression
			- Sex: Both - Age: Age-standardized (Percent)					
0	Afghanistan	AFG	1990	0.228979	0.721207	0.131001	4.835127	0.454202
1	Afghanistan	AFG	1991	0.228120	0.719952	0.126395	4.821765	0.447112
2	Afghanistan	AFG	1992	0.227328	0.718418	0.121832	4.801434	0.441190
3	Afghanistan	AFG	1993	0.226468	0.717452	0.117942	4.789363	0.435581

```
1 df2.head(10)
2 # Displaying the first few rows of the DataFrame df2
```

Entity Code Year DALYs (Disability-Adjusted Life Years) - Mental disorders - Sex: Both - Age: All Ages (Percent)

```

1 data = pd.merge(df1, df2)
2 data.head(10)
3 # Merging df1 and df2 into a single DataFrame and displaying the first 10 rows of the merged DataFrame
4

```

Entity	Code	Year	Prevalence - Schizophrenia - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Bipolar disorder - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Eating disorders - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Anxiety disorders - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Drug use disorders - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Depressive disorders - Sex: Both - Age: Age-standardized (Percent)
0	Afghanistan	AFG	1990	0.228979	0.721207	0.131001	4.835127	0.454202
1	Afghanistan	AFG	1991	0.228120	0.719952	0.126395	4.821765	0.447112
2	Afghanistan	AFG	1992	0.227328	0.718418	0.121832	4.801434	0.441190
3	Afghanistan	AFG	1993	0.226468	0.717452	0.117942	4.789363	0.435581
4	Afghanistan	AFG	1994	0.225567	0.717012	0.114547	4.784923	0.431822
5	Afghanistan	AFG	1995	0.224713	0.716686	0.111129	4.780851	0.428578
6	Afghanistan	AFG	1996	0.223690	0.716388	0.107786	4.777272	0.426393
7	Afghanistan	AFG	1997	0.222424	0.716143	0.103931	4.775242	0.423720

• Data Cleaning

```

1 data.isnull().sum()
2 # Checking the count of missing values (null values) in each column of the DataFrame data
3

```

Entity	0
Code	690
Year	0
Prevalence - Schizophrenia - Sex: Both - Age: Age-standardized (Percent)	0
Prevalence - Bipolar disorder - Sex: Both - Age: Age-standardized (Percent)	0
Prevalence - Eating disorders - Sex: Both - Age: Age-standardized (Percent)	0
Prevalence - Anxiety disorders - Sex: Both - Age: Age-standardized (Percent)	0
Prevalence - Drug use disorders - Sex: Both - Age: Age-standardized (Percent)	0
Prevalence - Depressive disorders - Sex: Both - Age: Age-standardized (Percent)	0
Prevalence - Alcohol use disorders - Sex: Both - Age: Age-standardized (Percent)	0
DALYs (Disability-Adjusted Life Years) - Mental disorders - Sex: Both - Age: All Ages (Percent)	0
dtype: int64	

```

1 data.drop('Code', axis=1, inplace=True)
2 # Dropping the 'Code' column from the DataFrame data

```

```

1 data.head(10)
2 # Displaying the first 10 rows of the DataFrame data

```

Entity	Year	Prevalence - Schizophrenia - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Bipolar disorder - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Eating disorders - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Anxiety disorders - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Drug use disorders - Sex: Both - Age: Age-standardized (Percent)	Prevalence - Depressive disorders - Sex: Both - Age: Age-standardized (Percent)
0	1990	0.228979	0.721207	0.131001	4.835127	0.454202	5.12529
1	1991	0.228120	0.719952	0.126395	4.821765	0.447112	5.11630
2	1992	0.227328	0.718418	0.121832	4.801434	0.441190	5.10655
3	1993	0.226468	0.717452	0.117942	4.789363	0.435581	5.10032
4	1994	0.225567	0.717012	0.114547	4.784923	0.431822	5.09942
5	1995	0.224713	0.716686	0.111129	4.780851	0.428578	5.09849
6	1996	0.223690	0.716388	0.107786	4.777272	0.426393	5.10058
7	1997	0.222424	0.716143	0.103931	4.775242	0.423720	5.10547

```

1 data.size, data.shape
2 # Calculating the size (total number of elements) and shape (number of rows, number of columns) of the DataFrame data
3
4 (68400, (6840, 10))

```

```

1 data.set_axis(['country', 'Year', 'Schizophrenia', 'Bipolar_disorder', 'Eating_disorder', 'Anxiety', 'drug_usage', 'depression', 'alcohol', 'mental_fitness'])
2 # Renaming the columns of the DataFrame data with more descriptive names

