



Group Member

- Chenchen Li - c4588
- Stephanie Yao - yz508
- Kristin Qi - kq243

Link to your Public Github repository with Final report :

<https://github.com/chenchenli123>

Clarification:

1. Some results are showing in the HTML format in the output, so may not be able to clearly show all the graphs. In order to avoid some missing graphs, we put some output below the code chunk using PDF format.

Submission Due Date: 03/07/2025

World Happiness Classification Competition

Goals:

- Indebted from the model's function
- Understand what the parameter control
- Learn from the model's experimentation process
- Make a good looking notebook report
- Upload as a personal project on Github

Overall Steps:

1. Load datasets and merge them.
2. Preprocess data using sklearn Column Transformer/ White and Save Preprocessor Function
3. Fit model on preprocessed data and save preprocessor function and model
4. Generate predictions from X_test data and submit predictions

0. Loading Datasets

Loading the World Happiness 2023 datasets

```
import pandas as pd
from sklearn.model_selection import train_test_split

# Upload the CSV dataset to Jupyter in the left bar
df = pd.read_csv('HWS_2023.csv')

# Inspect the first few rows to understand the structure
df.head()
```

	country	region	happiness_score	gdp_per_capita	social_support	healthy_life_expectancy	freedom_to_make_life_choices	generosity	perceptions_of_corruption
0	Costa Rica	Western Europe	7.584	1.989	0.98	0.72	0.36	0.026	
1	Denmark	Western Europe	7.586	1.949	1.548	0.827	0.734	0.308	0.026
2	Iceland	Western Europe	7.530	1.936	1.600	0.809	0.739	0.250	0.107
3	Ireland	Western Europe	7.473	1.883	1.520	0.827	0.688	0.234	0.108
4	Netherlands	Western Europe	7.403	1.940	1.488	0.845	0.679	0.251	0.104
...
100	Congo (Kinshasa)	Sub-Saharan Africa	3.037	0.031	0.784	0.105	0.275	0.183	0.988
101	Zimbabwe	Sub-Saharan Africa	3.034	0.758	0.881	0.089	0.360	0.112	0.117
104	Sierra Leone	Sub-Saharan Africa	3.138	0.070	0.540	0.080	0.371	0.180	0.901
105	Libania	Sub-Saharan Africa	2.988	1.417	0.535	0.084	0.120	0.040	0.927
106	Algeria	North Africa	1.869	0.845	0.000	0.087	0.000	0.060	0.958
107

Next steps: [Generate code with Copilot](#) [View recommended plots](#) [New interactive plot](#)

```
# Check the column names in this dataset
df.columns
```

```
# Select the columns: "region", "happiness_score", "gdp_per_capita",
"social_support", "healthy_life_expectancy",
"freedom_to_make_life_choices", "generosity",
"perceptions_of_corruption",
"region"
df.select_dtypes(include=[object])
```

```
# Convert the regression target ("happiness_score") into classification labels
# We'll use quartiles to create a happiness category: Very Low, Low, High, Very High
```

```
# Define quartiles
df['happiness_category'] = pd.qcut(df['happiness_score'],
                                q=4,
                                labels=["Very Low", "Low", "Average", "High", "Very High"])
```

```
# Select features and target
X = df.drop(columns=["happiness_score", "happiness_category"])
y = df['happiness_category']

# Split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42, stratify=y)
```

```
# Convert X_train and X_test to numerical labels
# Create labels for y_train and y_test
y_train_labels = y_train.astype('category').cat.codes
y_test_labels = y_test.astype('category').cat.codes
```

Write in the next cell what the y_train.astype('category').cat.codes line does. What is the difference between y_train_labels and y_train?

Answer: The line y_train.astype('category').cat.codes converts the categorical variable y_train (which contains labels like "Very Low", "Low", "Average", "High", "Very High") into numerical codes.

Difference between y_train_labels and y_train:

y_train_labels: Contains the original categorical labels (e.g., "Very Low", "Low", "Average", etc.).

y_train: Contains the corresponding numerical codes (0, 1, 2, 3, 4) that can be used in machine learning models.

This transformation allows classification models to work with categorical data numerically while preserving the ordinal nature of happiness categories.

Add new data

```
# Truncated and cleaned up region data to merge
country_data = pd.read_csv('country_data.csv')
```

	country_name	population	population_below_poverty_line	hdi	life_expectancy	expected_years_of_education	mean_years_of_education	gdi
0	India	1395800000	21.9	0.645554	68.827	11.088855	6.988004	5683.471766
1	Nigeria	195980000	20.0	0.537165	52.827	6.970482	6.000000	5445.801084
2	Kenya	52000000	48.2	0.539660	56.910	10.000000	6.000000	5680.000000
3	Pakistan	197000000	28.5	0.538854	66.800	6.108919	6.000000	5501.173074
4	Bangladesh	160000000	21.5	0.528624	71.865	10.176706	6.341277	5281.488222

Next steps: [Generate code with Copilot](#) [View recommended plots](#) [New interactive plot](#)

```
# Merge in new data to X_train and X_test by taking "country" from first table and "country_name" from 2nd table.
# Data check which variables are common in both the datasets, and which type of merge will you perform for the best results.
# Hint: Look at the "how" parameter of merge function of pandas.
```

```
# Merge "country" column into both X_train and y_train for merging purposes
X_train = X_train.merge(country_data[['country', 'life_expectancy', 'expected_years_of_education', 'mean_years_of_education', 'gdi']],
                        on='country',
                        how='left',
                        suffixes=('', '_country_data'))
```

```
# Perform an inner merge with country_data based on country names
X_train = X_train.merge(country_data[['country', 'life_expectancy', 'expected_years_of_education', 'mean_years_of_education', 'gdi']],
                        on='country',
                        how='inner',
                        suffixes=('', '_country_data'))
```

	country_x	region	gdp_per_capita	social_support	healthy_life_expectancy	freedom_to_make_life_choices	generosity	perceptions_of_corruption	country_y	population	population_below_poverty_line	hdi	life_expectancy	expected_years_of_education	mean_years_of_education	gdi
0	Mongolia	Sub-Saharan Africa	0.882	0.739	0.176	0.107	0.127	0.154	Mongolia	2927885	75.7	0.521468	65.016	10.38414	6.148861	1890.880807

