

```
In [15]: #Import Python libraries
import pandas as pd
import numpy as np

In [17]: #To ignore warnings
import warnings
warnings.filterwarnings('ignore')

In [18]: #Reading Dataset

In [19]: #Load the dataset
gd = pd.read_csv('gym.csv')
gd.head()
```

Age	Gender	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Duration (hours)	Calories_Burned	Workout_Type	Fat_Percentage	Water_Intake (liters)	Workout_Frequency (days/week)	Experience_Level
0	56	Male	88.3	1.71	180	157	60	1.69	1313.0	Yoga	12.6	3.5	4
1	46	Female	74.9	1.53	179	151	66	1.30	883.0	HIT	32.9	2.1	4
2	32	Female	68.1	1.66	167	122	64	1.11	677.0	Cardio	33.4	2.3	4
3	25	Male	53.2	1.70	190	164	56	0.59	532.0	Strength	28.8	2.1	3
4	38	Male	46.1	1.79	188	158	68	0.64	556.0	Strength	29.2	2.8	3

```
In [17]: pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)

In [18]:

In [19]: gd.shape
Out[19]: (973, 15)

In [19]: gd.head()
```

Age	Gender	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Duration (hours)	Calories_Burned	Workout_Type	Fat_Percentage	Water_Intake (liters)	Workout_Frequency (days/week)	Experience_Level
0	56	Male	88.3	1.71	180	157	60	1.69	1313.0	Yoga	12.6	3.5	4
1	46	Female	74.9	1.53	179	151	66	1.30	883.0	HIT	32.9	2.1	4
2	32	Female	68.1	1.66	167	122	64	1.11	677.0	Cardio	33.4	2.3	4
3	25	Male	53.2	1.70	190	164	56	0.59	532.0	Strength	28.8	2.1	3
4	38	Male	46.1	1.79	188	158	68	0.64	556.0	Strength	29.2	2.8	3

```
In [21]: gd.tail()
```

Age	Gender	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Duration (hours)	Calories_Burned	Workout_Type	Fat_Percentage	Water_Intake (liters)	Workout_Frequency (days/week)	Experience_Level
968	24	Male	87.1	1.74	187	158	67	1.57	1364.0	Strength	10.0	3.5	4
969	25	Male	66.6	1.61	184	166	56	1.38	1260.0	Strength	25.0	3.0	2
970	59	Female	60.4	1.76	194	120	53	1.72	929.0	Cardio	18.8	2.7	5
971	32	Male	126.4	1.83	198	146	62	1.10	883.0	HIT	28.2	2.1	3
972	46	Male	88.7	1.63	166	146	66	0.75	542.0	Strength	28.8	3.5	2

```
In [23]: gd.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 973 entries, 0 to 972
Data columns (total 15 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Age                   973 non-null    int64
 1   Gender                973 non-null    object
 2   Weight (kg)           973 non-null    float64
 3   Height (m)            973 non-null    float64
 4   Max_BPM               973 non-null    int64
 5   Avg_BPM               973 non-null    int64
 6   Resting_BPM           973 non-null    int64
 7   Session_Duration (hours) 973 non-null    float64
 8   Calories_Burned       973 non-null    float64
 9   Workout_Type          973 non-null    object
10   Fat_Percentage        973 non-null    float64
11   Water_Intake (liters) 973 non-null    float64
12   Workout_Frequency (days/week) 973 non-null    int64
13   Experience_Level       973 non-null    int64
14   BMI                   973 non-null    float64
dtypes: float64(11), int64(6), object(2)
memory usage: 114.2+ KB

In [25]: gd.unique()

Age
Gender
Weight (kg)
Height (m)
Max_BPM
Avg_BPM
Resting_BPM
Session_Duration (hours)
Calories_Burned
Workout_Type
Fat_Percentage
Water_Intake (liters)
Workout_Frequency (days/week)
Experience_Level
BMI
dtype: object

In [27]: gd.isnull().sum()

Age
Gender
Weight (kg)
Height (m)
Max_BPM
Avg_BPM
Resting_BPM
Session_Duration (hours)
Calories_Burned
Workout_Type
Fat_Percentage
Water_Intake (liters)
Workout_Frequency (days/week)
Experience_Level
BMI
dtype: int64

In [29]: # Check for missing values
print("Missing values in", gd.isnull().sum())

Missing values:
Age
Gender
Weight (kg)
Height (m)
Max_BPM
Avg_BPM
Resting_BPM
Session_Duration (hours)
Calories_Burned
Workout_Type
Fat_Percentage
Water_Intake (liters)
Workout_Frequency (days/week)
Experience_Level
BMI
dtype: int64

In [33]: # To separate the categorical and numerical variables for easy analysis
cat_cols = gd.select_dtypes(include='object').columns
num_cols = gd.select_dtypes(include='number').columns.tolist()