```

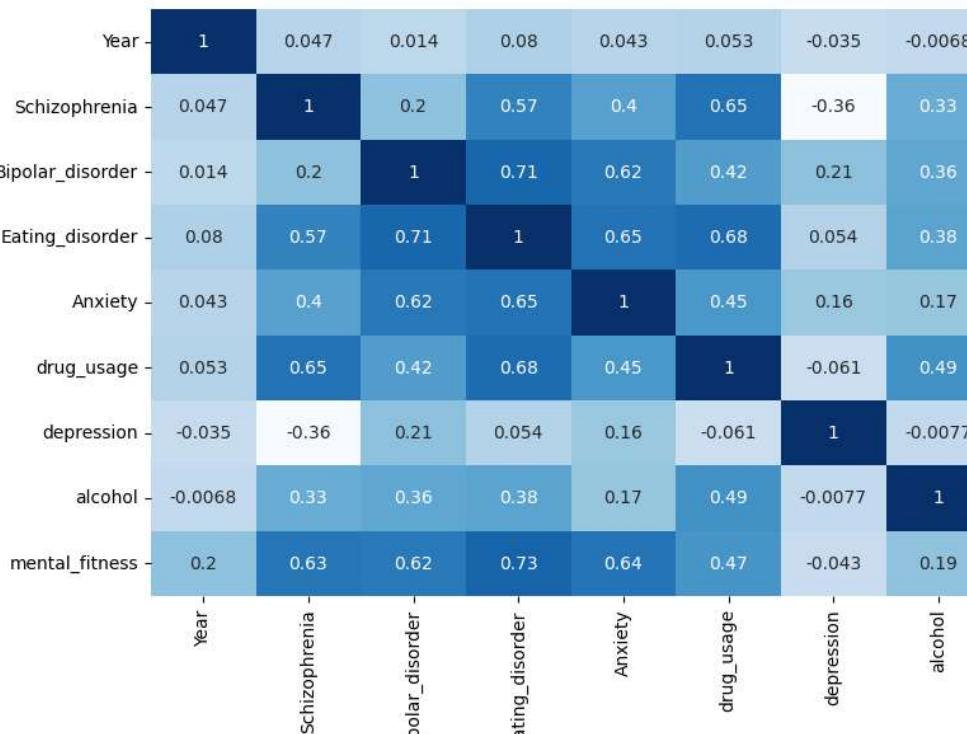
```
1 data.head(10)
```

	country	Year	Schizophrenia	Bipolar_disorder	Eating_disorder	Anxiety	drug_usage	depression
0	Afghanistan	1990	0.228979	0.721207	0.131001	4.835127	0.454202	5.125291
1	Afghanistan	1991	0.228120	0.719952	0.126395	4.821765	0.447112	5.116306
2	Afghanistan	1992	0.227328	0.718418	0.121832	4.801434	0.441190	5.106558
3	Afghanistan	1993	0.226468	0.717452	0.117942	4.789363	0.435581	5.100328
4	Afghanistan	1994	0.225567	0.717012	0.114547	4.784923	0.431822	5.099424
5	Afghanistan	1995	0.224713	0.716686	0.111129	4.780851	0.428578	5.098495
6	Afghanistan	1996	0.223690	0.716388	0.107786	4.777272	0.426393	5.100580
7	Afghanistan	1997	0.222424	0.716143	0.103931	4.775242	0.423720	5.105474
8	Afghanistan	1998	0.221129	0.716139	0.100343	4.777377	0.422491	5.113707
9	Afghanistan	1999	0.220065	0.716323	0.097946	4.782067	0.421215	5.120480

Visulaization

```
1 plt.figure(figsize=(12, 6))
2 sns.heatmap(data.corr(), annot=True, cmap='Blues')
3 plt.plot()
4 # Creating a heatmap of the correlation matrix for the columns in the DataFrame data
```