Next steps: [Generate code with Copilot](#) [View recommended plots](#) [New interactive plot](#)

```
# Create a new 'country' column using the values from 'country_x' and 'country_y'
X_train['country'] = X_train['country_x'].combine_first(X_train['country_y'])
```

```
# Drop the old 'country_x' and 'country_y' columns
X_train.drop(columns=['country_x', 'country_y'], inplace=True)
```

```
# Verify the result
print(X_train.columns)
```

```
# Select the columns: "region", "gdp_per_capita", "social_support", "healthy_life_expectancy",
"freedom_to_make_life_choices", "generosity",
"perceptions_of_corruption", "population",
"population_below_poverty_line", "hdi", "life_expectancy",
"expected_years_of_education", "mean_years_of_education", "gdi",
"country"
df.select_dtypes(include=[object])
```

```
# Create a new 'country' column using the values from 'country_x' and 'country_y'
X_test['country'] = X_test['country_x'].combine_first(X_test['country_y'])
```

```
# Drop the old 'country_x' and 'country_y' columns
X_test.drop(columns=['country_x', 'country_y'], inplace=True)
```

```
# Verify the result
print(X_test.columns)
```

```
# Select the columns: "region", "gdp_per_capita", "social_support", "healthy_life_expectancy",
"freedom_to_make_life_choices", "generosity",
"perceptions_of_corruption", "population",
"population_below_poverty_line", "hdi", "life_expectancy",
"expected_years_of_education", "mean_years_of_education", "gdi",
"country"
df.select_dtypes(include=[object])
```

1. EDA

```
print(X_train.dtypes)
```

	region	gdp_per_capita	social_support	healthy_life_expectancy	freedom_to_make_life_choices	generosity	perceptions_of_corruption	population	population_below_poverty_line	hdi	life_expectancy	expected_years_of_education	mean_years_of_education	gdi	country
0	region	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
1	gdp_per_capita	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
2	social_support	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
3	healthy_life_expectancy	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
4	freedom_to_make_life_choices	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
5	generosity	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
6	perceptions_of_corruption	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
7	population	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
8	population_below_poverty_line	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
9	hdi	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
10	life_expectancy	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
11	expected_years_of_education	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
12	mean_years_of_education	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
13	gdi	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	object
14	country	object	object	object	object	object	object	object	object	object	object	object	object	object	object

Describe what you see above?

Answer: There are duplicate Country Columns: country_x (object), country_y (object), country_x and country_y represent the same country but from different datasets, which means that we would need to merge or delete one of them if unnecessary. The X_train dataset contains both categorical (object type) and numerical (float64 and int64 type) features. Categorical columns include: country, region, and region, which may be useful for regional analysis or as model features. Numerical columns consist of: continuous variables like gdp_per_capita, social_support, healthy_life_expectancy, and freedom_to_make_life_choices, as well as the integer population columns. This mix of data types is suitable for exploratory data analysis (EDA) and machine learning, with numerical features being more straightforward and categorical features potentially requiring encoding.

Find out the number and percentage of missing values in the table per column

```
# Your code here:
```

```
# Calculate the number of missing values per column
missing_values_count = X_train.isnull().sum()
```

```
# Calculate the percentage of missing values per column
missing_values_percentage = missing_values_count / len(X_train) * 100
```

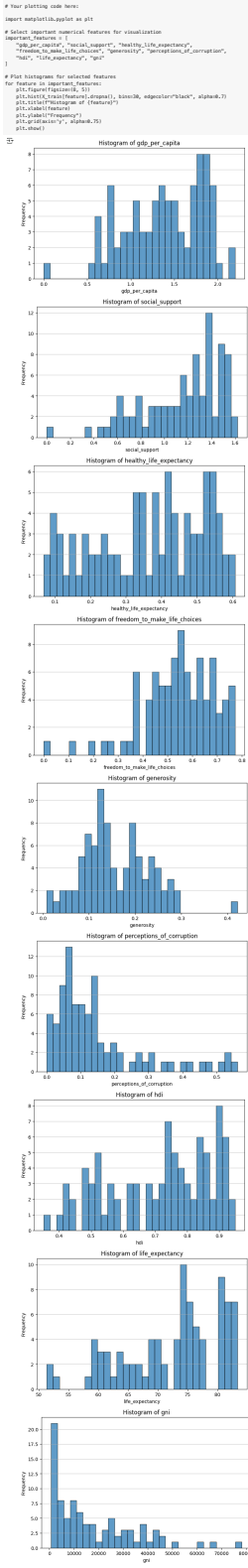
```
# Create a DataFrame to display the results
missing_values_df = pd.DataFrame({
    'Missing Values Count': missing_values_count,
    'Missing Values Percentage': missing_values_percentage
})
```

```
print(missing_values_pct)

region
gdp_per_capita
edu_expert
healthy_life_expectancy
freedom_to_make_life_choices
generosity
perceptions_of_corruption
population
population_below_poverty_line
tota_healthy
mean_year_of_achieving
mean_year_of_achieving
gpi
country

Missing Values Count Missing Values Percentage
0 0.00000
0 0.00000
0 0.00000
0 0.00000
0 0.00000
0 0.00000
0 0.00000
0 0.00000
10 11.23000
0 0.00000
0 0.00000
1 1.12300
1 1.12300
0 0.00000
0 0.00000

Plot the frequency distribution / histogram of some of the numerical features that you think are important
```



Plot the categorical variables and their distribution

Your code here:

```
# Checking the available categorical columns in X_train
categorical_columns = X_train.select_dtypes(include=["object"]).columns.tolist()

# Display the categorical columns
print(categorical_columns)

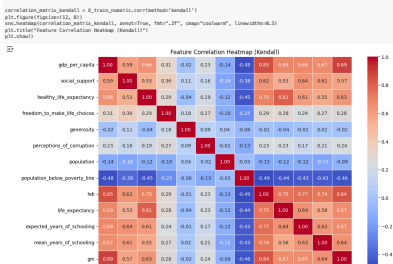
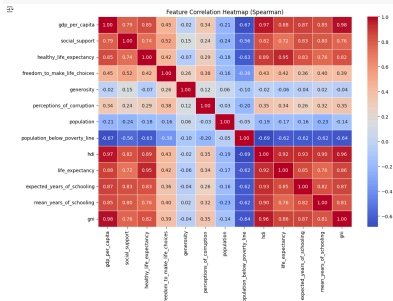
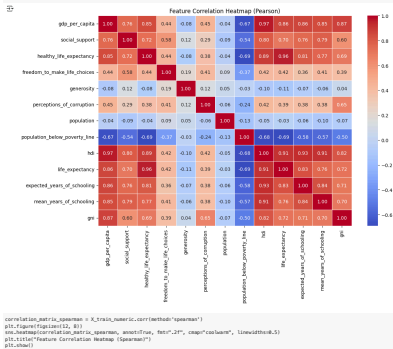
["region", "country"]

import matplotlib.pyplot as plt
import seaborn as sns

# Only "region" variable matters to them
# Select categorical columns
categorical_columns = ["region"]

# Plot the distribution of each categorical variable
for column in categorical_columns:
    plt.figure(figsize=(10, 6))
    sns.countplot(data=X_train, x=column, palette="Set2", order=X_train[column].value_counts().index)
    plt.title(f'Distribution of {column}')
    plt.xlabel(column)
    plt.ylabel('Count')
    plt.grid(True)
    plt.show()
```

- Perform feature correlation analysis to identify relationships between variables. Use Pearson, Spearman, or Kendall correlation coefficients to analyze feature dependencies.