print("Categorical Variables:")
print(cat_cols)
print("Numerical Variables:")
print(num_cols)

Categorical Variables:
Index(['Gender', 'Workout_Type'], dtype='object')
Numerical Variables:
Index(['Age', 'Weight (kg)', 'Height (m)', 'Max_BPM', 'Avg_BPM', 'Resting_BPM', 'Session_Duration (hours)', 'Calories_Burned', 'Fat_Percentage', 'Water_Intake (liters)', 'Workout_Frequency (days/week)', 'Experience_Level', 'BMI'])

In [35]: # EDA Univariate Analysis

In [39]: # Libraries to plot the charts
import matplotlib.pyplot as plt
import seaborn as sns

In [51]: # Histograms of numerical features
fig=plt.figure(figsize=(12, 8), bins=5, edgecolor='black')
plt.show()
```

```
In [47]: # Boxplot of Calories Burned by Workout Type
plt.figure(figsize=(10, 6))
sns.boxplot(x='Workout_Type', y='Calories_Burned', data=gd)
plt.xticks(rotation=45)

Out[47]: (Text(0, 0, 'Yoga'),
Text(1, 0, 'HIT'),
Text(2, 0, 'Cardio'),
Text(3, 0, 'Strength'))
```

```
In [119]: # Define age bins
age_bins = [0, 18, 25, 35, 45, 55, 65, 100]
age_labels = ['18-25', '26-35', '36-45', '46-55', '56-65', '66-70', '71+']

# Define session duration bins
session_duration_bins = [0, 1, 2, 3, 4, 5, 24] # Adjust based on your data range
session_duration_labels = ['1-30 Min', '30 Min-1 hour', '1-2 hours', '2-3 hours', '3-4 hours', '4-5 hours']

# Create new grouped columns
gd['Age_Group'] = pd.cut(gd['Age'], bins=age_bins, labels=age_labels, right=False)
gd['Session_Duration_Group'] = pd.cut(gd['Session_Duration (hours)'], bins=session_duration_bins, labels=session_duration_labels, right=False)

# Create subplots
fig, axes = plt.subplots(2, 2, figsize=(18, 18))
fig.suptitle('Bar plot for all categorical variables in the dataset', fontsize=16, fontweight='bold')

# Plot for each categorical variable
sns.countplot(ax=axes[0, 0], x='Workout_Type', data=gd, color='blue',
              order=gd['Workout_Type'].value_counts().index)
axes[0, 0].set_title('Workout Type Distribution', fontweight='bold')
axes[0, 0].set_xticklabels(axes[0, 0].get_xticklabels(), fontweight='bold')

sns.countplot(ax=axes[0, 1], x='Age_Group', data=gd, color='blue',
              order=gd['Age_Group'].value_counts().index)
axes[0, 1].set_title('Age Group Distribution', fontweight='bold')
axes[0, 1].set_xticklabels(axes[0, 1].get_xticklabels(), fontweight='bold')

sns.countplot(ax=axes[1, 0], x='Session_Duration_Group', data=gd, color='blue',
              order=gd['Session_Duration_Group'].value_counts().index)
axes[1, 0].set_title('Session Duration Distribution', fontweight='bold')
axes[1, 0].set_xticklabels(axes[1, 0].get_xticklabels(), fontweight='bold')

sns.countplot(ax=axes[1, 1], x='Workout_Frequency (days/week)', data=gd, color='blue',
              order=gd['Workout_Frequency (days/week)'].value_counts().index)
axes[1, 1].set_title('Workout Frequency Distribution', fontweight='bold')
axes[1, 1].set_xticklabels(axes[1, 1].get_xticklabels(), fontweight='bold')

plt.tight_layout()
plt.subplots_adjust(top=0.9, hspace=0.2) # Adjust to prevent title overlap
plt.show()
```

```
In [133]: # Define age bins
age_bins = [0, 18, 25, 35, 45, 55, 65, 100]
age_labels = ['18-25', '26-35', '36-45', '46-55', '56-65', '66-70', '71+']

# Define session duration bins
session_duration_bins = [0, 1, 2, 3, 4, 5, 24] # Adjust based on your data range
session_duration_labels = ['1-30 Min', '30 Min-1 hour', '1-2 hours', '2-3 hours', '3-4 hours', '4-5 hours']

# Define custom colors for each category
colors_workout = ['red', 'yellow', 'green', 'blue', 'purple'] # Adjust colors if needed
colors_age = ['red', 'yellow', 'green', 'blue', 'purple', 'pink', 'green']
colors_session = ['red', 'yellow', 'green', 'blue', 'purple', 'pink']
colors_workout_frequency = ['red', 'yellow', 'green', 'blue', 'purple', 'pink']

# Create new grouped columns
gd['Age_Group'] = pd.cut(gd['Age'], bins=age_bins, labels=age_labels, right=False)