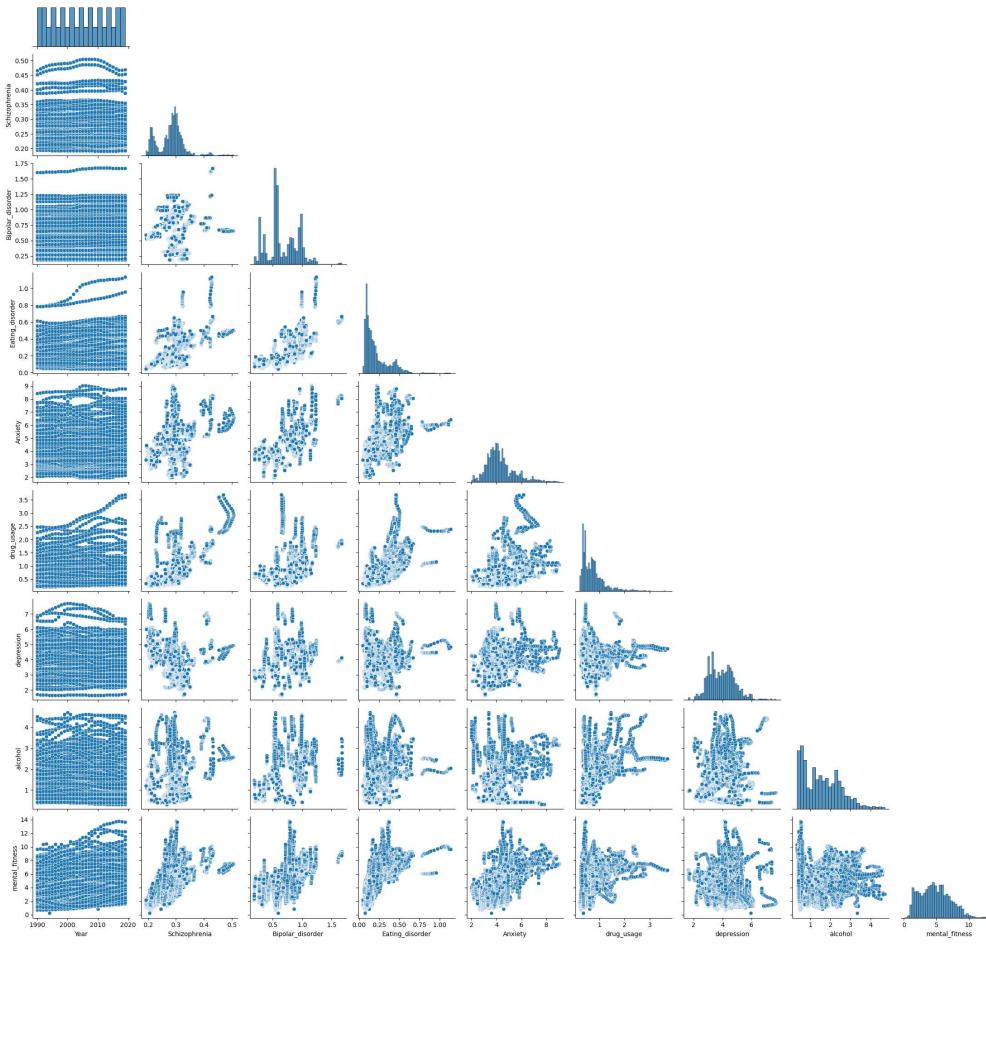
```
[]
```



Key Points:

Eating_disorder is positively Correlated to mental_fitness and vice-versa as our eating choice affect our mental health

```
1 sns.pairplot(data, corner=True)
2 plt.show()
3 # Creating a pairwise scatter plot matrix for the columns in the DataFrame data
```



```

1 mean = data['mental_fitness'].mean()
2 # Calculating the mean value of the 'mental_fitness' column in the DataFrame data
3 mean
4 # Displaying the calculated mean value