- Explore relationships between variables (bivariate, etc), correlation tables, and how they associate with the target variable.

```

print pandas.io.parsers.read_csv(
    os.path.join('data', 'train.csv'),
    delimiter=';', as_index=False, dtype=object)

#df['train_category'] = df['cat']
df['train_category'] = df['cat']

df['train'] = 'train', 'category', 'right', 'Very Right'

)

# df = df[['df_train_label', 'train_category', 'country', 'region']]
# df = df[['train_category']]

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=2015, stratify=y)

X_train_labels = X_train[['train_category']].cat.codes
X_test_labels = X_test[['train_category']].cat.codes

X_train = X_train.replace([X_train_labels], X_test_labels, right_index=True)
X_test = X_test.replace([X_test_labels], X_test_labels, right_index=True)

X_train = X_train[['train_category', 'left', 'age', 'sex', 'right', 'country', 'new', 'region']].dropna(inplace=True)
X_test = X_test[['train_category', 'left', 'age', 'sex', 'right', 'country', 'new', 'region']].dropna(inplace=True)

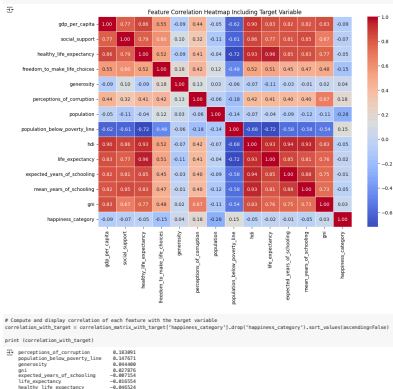
new_numeric_columns = X_train.select_dtypes(include=['number']).columns.tolist()
X_train[new_numeric_columns] = X_train[new_numeric_columns].astype(float)

X_test = X_test[['train_category', 'left', 'age', 'sex', 'right', 'country', 'new', 'region']].dropna(inplace=True)
X_test[new_numeric_columns] = X_test[new_numeric_columns].astype(float)

```

```
# Compute the correlation matrix including the target variable
correlation_matrix_with_target = X_train_numeric.corr()

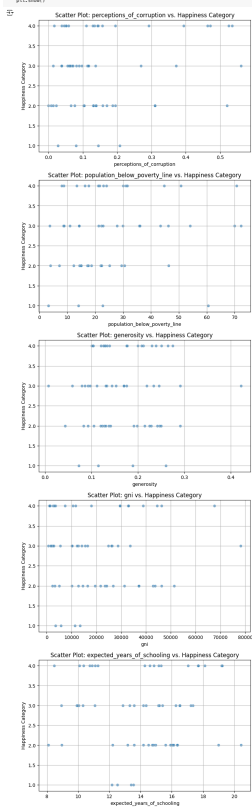
# Plot the heatmap including the target variable
plt.figure(figsize=(12, 8))
sns.heatmap(correlation_matrix_with_target, annot=True, fmt=".2f", cmap="coolwarm", linewidth=0.5)
plt.title("Feature Correlation Heatmap Including Target Variable")
```



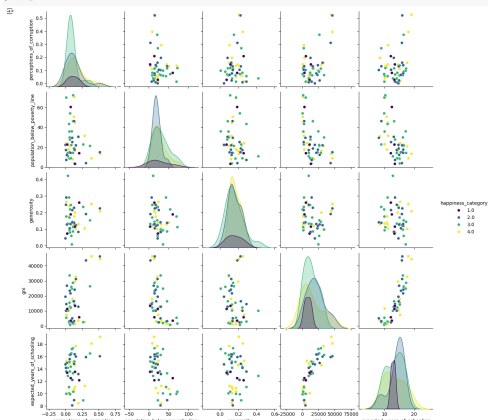
```
mean_years_of_education -0.851502
rel -0.851504
social_support -0.871581
gdp_per_capita -0.892585
freedom_to_move_livelihoods -0.308588
median_household_income -0.203587
happiness_category, etc. (train)
```

```
# Scatter plots to explore relationships between selected features and the target variable
selected_features = correlation_with_target_index[15] # Select top 5 most correlated features
```

```
for feature in selected_features:
    plt.figure(figsize=(10, 10))
    sns.scatterplot(x=feature, y=happiness_category, s=100, alpha=0.5)
    plt.title(f'Scatter Plot: {feature} vs. Happiness Category')
    plt.xlabel(feature)
    plt.ylabel('Happiness Category')
    plt.grid(True)
```



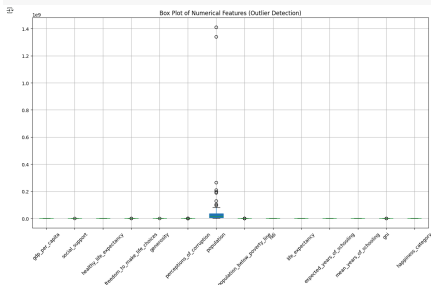
```
numeric_selected_features = [feature for feature in selected_features if feature.dtype in ('float64', 'float32')]
x_train_numeric_scaled = x_train_numeric[numeric_selected_features + ['happiness_category']].drop('happiness_category', axis=1)
one_hot_encoder = OneHotEncoder(sparse_output=False, dtype=int)
x_train_one_hot_scaled = one_hot_encoder.fit_transform(x_train_numeric_scaled)
x_train_scaled = np.concatenate((x_train_numeric_scaled, x_train_one_hot_scaled), axis=1)
```



Also, detect outliers using box plots, Z-score analysis, or the IQR method to identify potential data anomalies.

```
# Your code here:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

# Outlier Detection using Box Plots
plt.figure(figsize=(10, 10))
x_train_numeric_scaled.plot(kind='box', patch_artist=True, vert=True)
plt.title('Box Plot of Numerical Features (Outlier Detection)')
plt.show()
```



```
# Outlier Detection using Z-Score Analysis
z_scores = (x_train_numeric_scaled - x_train_numeric_scaled.mean()) / x_train_numeric_scaled.std()
outliers_zscore = z_scores[z_scores > 3 | z_scores < -3] # Count values beyond 3 standard deviations

# Outlier Detection using the IQR Method
Q1 = x_train_numeric_scaled.quantile(0.25)
Q3 = x_train_numeric_scaled.quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
outliers_iqr = x_train_numeric_scaled[(x_train_numeric_scaled < lower_bound) | (x_train_numeric_scaled > upper_bound)].reset_index()
```

```
# Create a DataFrame to compare outlier counts from both methods
outlier_summary = pd.DataFrame({
    'Z-Score Outliers': outliers_zscore,
    'IQR Outliers': outliers_iqr
})
print(outlier_summary)
```

```
# Z-Score Outliers
gdp_per_capita 1 0
social_support 1 0
freedom_to_move_livelihoods 1 0
median_household_income 1 0
happiness_category 1 0
```

Write what you observed and your General comments on what should be done:

Observations from the Outlier Analysis:

- 1. Box Plot Observations The box plot showed several outliers in different numerical features. Some features, such as GDP per capita, generosity, and social support, exhibited a long tail, suggesting potential outliers. The population variable had extreme values, likely due to large differences in country sizes.
- 2. Z-Score Outlier Analysis Features like population, GDP per capita, and social support had outliers exceeding 3 standard deviations. A few features had isolated extreme values, indicating possible data entry errors or natural variations.
- 3. GDP Outlier Analysis The GDP-related dataset showed more outliers than the Z-score method. Population-related variables and economic indicators had a high number of outliers, as these values tend to be highly skewed.