gd['Session_Duration_Group'] = pd.cut(gd['Session_Duration (hours)'], bins=session_duration_bins, labels=session_duration_labels, right=False)

# Create subplots
fig, axes = plt.subplots(2, 2, figsize=(18, 18))
fig.suptitle('Bar plot for all categorical variables in the dataset', fontsize=16, fontweight='bold')

# Plot for each categorical variable with different colors
sns.countplot(ax=axes[0, 0], x='Workout_Type', data=gd, palette=colors_workout,
              order=gd['Workout_Type'].value_counts().index)
axes[0, 0].set_title('Workout Type Distribution', fontweight='bold')
axes[0, 0].set_xticklabels(axes[0, 0].get_xticklabels(), fontweight='bold')

sns.countplot(ax=axes[0, 1], x='Age_Group', data=gd, palette=colors_age,
              order=gd['Age_Group'].value_counts().index)
axes[0, 1].set_title('Age Group Distribution', fontweight='bold')
axes[0, 1].set_xticklabels(axes[0, 1].get_xticklabels(), fontweight='bold')

sns.countplot(ax=axes[1, 0], x='Session_Duration_Group', data=gd, palette=colors_session,
              order=gd['Session_Duration_Group'].value_counts().index)
axes[1, 0].set_title('Session Duration Distribution', fontweight='bold')
axes[1, 0].set_xticklabels(axes[1, 0].get_xticklabels(), fontweight='bold')

sns.countplot(ax=axes[1, 1], x='Workout_Frequency (days/week)', data=gd, palette=colors_workout_frequency,
              order=gd['Workout_Frequency (days/week)'].value_counts().index)
axes[1, 1].set_title('Workout Frequency Distribution', fontweight='bold')
axes[1, 1].set_xticklabels(axes[1, 1].get_xticklabels(), fontweight='bold')

plt.tight_layout()
plt.subplots_adjust(top=0.9, hspace=0.2) # Adjust to prevent title overlap
plt.show()
```

```
In [151]: # EDA Bivariate Analysis

In [180]: # Scatter Plot: Session Duration vs Calories Burned
plt.figure(figsize=(10, 4))
sns.scatterplot(x=gd['Session_Duration (hours)'], y=gd['Calories_Burned'], hue=gd['Gender'])
plt.title('Session Duration vs Calories Burned')
plt.xlabel('Session Duration (hours)')
plt.ylabel('Calories Burned')
plt.show()

# Scatter Plot: Avg BPM vs Calories Burned
plt.figure(figsize=(10, 4))
sns.scatterplot(x=gd['Avg_BPM'], y=gd['Calories_Burned'], hue=gd['Workout_Type'])
plt.title('Avg BPM vs Calories Burned')
plt.xlabel('Average BPM')
plt.ylabel('Calories Burned')
plt.show()
```

```
In [181]: # Box Plots: Workout Type vs Calories Burned
plt.figure(figsize=(5, 3))
sns.boxplot(x=gd['Workout_Type'], y=gd['Calories_Burned'])
plt.xticks(rotation=45)
plt.title('Workout Type vs Calories Burned')
plt.show()

# Box Plots: Gender vs Calories Burned
plt.figure(figsize=(5, 3))
sns.boxplot(x=gd['Gender'], y=gd['Calories_Burned'])
plt.title('Gender vs Calories Burned')
plt.show()
```

```
In [186]: # Create a figure with 3 subplots in a single row
fig, axes = plt.subplots(1, 3, figsize=(18, 5))

# Bar Plot: Gender vs Total Calories Burned
sns.barplot(x=gd['Gender'], y='Calories_Burned', data=gd, estimator='sum', ci=None, ax=axes[0])
axes[0].set_title('Total Calories Burned by Gender')
axes[0].set_xlabel('Total Calories Burned')

# Bar Plot: Workout Type vs Average Calories Burned
sns.barplot(x=gd['Workout_Type'], y='Calories_Burned', data=gd, estimator='mean', ci=None, ax=axes[1])
axes[1].set_title('Avg Calories Burned per Workout Type')
axes[1].set_xlabel('Average Calories Burned')

# Bar Chart: Experience Level vs Average Calories Burned
sns.barplot(x=gd['Experience_Level'], y='Calories_Burned', data=gd, estimator='mean', ci=None, ax=axes[2])
axes[2].set_title('Calories Burned by Experience Level')
axes[2].set_xlabel('Experience Level (1 = Beginner, 3 = Advanced)')

# Adjust layout
plt.tight_layout()
plt.show()
```

```
In [193]: # EDA Multivariate Analysis

In [195]: # Compute correlation matrix
corr_matrix = gd.corr(numeric_only=True)

# Plot heatmap of correlations
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f', linewidth=0.5)
plt.title('Correlation Heatmap')
plt.show()
```