```

4.8180618117506135

```

1 fig = px.pie(data, values='mental_fitness', names='Year')
2 fig.show()
3 # Creating a pie chart to visualize the distribution of mental fitness across different years

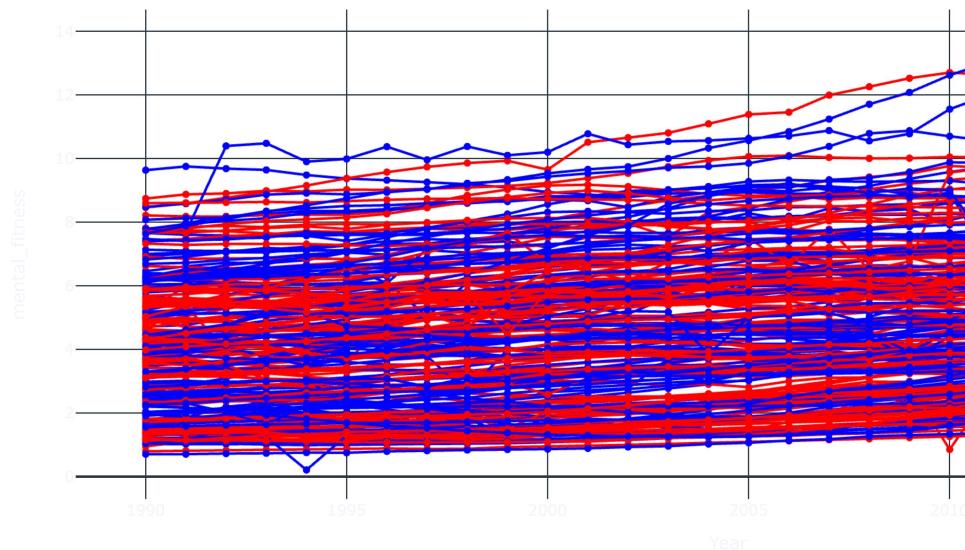
```



```

1 fig = px.line(data, x='Year', y='mental_fitness', color='country', markers=True, color_discrete_sequence=['red', 'blue'], template='plotly_dark')
2 fig.show()
3 # Creating a line plot to visualize the trend of mental fitness over time, differentiated by country
4

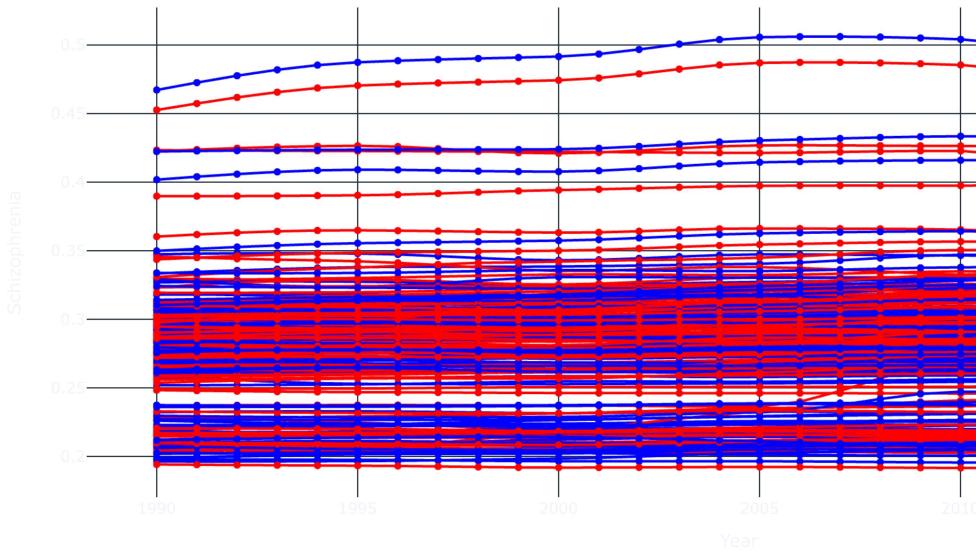
```



```

1 fig = px.line(data, x='Year', y='Schizophrenia', color='country', markers=True, color_discrete_sequence=['red', 'blue'], template='plotly_dark')
2 fig.show()
3 # Creating a line plot to visualize the trend of schizophrenia prevalence over time, differentiated by country

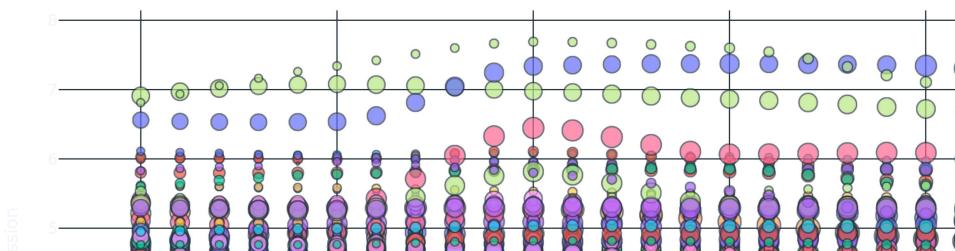
```



```

1 fig = px.scatter(data, x='Year', y='depression', color='country', size='mental_fitness', hover_data=['Schizophrenia'], template='plotly_dark')