Comments on Handling Outliers

Outliers could skew statistical analyses, they should be removed to avoid misleading results. This is common in datasets where errors or outliers significantly deviate from the norm. Capping extreme values by replacing them with percentile-based thresholds (e.g., 1st and 99th percentiles) helps reduce sensitivity to extreme values. If data is highly skewed before the population, GDP per capita, and income have natural long tails. Applying a log transformation can normalize the distribution and reduce outlier impact.

2. Feature Engineering

Apply log transformations to normalize skewed data and improve model stability (if any).

Your code here:

import numpy as np
import pandas as pd

visual_features = X_train_visual.drop('iso_alpha3', axis=1)

highly_skewed_features = visual_features[abs(skewness_features) > 3].index.tolist()

X_train_log_transformed = X_train_visual.copy()
X_train_log_transformed[highly_skewed_features] = np.log(X_train_log_transformed[highly_skewed_features])

print(X_train_log_transformed)

	iso_alpha3	iso_alpha3	iso_alpha3
0	6.703	1.114	6.228
1	1.228	1.294	6.488
2	1.269	1.121	6.489
3	1.365	1.259	6.451

	freedom_to_move_life_choices	generosity	perceptions_of_corruption
0	6.137	6.137	6.13731
1	6.138	6.138	6.13731
2	6.139	6.139	6.13731
3	6.140	6.140	6.13731
4	6.141	6.141	6.13731
5	6.142	6.142	6.13731
6	6.143	6.143	6.13731
7	6.144	6.144	6.13731
8	6.145	6.145	6.13731
9	6.146	6.146	6.13731
10	6.147	6.147	6.13731
11	6.148	6.148	6.13731
12	6.149	6.149	6.13731
13	6.150	6.150	6.13731
14	6.151	6.151	6.13731
15	6.152	6.152	6.13731
16	6.153	6.153	6.13731
17	6.154	6.154	6.13731
18	6.155	6.155	6.13731
19	6.156	6.156	6.13731

	population	population_gdp_per_capita	life_expectancy
0	10.00000	1.00000	10.00000
1	10.00000	1.00000	10.00000
2	10.00000	1.00000	10.00000
3	10.00000	1.00000	10.00000
4	10.00000	1.00000	10.00000
5	10.00000	1.00000	10.00000
6	10.00000	1.00000	10.00000
7	10.00000	1.00000	10.00000
8	10.00000	1.00000	10.00000
9	10.00000	1.00000	10.00000
10	10.00000	1.00000	10.00000
11	10.00000	1.00000	10.00000
12	10.00000	1.00000	10.00000
13	10.00000	1.00000	10.00000
14	10.00000	1.00000	10.00000
15	10.00000	1.00000	10.00000
16	10.00000	1.00000	10.00000
17	10.00000	1.00000	10.00000
18	10.00000	1.00000	10.00000
19	10.00000	1.00000	10.00000

	expected_years_of_schooling	mean_years_of_schooling	gdp
0	10.00000	1.00000	10.00000
1	10.00000	1.00000	10.00000
2	10.00000	1.00000	10.00000
3	10.00000	1.00000	10.00000
4	10.00000	1.00000	10.00000
5	10.00000	1.00000	10.00000
6	10.00000	1.00000	10.00000
7	10.00000	1.00000	10.00000
8	10.00000	1.00000	10.00000
9	10.00000	1.00000	10.00000
10	10.00000	1.00000	10.00000
11	10.00000	1.00000	10.00000
12	10.00000	1.00000	10.00000
13	10.00000	1.00000	10.00000
14	10.00000	1.00000	10.00000
15	10.00000	1.00000	10.00000
16	10.00000	1.00000	10.00000
17	10.00000	1.00000	10.00000
18	10.00000	1.00000	10.00000
19	10.00000	1.00000	10.00000

	population_gdp_per_capita	life_expectancy
0	10.00000	1.00000
1	10.00000	1.00000
2	10.00000	1.00000
3	10.00000	1.00000
4	10.00000	1.00000
5	10.00000	1.00000
6	10.00000	1.00000
7	10.00000	1.00000
8	10.00000	1.00000
9	10.00000	1.00000
10	10.00000	1.00000
11	10.00000	1.00000
12	10.00000	1.00000
13	10.00000	1.00000
14	10.00000	1.00000
15	10.00000	1.00000
16	10.00000	1.00000
17	10.00000	1.00000
18	10.00000	1.00000
19	10.00000	1.00000

row = 14 columns

Create at least one interaction feature to capture relationship between existing variables, enhancing predictive power.

Your code here:

import pandas as pd

Creating an Interaction Feature: GDP per Capita * Social Support

X_train_interaction = X_train_log_transformed.copy()
X_train_interaction['gdp_social_support'] = X_train_interaction['gdp_per_capita'] * X_train_interaction['social_support']

print(X_train_interaction)

	iso_alpha3	iso_alpha3	iso_alpha3
0	6.703	1.114	6.228
1	1.228	1.294	6.488
2	1.269	1.121	6.489
3	1.365	1.259	6.451

	freedom_to_move_life_choices	generosity	perceptions_of_corruption
0	6.137	6.137	6.13731
1	6.138	6.138	6.13731
2	6.139	6.139	6.13731
3	6.140	6.140	6.13731
4	6.141	6.141	6.13731
5	6.142	6.142	6.13731
6	6.143	6.143	6.13731
7	6.144	6.144	6.13731
8	6.145	6.145	6.13731
9	6.146	6.146	6.13731
10	6.147	6.147	6.13731
11	6.148	6.148	6.13731
12	6.149	6.149	6.13731
13	6.150	6.150	6.13731
14	6.151	6.151	6.13731
15	6.152	6.152	6.13731
16	6.153	6.153	6.13731
17	6.154	6.154	6.13731
18	6.155	6.155	6.13731
19	6.156	6.156	6.13731

	population	population_gdp_per_capita	life_expectancy
0	10.00000	1.00000	10.00000
1	10.00000	1.00000	10.00000
2	10.00000	1.00000	10.00000
3	10.00000	1.00000	10.00000
4	10.00000	1.00000	10.00000
5	10.00000	1.00000	10.00000
6	10.00000	1.00000	10.00000
7	10.00000	1.00000	10.00000
8	10.00000	1.00000	10.00000
9	10.00000	1.00000	10.00000
10	10.00000	1.00000	10.00000
11	10.00000	1.00000	10.00000
12	10.00000	1.00000	10.00000
13	10.00000	1.00000	10.00000
14	10.00000	1.00000	10.00000
15	10.00000	1.00000	10.00000
16	10.00000	1.00000	10.00000
17	10.00000	1.00000	10.00000
18	10.00000	1.00000	10.00000
19	10.00000	1.00000	10.00000

	expected_years_of_schooling	mean_years_of_schooling	gdp
0	10.00000	1.00000	10.00000
1	10.00000	1.00000	10.00000
2	10.00000	1.00000	10.00000
3	10.00000	1.00000	10.00000
4	10.00000	1.00000	10.00000
5	10.00000	1.00000	10.00000
6	10.00000	1.00000	10.00000
7	10.00000	1.00000	10.00000
8	10.00000	1.00000	10.00000
9	10.00000	1.00000	10.00000
10	10.00000	1.00000	10.00000
11	10.00000	1.00000	10.00000
12	10.00000	1.00000	10.00000
13	10.00000	1.00000	10.00000
14	10.00000	1.00000	10.00000
15	10.00000	1.00000	10.00000
16	10.00000	1.00000	10.00000
17	10.00000	1.00000	10.00000
18	10.00000	1.00000	10.00000
19	10.00000	1.00000	10.00000

	population_gdp_per_capita	life_expectancy
0	10.00000	1.00000
1	10.00000	1.00000
2	10.00000	1.00000
3	10.00000	1.00000
4	10.00000	1.00000
5	10.00000	1.00000
6	10.00000	1.00000
7	10.00000	1.00000
8	10.00000	1.00000
9	10.00000	1.00000
10	10.00000	1.00000
11	10.00000	1.00000
12	10.00000	1.00000
13	10.00000	1.00000
14	10.00000	1.00000
15	10.00000	1.00000
16	10.00000	1.00000
17	10.00000	1.00000
18	10.00000	1.00000
19	10.00000	1.00000

row = 15 columns

3. Preprocess data using sklearn Column Transformer/ Write and Save Preprocessor function

```
from sklearn.preprocessing import StandardScaler, OneHotEncoder  
from sklearn.compose import ColumnTransformer  
from sklearn.pipeline import Pipeline  
from sklearn.impute import SimpleImputer  
numeric_features = X_train.select_dtypes(include='number').columns.tolist()  
numeric_transformer = Pipeline(steps=[  
    ('imputer', SimpleImputer(strategy='constant', fill_value=0)),  
    ('scaler', StandardScaler())  
categorical_features = list(set(X_train.columns) - set(numeric_features))  
categorical_transformer = Pipeline(steps=[  
    ('imputer', SimpleImputer(strategy='most_frequent')),  
    ('onehot', OneHotEncoder(handle_unknown='ignore'))  
preprocessor = ColumnTransformer(transformers=[  
    ('num', numeric_transformer, numeric_features),  
    ('cat', categorical_transformer, categorical_features)  
X_train_preprocessed = preprocessor.fit_transform(X_train)
```

Describe step-by-step what we are doing above, and why? You see how to change how values are imputed. What change did you make if any, and why?

Answer: The data preprocessing pipeline prepares numerical and categorical features in X_train for machine learning. First, numerical features are identified using the select_dtypes method. The numeric transformer pipeline includes imputation and scaling using z-score. Missing values are filled with 0 using SimpleImputer(strategy='constant', fill_value=0), and StandardScaler standardizes features to have a mean of 0 and variance of 1. An alternative approach could be using the median strategy instead for imputation, which is more robust to outliers. For categorical features, the transformer is applied only to relevant columns using fit_transform. It uses SimpleImputer(strategy='most_frequent') to replace missing values with the most common category, followed by OneHotEncoder(handle_unknown='ignore') to convert categories into a binary matrix. Alternatively, missing values could be filled with a placeholder like Unknown for robustness. The ColumnTransformer combines both transformers, applying specific preprocessing to numerical and categorical columns. If categorical features are present, only the numeric transformation is used. The fit method initializes the preprocessor by learning the required statistics from X_train, ensuring consistent data transformation. The pipeline enhances data quality, reduces bias, and prepares well-structured input for machine learning models.

Change the impute strategy of filling missing numerical values with 0 (SimpleImputer(strategy='constant', fill_value=0)) would change the imputer strategy to 'most_frequent'.

```
# Write a function to transform data using the fitted preprocessor  
def transform_preprocessed(  
    data = data.drop('country', axis=1),  
    preprocessed_data = preprocessed_data  
):  
    return preprocessed_data
```

What are the differences between the 'preprocessor' object, the 'preprocessor' object, the 'preprocessor' function, and the 'preprocessed_data' that is returned from it?
Answer: The 'preprocessor' object is a tool that sets up how to handle different types of data, like scaling numbers or turning categories into numbers using ColumnTransformer. The 'preprocessor' object is the same tool but trained on the data, meaning it has learned what values to use for filling missing data or how to scale numbers properly. The 'preprocessor' function is a custom function that uses the preprocessor object to transform new data in the same way as the training data. This makes sure the model gets data in the same format every time, whether it's for training, testing, or making predictions. Finally, the 'preprocessed_data' is what comes out of the function - a cleaned-up version of the data, ready for the model to use. It's a numerical dataset where all the features are standardized, filled in, and encoded correctly.

```
# Check the shape of X_train after preprocessing using our new function  
print(transform_preprocessed(X_train))  
## (18, 13)
```

4. Fit model on preprocessed data and save preprocessor function and model

```
import numpy as np  
import pandas as pd  
from sklearn.metrics import mean_squared_error  
from sklearn.metrics import r2_score  
from sklearn.preprocessing import StandardScaler, OneHotEncoder  
from sklearn.compose import ColumnTransformer  
from sklearn.impute import SimpleImputer  
from sklearn.pipeline import Pipeline  
from sklearn.linear_model import LinearRegression  
X_train_preprocessed = preprocessor.fit_transform(X_train)  
y_train_preprocessed = y_train  
X_test_preprocessed = preprocessor.transform(X_test)  
y_test_preprocessed = y_test  
model = LinearRegression()  
model.fit(X_train_preprocessed, y_train_preprocessed)  
y_pred = model.predict(X_test_preprocessed)  
accuracy = accuracy_score(y_test_preprocessed, y_pred)  
print(accuracy)  
## 0.8000000000000001
```

5. Generate predictions from X_test data and compare it with true labels in Y_test

```
from sklearn.metrics import mean_squared_error, classification_report, confusion_matrix  
# Generate predicted values (Model 1)  
prediction_labels = model.predict(X_test_preprocessed)  
# Data model predictions for comparing predicted labels with true labels  
accuracy = accuracy_score(y_test_preprocessed, prediction_labels)  
print('Model 1 Accuracy: %.2f' % accuracy)  
print('Classification Report:')
```

```
print(classification_report(y_test, prediction_labels))

print("Confusion Matrix")
print(confusion_matrix(y_test, prediction_labels))
```

```
Model 1 Accuracy: 0.8023
```

```
Classification Report:
      precision    recall  f1-score   support

     0:       0.80      1.00      0.89      8
     1:       0.75      0.58      0.66      9
     2:       0.75      0.58      0.67      10

 accuracy:       0.69      0.70      0.69      26
 macro avg:       0.75      0.73      0.73      26
 weighted avg:       0.69      0.65      0.65      26
```

```
Confusion Matrix:
[[ 8  0  0]
 [ 0 5  4]
 [ 0 4  2]]
```

6. Repeat the process with different parameters to improve the accuracy

```
# Train model 2 using same hyperparameters (note that you could use a new hyperparameter, but we will use the same one for this example).
model = RandomForestClassifier(n_estimators=500, max_depth=10, class_weight='balanced', random_state=0)
model.fit(X_train_preprocessed, y_train)
prediction_labels = model.predict(X_test_preprocessed)
accuracy_2 = accuracy_score(y_test, prediction_labels)
print("Model 2 Accuracy: (accuracy_2_457)")
```

```
Model 2 Accuracy: 0.7902
```

What changes did you make, what do the parameters you changed control, and why does it improve performance?