2 fig.show()
3 # Creating a scatter plot to visualize the relationship between depression and mental fitness, differentiated by country

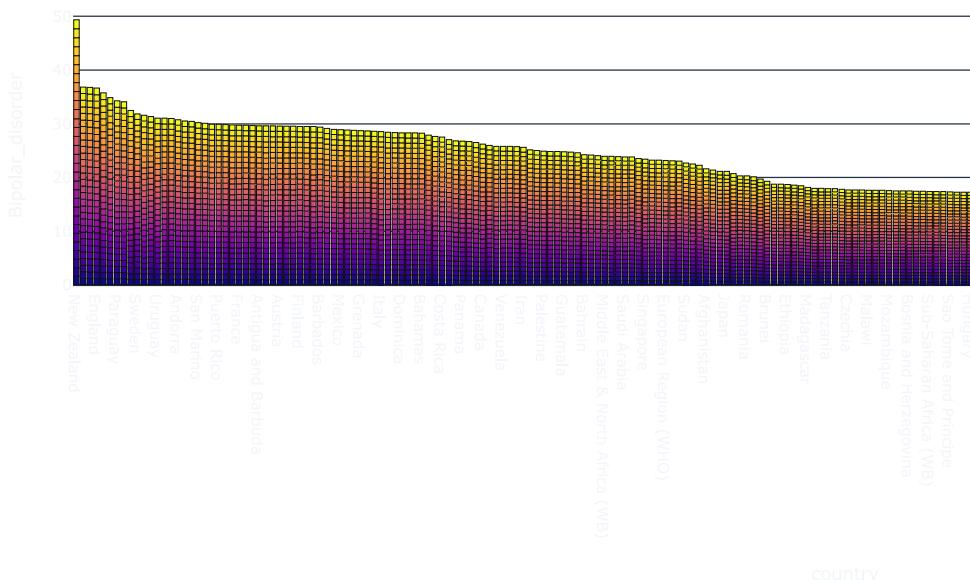
```



```

1 fig = px.bar(data, x='country', y='Bipolar_disorder', color='Year', barmode='group', template='plotly_dark')
2 fig.update_layout(xaxis={'categoryorder':'total descending'})
3 fig.show()
4 # Creating a grouped bar chart to compare the prevalence of Bipolar Disorder across countries and years

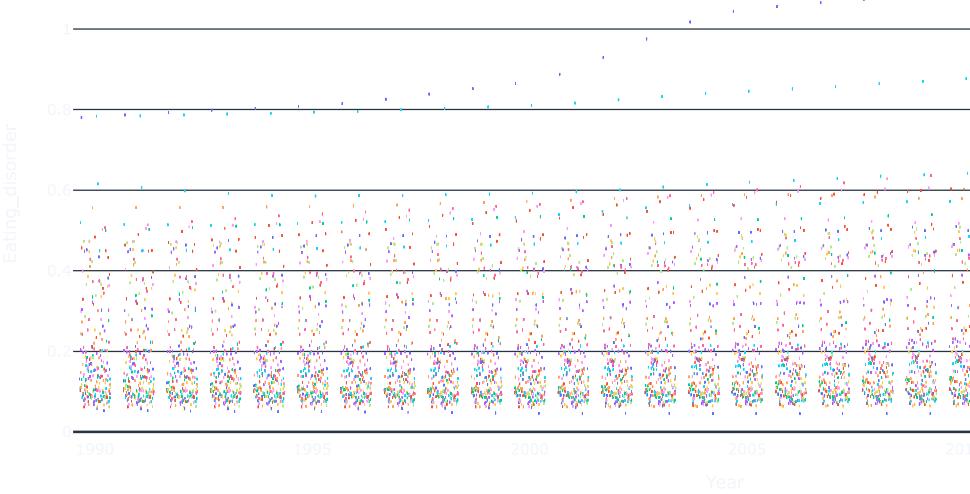
```



```

1 fig = px.box(data, x='Year', y='Eating_Disorder', color='country', template='plotly_dark')
2 fig.show()
3 # Creating a box plot to visualize the distribution of Eating Disorder prevalence across different years and countries

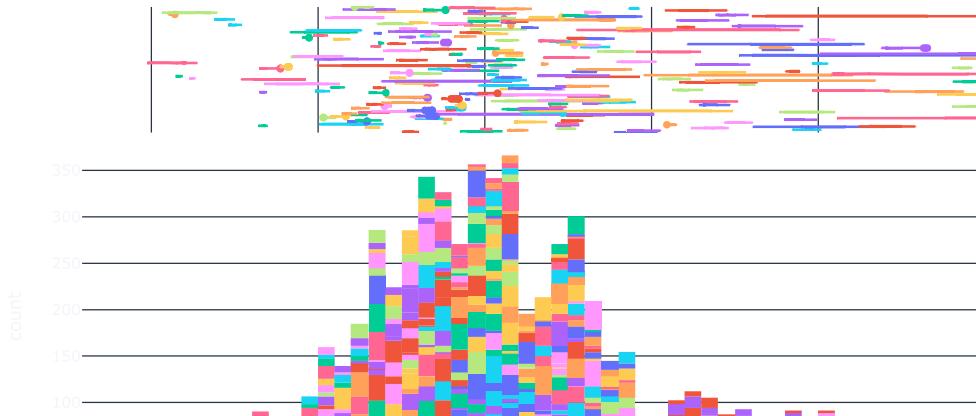
```



```

1 fig = px.histogram(data, x='Anxiety', color='country', marginal='box', template='plotly_dark')
2 fig.show()
3 # Creating a histogram with overlaid box plots to visualize the distribution of Anxiety levels across different countries

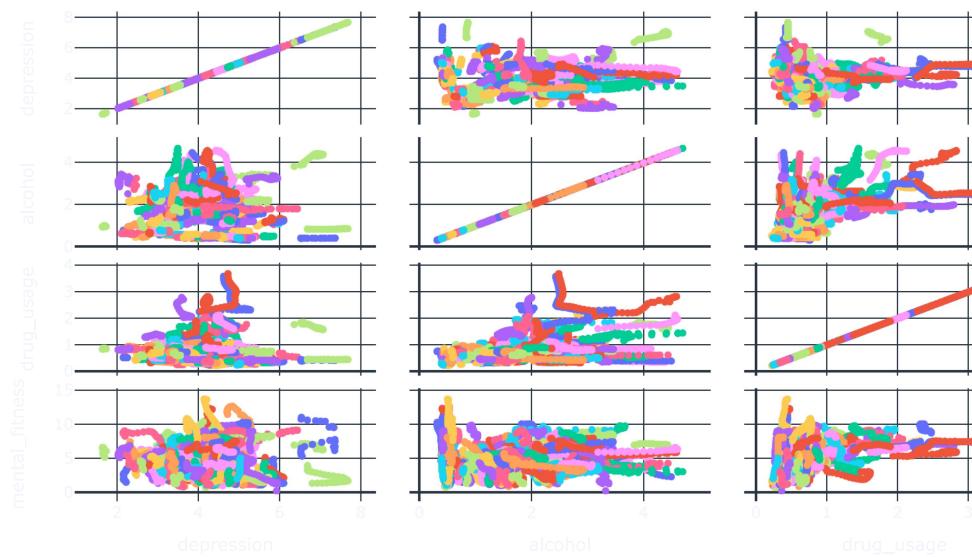
```



```

1 fig = px.scatter_matrix(data, dimensions=['depression', 'alcohol', 'drug_usage', 'mental_fitness'], color='country', template='plotly_dark')
2 fig.show()
3 # Creating a scatter matrix plot to visualize the relationships between multiple variables (depression, alcohol, drug usage, mental fitness) across different

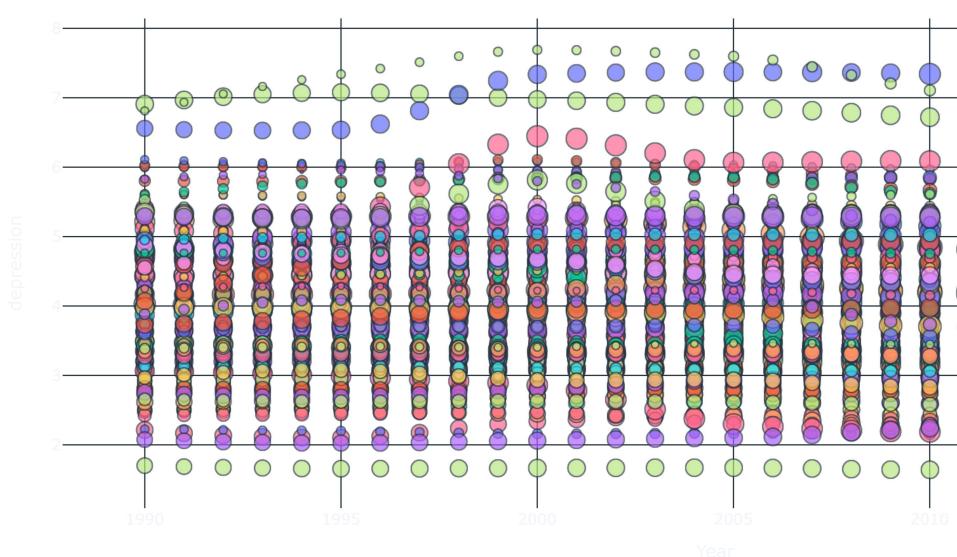
```



```

1 fig = px.scatter(data, x='Year', y='depression', size='mental_fitness', color='country', template='plotly_dark')
2 fig.show()
3 # Creating a scatter plot to visualize the relationship between depression, mental fitness, and year, differentiated by country

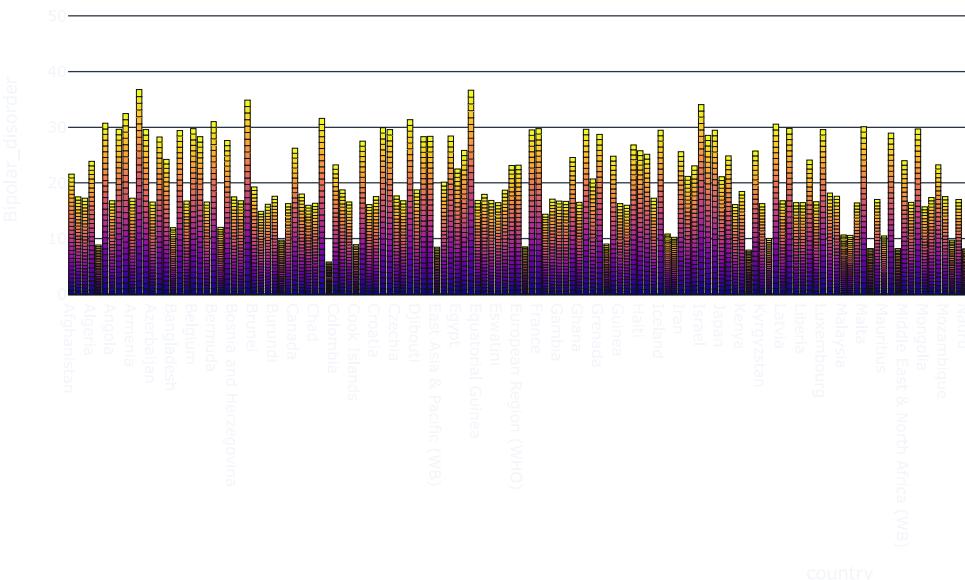
```



```

1 fig = px.bar(data, x='country', y='Bipolar_Disorder', color='Year', barmode='group', template='plotly_dark')
2 fig.show()
3 # Creating a grouped bar plot to compare the prevalence of Bipolar Disorder across different countries and years