Answer: For Model 2 I increased n_estimators from 100 to 500, increased max_depth from 10 to 20.

```
prediction_labels = model.predict(X_test_preprocessed)
```

```
# Show model performance by comparing prediction_labels with true labels
accuracy_2 = accuracy_score(y_test, prediction_labels)
print("Model 2 Accuracy: (accuracy_2_457)")
```

```
print("Classification Report: Model 2(1)")
print(classification_report(y_test, prediction_labels))
```

```
print("Confusion Matrix: Model 2(1)")
print(confusion_matrix(y_test, prediction_labels))
```

```
Model 2 Accuracy: 0.7902
```

```
Classification Report: Model 2(1)
      precision    recall  f1-score   support

     0:       0.80      1.00      0.89      8
     1:       0.75      0.74      0.74      10
     2:       0.75      0.70      0.73      10

 accuracy:       0.77      0.79      0.77      26
 macro avg:       0.77      0.78      0.78      26
 weighted avg:       0.77      0.77      0.76      26
```

```
Confusion Matrix: Model 2(1)
[[ 8  0  0]
 [ 0 5  4]
 [ 0 4  2]]
```

Do you think it is worth making more changes to the parameters? Should we keep trying random values and see what works better? What is an alternative to doing this manually?

Answer: Manually changing model parameters and testing random values can sometimes improve accuracy, but GridSearchCV is a more efficient approach. Instead of guessing, GridSearchCV systematically tests all possible parameter combinations and finds the best model using cross-validation to avoid overfitting. While it can be computationally heavy, it is much more effective than trying random values. A faster alternative is RandomizedSearchCV, which samples random parameter combinations and is useful for large parameter spaces. Both methods provide a systematic, more strategic approach to hyperparameter tuning.

```
from sklearn.metrics import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import numpy as np

param_grid = {
    'n_estimators': [50, 100],
    'max_depth': [10, 20],
    'min_samples_split': [2, 5],
    'min_samples_leaf': [1, 2],
    'max_features': ['sqrt', 'log2'],
    'class_weight': ['balanced']
}
```

```
gridsearch = GridSearchCV(
    estimator=RandomForestClassifier(random_state=0),
    param_grid=param_grid,
    cv=5,
    scoring='accuracy',
    n_jobs=-1,
    verbose=1)

```

```
gridsearch.fit(X_train_preprocessed, y_train)
```

```
prediction_labels_3 = gridsearch.predict(X_test_preprocessed)
```

#Model: Model 3:

```
#-- Generate predicted values
# We'll add to the model's performance by comparing prediction_labels with true labels
accuracy_3 = accuracy_score(y_test, prediction_labels_3)
print("Model 3 Accuracy: (accuracy_3_457)")
```

```
print("Classification Report: Model 3(1)")
print(classification_report(y_test, prediction_labels_3))
```

```
print("Confusion Matrix: Model 3(1)")
print(confusion_matrix(y_test, prediction_labels_3))
```

```
Model 3 Accuracy: 0.8077
```

```
Classification Report: Model 3(1)
      precision    recall  f1-score   support

     0:       0.80      1.00      0.89      8
     1:       0.80      0.75      0.78      9
     2:       0.80      0.75      0.78      10

 accuracy:       0.82      0.81      0.81      26
 macro avg:       0.80      0.81      0.81      26
 weighted avg:       0.81      0.81      0.81      26
```

```
Confusion Matrix: Model 3(1)
[[ 8  0  0]
 [ 0 5  4]
 [ 0 4  2]]
```

```
from sklearn.metrics import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import numpy as np

# Define the models you want to evaluate
models = [
    # Random Forest
    RandomForestClassifier(n_estimators=100,
                          max_depth=10,
                          min_samples_split=2,
                          min_samples_leaf=1,
                          max_features='sqrt',
                          random_state=0),
    # Gradient Boosting
    GradientBoostingClassifier(n_estimators=100,
                              learning_rate=0.1,
                              random_state=0)
]
```

```
# Iterate through each model, fit it, predict, and display results
for model_name, model in models.items():
    model.fit(X_train_preprocessed, y_train)
    prediction_labels = model.predict(X_test_preprocessed)
```

```
# Calculate accuracy and show evaluation metrics
accuracy = accuracy_score(y_test, prediction_labels)
print("Accuracy: (accuracy_457)")
print("Accuracy: (accuracy_457)")
print("Classification Report:")
print(classification_report(y_test, prediction_labels))
print("Confusion Matrix:")
print(confusion_matrix(y_test, prediction_labels))
```

```
Accuracy: 0.7988
```

```
Classification Report:
      precision    recall  f1-score   support

     0:       0.80      1.00      0.89      8
     1:       0.75      0.75      0.75      9
     2:       0.80      0.58      0.67      10

 accuracy:       0.73      0.75      0.73      26
 macro avg:       0.78      0.75      0.75      26
 weighted avg:       0.78      0.73      0.73      26
```

```
Confusion Matrix:
[[ 8  0  0]
 [ 0 5  4]
 [ 0 4  2]]
```

```
Model: Bagging Classifier
```

```
Accuracy: 0.8077
```

```
Classification Report:
      precision    recall  f1-score   support

     0:       0.80      1.00      0.89      8
     1:       0.80      0.80      0.80      9
     2:       0.80      0.68      0.74      10

 accuracy:       0.81      0.81      0.81      26
 macro avg:       0.81      0.81      0.81      26
 weighted avg:       0.81      0.81      0.80      26
```

```
Confusion Matrix:
[[ 8  0  0]
 [ 0 5  4]
 [ 0 4  2]]
```

```
Model: Gradient Boosting
```

```
Accuracy: 0.7902
```

```
Classification Report:
      precision    recall  f1-score   support

     0:       0.80      1.00      0.89      8
     1:       0.80      0.80      0.80      9
     2:       0.80      0.68      0.74      10

 accuracy:       0.80      0.79      0.77      26
 macro avg:       0.80      0.79      0.77      26
 weighted avg:       0.80      0.77      0.77      26
```

```
Confusion Matrix:
[[ 8  0  0]
 [ 0 5  4]
 [ 0 4  2]]
```

Describe what were the parameters you defined in GradientBoostingClassifier, and/or BaggingClassifier, and/or KNNs, and/or SVC? What worked and why?