```



```
1 df= data
2 df.head(10)
```

	country	Year	Schizophrenia	Bipolar_disorder	Eating_disorder	Anxiety	drug_usage	depression
0	Afghanistan	1990	0.228979	0.721207	0.131001	4.835127	0.454202	5.125291
1	Afghanistan	1991	0.228120	0.719952	0.126395	4.821765	0.447112	5.116306
2	Afghanistan	1992	0.227328	0.718418	0.121832	4.801434	0.441190	5.106558
3	Afghanistan	1993	0.226468	0.717452	0.117942	4.789363	0.435581	5.100328
4	Afghanistan	1994	0.225567	0.717012	0.114547	4.784923	0.431822	5.099424
5	Afghanistan	1995	0.224713	0.716686	0.111129	4.780851	0.428578	5.098495
6	Afghanistan	1996	0.223690	0.716388	0.107786	4.777272	0.426393	5.100580
7	Afghanistan	1997	0.222424	0.716143	0.103931	4.775242	0.423720	5.105474
8	Afghanistan	1998	0.221129	0.716139	0.100343	4.777377	0.422491	5.113707
9	Afghanistan	1999	0.220065	0.716323	0.097946	4.782067	0.421215	5.120480

```
1 # Display information about the DataFrame
2 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 6840 entries, 0 to 6839
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   country          6840 non-null   object 
 1   Year              6840 non-null   int64  
 2   Schizophrenia    6840 non-null   float64
 3   Bipolar_disorder 6840 non-null   float64
 4   Eating_disorder  6840 non-null   float64
 5   Anxiety           6840 non-null   float64
 6   drug_usage        6840 non-null   float64
 7   depression         6840 non-null   float64
 8   alcohol            6840 non-null   float64
 9   mental_fitness    6840 non-null   float64
dtypes: float64(8), int64(1), object(1)
memory usage: 587.8+ KB
```

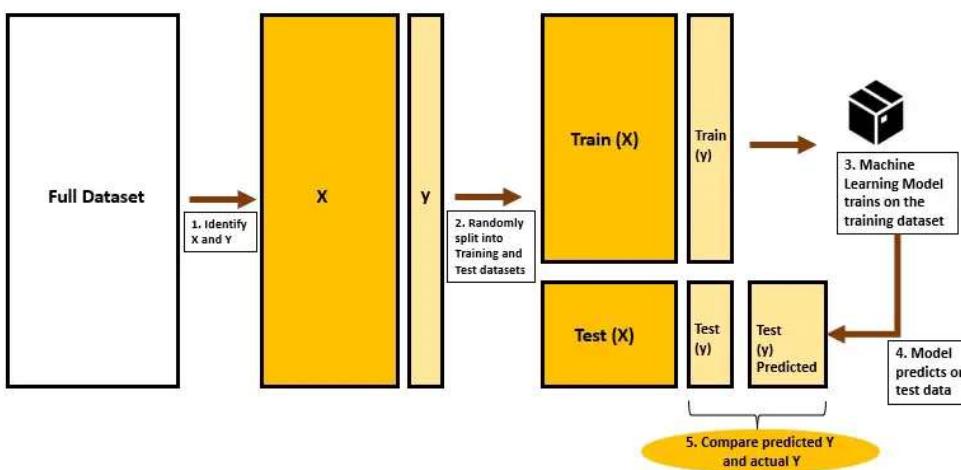
```
1 from sklearn.preprocessing import LabelEncoder
2
3 # Create an instance of LabelEncoder
4 label_encoder = LabelEncoder()
5
6 # Iterate over the columns of the DataFrame
7 for column in df.columns:
8     # Check if the column data type is 'object' (categorical)
9     if df[column].dtype == 'object':
10         # Use the label encoder to transform the categorical data into numerical labels
11         df[column] = label_encoder.fit_transform(df[column])
12
13 # Comment: The above code snippet utilizes the LabelEncoder from scikit-learn to encode categorical variables in the DataFrame. It iterates over each column
```

```
1 # Retrieve the shape of the DataFrame
2 df_shape = df.shape
3
4 # Print the shape of the DataFrame
5 print(df_shape)
```

(6840, 10)

• Split Data - (6840,10)