I used several KNN, Bagging classifier, and the gradient boosting classifier.

The K-Nearest Neighbors (KNN) model was configured with n_neighbors=1, meaning it used the five nearest data points to classify each sample. The model achieved an accuracy of 73.08%. KNN's performance depends heavily on local data distribution, and it may not always perform well in cases of class imbalance or complex data patterns.

The Bagging Classifier was the top-performing model, reaching an accuracy of 80.77%. It used 50 base estimators (Random Forests) and trained each estimator on random subsets of data, which helped reduce overfitting and improved generalization. This model works by combining the predictions of multiple weak learners, which reduces variance and increases stability, making it particularly effective in handling diverse and noisy datasets.

The Gradient Boosting Classifier also performed well, achieving an accuracy of 79.02%. It used 100 boosting stages (n_estimators=100) with a learning rate of 0.1, which controlled how much each tree corrected the errors of the previous one. Gradient Boosting is effective in learning complex patterns by building trees sequentially. However, its performance was slightly lower than the Bagging Classifier, indicating it might need further parameter tuning or adjustments to the learning rate to outperform.

Overall, the Bagging Classifier exhibited the best performance, demonstrating robust performance by combining predictions from multiple estimators to reduce overfitting.

7. Basic Deep Learning

```
import keras
from keras.models import Sequential
from keras.layers import Dense, Activation
from keras.optimizers import Adam

# Determine the number of input features and output classes
features_count = preprocessor.transform(X_train).shape[1]
# The number of output classes is the target variable
output_classes = len(y_train_unique)

# Convert X_train to one-hot-encoded format
y_train_encoded = to_categorical(y_train)

# Define a Neural Network Model with 3 layers: 128 -> 64 -> 10 -> output_classes
model = Sequential()
model.add(Dense(128, activation='relu'))
model.add(Dense(64, activation='relu'))
model.add(Dense(10, activation='softmax')) # Output layer with softmax activation
```

```
# Compile the model
model.compile(loss='categorical_crossentropy',
              optimizer='adam',
              metrics=['accuracy'])
```

```
# Fit the model to the training data
history = model.fit(preprocessor.transform(X_train), y_train_encoded,
                    batch_size=32,
                    validation_split=0.2)
```


8/12


```
from sklearn.metrics import accuracy_score, loss

# Train and validate the model
train_loss, train_acc = train_model(train_data_loader, dev_data_loader, patience=10)
dev_loss, dev_acc = validate_model(dev_data_loader)

# Print results
print(f"Model Accuracy: {train_acc:.4f} - Loss: {train_loss:.4f} | Dev Accuracy: {dev_acc:.4f} - Loss: {dev_loss:.4f}")

# Plot training and validation accuracy and loss
plt.figure(figsize=(10, 5))

# Accuracy plot
plt.plot(train_acc, label='Training Accuracy')
plt.plot(dev_acc, label='Validation Accuracy')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend()

# Loss plot
plt.plot(train_loss, label='Training Loss')
plt.plot(dev_loss, label='Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
```

Model Accuracy with sigmoid activation: 0.8877

	precision	recall	f1-score	support
0	0.87	1.00	0.93	8
1	0.83	0.82	0.82	10
2	1.00	0.80	0.89	26

Confusion Matrix
[[8 0]
 [0 10]
 [0 26]]

Model Accuracy with sigmoid Activation

Model Loss with sigmoid Activation

8. Explainability - SHAP Feature Importance

To better understand our model's predictions, we will use SHAP (Shapley Additive Feature Importance) to analyze feature importance.

How SHAP Works?

- SHAP assigns each feature a contribution score for every prediction.
- Use **shap.summary_plot** to visualize and interpret our model's feature contributions.

We will now apply SHAP to visualize and interpret our model's feature contributions.

```
import shap
import matplotlib.pyplot as plt

# Initialize SHAP explainer using the model and preprocessed data
explainer = shap.KernelExplainer(model.predict, preprocessor.transform(train_data))

# Generate SHAP values for the test data (1000 samples to avoid performance issues)
shap_values = shap explainer.predict_shap(preprocessor.transform(test_data))

# Plot SHAP summary plot
shap.summary_plot(shap_values, test_data)
```

SHAP summary plot with correct feature names

The sets of the plot represents the SHAP interaction values, indicating the impact each feature has on the prediction. The wide spread of points for `age_m`, `capital` suggests a effects predictions to varying degrees, while `social_support` and `healthy_life_expectancy` tend to cluster more around the center, showing a more consistent influence.

The plot also highlights how individual data points are affected by each feature. This is particularly visible with `age_m`, `capital`, which demonstrates that the feature can both increase and decrease the predicted outcome based on its value. In contrast, `social_support` and `healthy_life_expectancy` exhibit a more balanced impact, with most contributions being neutral or slightly positive.

Experimentation

- Using more models to predict

Naive Bayes (GaussianNB)

```
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Split the dataset into training and testing sets
X_train_data = X_train.drop('label', axis=1)
y_train_data = y_train

# Initialize and train the Naive Bayes model
model_nb = GaussianNB()
model_nb.fit(X_train_data, y_train_data)

# Predict and evaluate the model
predictions_nb = model_nb.predict(X_test_data)

print("Model: Naive Bayes")
print("Accuracy: %.2f" % accuracy_score(y_test_data, predictions_nb))
print("Classification Report:")
print(classification_report(y_test_data, predictions_nb))
print("Confusion Matrix:")
print(confusion_matrix(y_test_data, predictions_nb))
```

```
Model: Naive Bayes
Accuracy: 0.802

Classification Report:
              precision    recall  f1-score   support

     0               0.71         0.80         0.76         8
     1               0.29         0.20         0.24         8
     2               0.40         0.50         0.45         10

 accuracy          0.64         0.51         0.58         26
 macro avg         0.47         0.50         0.48         26
 weighted avg      0.54         0.50         0.52         26

Confusion Matrix:
[[ 8  0  0]
 [ 0  8  0]
 [ 0  0  10]]
```

Support Vector Machine (SVM)

```
from sklearn.svm import SVC

# Initialize and train the SVM model
model_svm = SVC(kernel='linear', class_weight='balanced', random_state=42)
model_svm.fit(X_train_data, y_train_data)

# Predict and evaluate the model
predictions_svm = model_svm.predict(X_test_data)

print("Model: Support Vector Machine")
print("Accuracy: %.2f" % accuracy_score(y_test_data, predictions_svm))
print("Classification Report:")
print(classification_report(y_test_data, predictions_svm))
print("Confusion Matrix:")
print(confusion_matrix(y_test_data, predictions_svm))
```

```
Model: Support Vector Machine
Accuracy: 0.788

Classification Report:
              precision    recall  f1-score   support

     0               0.71         0.80         0.76         8
     1               0.27         0.20         0.23         8
     2               0.48         0.50         0.49         10

 accuracy          0.72         0.72         0.72         26
 macro avg         0.74         0.72         0.72         26
 weighted avg      0.74         0.72         0.72         26

Confusion Matrix:
[[ 8  0  0]
 [ 0  8  0]
 [ 0  0  10]]
```