- In this step, the data is split into two parts: training and testing datasets. This allows us to train our model on the training dataset and evaluate its accuracy on unseen test data.
- The shape of the split data is (6840, 10), indicating that the dataset contains 6,840 rows and 10 columns.



```
1 # Splitting the Data into Training and Testing Sets
2
3 # In this step, the dataset is split into input features (X) and target variable (y). The target variable, 'mental_fitness', is dropped from the input features.
4 # The target variable, 'mental_fitness', is assigned to y.
5 # The data is then split into training and testing sets using the train_test_split function from the sklearn.model_selection module. The test_size parameter
6 # The random_state parameter is set to 2 to ensure reproducibility of the split.
7
8 import sklearn.model_selection as ms
9
10 X = df.drop('mental_fitness', axis=1)
11 y = df['mental_fitness']
12
13 xtrain, xtest, ytrain, ytest = ms.train_test_split(X, y, test_size=0.20, random_state=2)

1 # Training (6840, 10)
2 # 6840*80/100 = 5472
3 # 6840*20/100 = 1368
4 # Training (5472, 10)
5 # The training data consists of 5,472 samples (80% of the total data) with 10 input features.
6
7 # Testing (1368, 10)
8 # The testing data consists of 1,368 samples (20% of the total data) with 10 input features.
9
10 print("xtrain:", xtrain.shape)
11 print("xtest:", xtest.shape)
12 print("ytrain:", ytrain.shape)
13 print("ytest:", ytest.shape)
14
15 # Comment: The code above prints the shape of the training and testing data.
16 # The training data (xtrain) has a shape of (5472, 10), indicating 5,472 samples and 10 input features.
17 # The testing data (xtest) has a shape of (1368, 10), indicating 1,368 samples and 10 input features.
18 # The target variable for training (ytrain) has a shape matching the number of samples in xtrain, and the target variable for testing (ytest) has a shape matching the number of samples in xtest.
19
20 # Ensure that the shapes printed match the expected dimensions of your training and testing data.

xtrain: (5472, 9)
xtest: (1368, 9)

ytrain: (5472,)
ytest: (1368,)
```

• Key Points:

- Usually we take more and more data in training so it's easy for the model to learn with more data
- Usually, it is beneficial to use a larger amount of data during model training. This allows the model to learn from a more extensive and diverse set of examples, leading to better performance and generalization.

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import mean_squared_error, r2_score
3
4 # Create a Linear Regression model and fit it to the training data
5 lr = LinearRegression()
6 lr.fit(xtrain, ytrain)
7
8 # Make predictions on the training data
9 ytrain_pred = lr.predict(xtrain)
10
```

```

11 # Calculate Mean Squared Error (MSE) and R-squared score for training data
12 mse = mean_squared_error(ytrain, ytrain_pred)
13 r2 = r2_score(ytrain, ytrain_pred)
14
15 # Print the performance metrics of the Linear Regression model on the training set
16 print("Linear Regression model performance on the training set")
17 print("-----")
18 print("MSE: {}".format(mse))
19 print("RMSE: {}".format(np.sqrt(mse))) # Corrected the variable name to 'mse'
20 print("R2 Score: {}".format(r2))
21
22
23 #Float format wit 2 Decimals:
24 #print("MSE: {:.2f}".format(mse))
25 #print("RMSE: {:.2f}".format(np.sqrt(mse))) # Corrected the variable name to 'mse'
26 #print("R2 Score: {:.2f}".format(r2))

    Linear Regression model performance on the training set
-----
MSE: 1.389959372405798
RMSE: 1.1789653821914357
R2 Score: 0.7413245790025275

```

```

1 from sklearn.ensemble import RandomForestRegressor
2 from sklearn.metrics import mean_squared_error, r2_score
3
4 # Create a Random Forest Regressor model and fit it to the training data
5 rf = RandomForestRegressor()
6 rf.fit(xtrain, ytrain)
7
8 # Make predictions on the training data
9 ytrain_pred = rf.predict(xtrain)
10
11 # Calculate Mean Squared Error (MSE), RMSE, and R-squared score for training data
12 mse = mean_squared_error(ytrain, ytrain_pred)
13 rmse = np.sqrt(mse)
14 r2 = r2_score(ytrain, ytrain_pred)
15
16 # Print the performance metrics of the Random Forest Regressor model on the training set
17 print("Random Forest Regressor model performance on the training set")
18 print("-----")
19 print("MSE: {}".format(mse))
20 print("RMSE: {}".format(rmse))
21 print("R2 Score: {}".format(r2))

    Random Forest Regressor model performance on the training set
-----
MSE: 0.0048916602143839775
RMSE: 0.06994040473420195
R2 Score: 0.999089648021048

```

• Evaluation:

- In this part, we will compare the scores of above two models.

```

1 # Compare the scores of the Linear Regression and Random Forest Regressor models
2
3 # Calculate the scores of Linear Regression model
4 lr_scores = lr.score(xtest, ytest)
5
6 # Calculate the scores of Random Forest Regressor model
7 rf_scores = rf.score(xtest, ytest)
8
9 print("Evaluation of Models")
10 print("-----")
11 print("Linear Regression Score: {:.2f}".format(lr_scores))
12 print("Random Forest Regressor Score: {:.2f}".format(rf_scores))

    Evaluation of Models
-----
Linear Regression Score: 0.76
Random Forest Regressor Score: 0.99

```

The evaluation results showed that the Random Forest Regressor model outperformed the Linear Regression model, achieving a higher score of 0.99 compared to 0.76.

This indicates that the Random Forest Regressor model provided a better fit to the data and could make more accurate predictions regarding mental fitness.