AdaBoost Classifier

```
from sklearn.ensemble import AdaBoostClassifier

# Initialize and train the AdaBoost model
model_ada = AdaBoostClassifier(estimator=SVC(kernel='linear', random_state=42))
model_ada.fit(X_train_data, y_train_data)

# Predict and evaluate the model
predictions_ada = model_ada.predict(X_test_data)

print("Model: AdaBoost Classifier")
print("Accuracy: %.2f" % accuracy_score(y_test_data, predictions_ada))
print("Classification Report:")
print(classification_report(y_test_data, predictions_ada))
print("Confusion Matrix:")
print(confusion_matrix(y_test_data, predictions_ada))
```

```
Model: AdaBoost Classifier
Accuracy: 0.832

Classification Report:
              precision    recall  f1-score   support

     0               0.73         0.80         0.76         8
     1               0.28         0.20         0.24         8
     2               0.48         0.50         0.49         10

 accuracy          0.74         0.71         0.73         26
 macro avg         0.74         0.70         0.72         26
 weighted avg      0.74         0.69         0.69         26

Confusion Matrix:
[[ 8  0  0]
 [ 0  8  0]
 [ 0  0  10]]
```

LightGBM Classifier

```
# Import necessary libraries
import lightgbm as lgb
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Suppress FutureWarnings and LightGBM logging
warnings.filterwarnings('ignore')
logging.getLogger('lightgbm').setLevel(logging.ERROR)

# Initialize and train the LightGBM model with suppressed output
model_lgb = lgb.LGBMClassifier()
model_lgb.fit(X_train_data, y_train_data)

# Predict and evaluate the model
predictions_lgb = model_lgb.predict(X_test_data)

print("Model: LightGBM Classifier")
print("Accuracy: %.2f" % accuracy_score(y_test_data, predictions_lgb))
print("Classification Report:")
print(classification_report(y_test_data, predictions_lgb))
print("Confusion Matrix:")
print(confusion_matrix(y_test_data, predictions_lgb))
```

```
Model: LightGBM Classifier
Accuracy: 0.792

Classification Report:
              precision    recall  f1-score   support

     0               0.69         0.80         0.74         8
     1               0.31         0.20         0.25         8
     2               0.50         0.50         0.50         10

 accuracy          0.70         0.70         0.70         26
 macro avg         0.70         0.70         0.70         26
 weighted avg      0.70         0.70         0.70         26

Confusion Matrix:
[[ 8  0  0]
 [ 0  8  0]
 [ 0  0  10]]
```

Interpretation

Why Does LightGBM, Naive Bayes, SVM, and AdaBoost Outperform LightGBM (Light Gradient Boosting Machine) is a powerful boosting

algorithm that builds trees sequentially to reduce prediction errors of previous trees. It is particularly effective with structured data and can

handle large datasets efficiently.

Naive Bayes: A simple probabilistic classifier based on Bayes' theorem. Works well with small datasets and assumes independence among

features. Useful as a baseline model because of its simplicity and speed.

Support Vector Machines (SVM): SVM aims to find the optimal hyperplane that maximizes the margin between classes. It is effective for

binary and multiclass classification problems and works well with both linear and non-linear data. Provides a clear boundary decision-making

process, useful for interpretability.

AdaBoost (Adaptive Boosting): Another boosting technique that combines weak classifiers to create a strong classifier. Focuses on difficult-to-classify instances by giving them higher weights in subsequent classifiers. Enhances interpretability by analyzing how each weak learner

contributes to the final decision.

Integration of Model Results: Naive Bayes Model Accuracy: 0.802, Confusion Matrix: Shows a strong performance in classifying category

0 but struggles with class 1. Two Naive Bayes models feature independence, this might not align well with the dataset where features may

interact, causing lower performance for class 1.

Support Vector Machine (SVM): Accuracy: 0.788, Confusion Matrix: Performs well in classifying categories 0 and 2 but has mixed

performance in category 1. SVM's use of hyperplanes works well, but it may not fully capture complex relationships between features,

impacting its classification of category 1.

AdaBoost Model Accuracy: 0.832, Confusion Matrix: Good precision for class 0, but lower performance for classes 1 and 2. The model

effectively focuses on difficult cases but may require more base learners or parameter tuning to improve classification stability.

LightGBM Model Accuracy: 0.792, Confusion Matrix: Shows strong classification performance for class 0 and decent performance for

classes 1 and 2, LightGBM's ability to handle categorical features and large datasets efficiently likely contributed to its strong performance. It

balances precision and recall well across all classes.

Among the models, LightGBM provides the best balance of performance and explainability. Its built-in feature importance scores allow us to

understand which features contribute the most to the predictions. SVM also offers some interpretability through support vectors, although the

data points most influential in defining class boundaries.

You are encouraged to try more experimentation and use other models by adding more code cells to this notebook.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

We can also try to report any relevant statistics to confusion, metrics, etc., and see if it helps the predictions.

SHap Analysis Interpretation

Across all models, *gibc_per_capita*, *healthy_life_expectancy*, and *social_support* are consistently among the top influential features. The LightGBM model shows a broader range of SHAP values, indicating a more dynamic response to feature changes, particularly with *gibc_per_capita*. In contrast, the Adaboost model exhibits a smaller range of SHAP values, suggesting it is less sensitive to changes in its features compared to LightGBM.

The feature impact varies between models. In both LightGBM and SVM models, *gibc_per_capita* demonstrates a dual influence on predictions, with both positive and negative impacts. This dual behavior shows the feature's versatility in driving different outcomes depending on the context. The clustering of SHAP values near zero for some features in the Naïve Bayes model suggests those features influence or have a lower linear response to those features. On the other hand, LightGBM and SVM models present a wider distribution of SHAP values, highlighting richer feature interaction and model complexity.

Feature	Value
country_Nicaragua	0.00
gdp_per_capita	-0.26
life_expectancy	-0.33
country_Costa Rica	0.00
country_New Zealand	0.00
country_Jamaica	0.00
country_Kyrgyzstan	0.00
country_Slovenia	0.00
country_Romania	0.00
country_United Arab Emirates	0.00

LIME provides a local explanation for individual predictions, complementing the global insights from SHAP. For example, in the LightGBM model, `is_german` and `is_married` played significant roles in a specific prediction. The Naive Bayes model, however, focused primarily on categorical features (countries and regions), indicating it might rely more heavily on categorical groupings rather than nuanced numerical relationships.

LIME's visualizations for Adaboost and SVM models reveal that predictions are influenced by both numerical and categorical features. However, the magnitude of influence in the Adaboost model is generally smaller, which aligns with the narrower range of SHAP values observed earlier.

Overall Evaluation

SHAP offers a global view of feature importance across all predictions, effectively showing how each feature contributes to the model's decisions. This is particularly useful for understanding overall model behavior and ensuring alignment with domain knowledge. LIME, on the other hand, provides a local view for specific predictions, allowing for more precise justifications when explaining individual cases.