#### **HEART DISEASE PREDICTION**

#### **OBJECTIVE:**

TO PERFORM THE EXPLORATORY DATA ANALYSIS IN THE DATASET TO HELP LOOK AT DATA BEFORE MAKING ANY ASSUMPTIONS. IT CAN HELP IDENTIFY OBVIOUS ERRORS, AS WELL AS BETTER UNDERSTAND PATTERNS WITHIN DATA, DETECT OUTLIERS OR ANOMALOUS EVENTS, FIND INTERESTING RELATIONS AMONG THE VARIABLES IN THE DATASET AND INSIGHTS ABOUT THE FACTORS AFFECTING THE HEART ATTACK RISK.

#### **ABOUT THE DATASET:**

THIS DATASET IS USED TO PREDICT HEART
DIESASE. PATIENTS ARE LIKELY TO BE DIAGNOSED
WITH ANY CARDIOVASCULAR HEART DISEASES
BASED ON THEIR MEDICAL ATTRIBUTES SUCH AS
AGE, BLOOD PRESSURE, BMI, CHOLESTEROL,

TRIGLYCERIADS, SMOKING, ALCOHOL
CONSUMPTION, STRESS LEVEL, DIABETICS, HEART

## RATE etc.

# **EXPLORATORY ANALYSIS**

```
In [105]: import numpy as np
import pandau am pd
import mathlotlib.pyplot am plt
import semborn am ab
import warnings
warnings.filterwarnings('ignore')
%mathlotlib inline
```

```
In [3]: data = pd.read_csv(r*Documents\heart_attack_prediction_dataset.csv*)

In [4]: data
```

]:		Patient ID	Age	Sex	Cholesterol	Blood Pressure	Heart Rate	Diabetes	Family History	Smoking	Obesity	-	Sedentary Hours Per Day	income	DIMI	Triglycerides	Physical Activity Days Per Week
	0 1	BMW7612	67.0	Male	208	158/88	72	0	0	1	0		6.615001	261484	31.251233	288	0
	1	GZE1114	21.0	Male	389	105/93	98	1	1	1	- 1		4.963459	285768	27.194973	235	1
	2	BN89906	21.0	Female	324	174/99	72	1	0	0	0		9.463426	235262	28.176571	587	4
	3	JLN3497	84.0	Male	303	163/100	73	1	1	1	0		7.640901	125640	36.464704	378	3
	4	GF08847	66.0	Male	318	91/00	93	1	1	1	1		1.514821	190555	21.009144	231	1
			-		-	-	-		-		-		-	-	-		
87	58	MSV9918	60.0	Male	121	94/76	61	1	1	1	0		10.806373	295429	19.655895	67	7
87	59	Q5V6764	28.0	Female	120	157/102	73	1	0	0	1		3.833838	217881	23.993866	617	4
87	50	XXA5925	4T.0	Male	250	161/75	105	0	1	1	- 1		2.375214	36966	35.406146	527	4
87	81	EPE6801	36.0	Male	178	119/67	60	1	0	1	0		0.029104	299943	27.294020	114	2
27	**	TANKS	25.0	Female	266	TREAT	75						0.005734	247338	22.014161	188	,

In [6]:																		
		Patient ID	Age	Sex	Cholesterol	Blood Pressure	Heart Rate	Diabetes	Family History	Smoking	Obesity	-	Sedentary Hours Per Day	Income	BMI	Triglycerides	Activity Duya Per Week	H
	0	BMW7812	67.0	Male	208	158/68	72	0	0	1	0		6.615801	261404	31.251233	296	0	
	1	CZE1114	21.0	Male	309	165/93	96	1	1	1	1		4.963459	265766	27.194973	235	1	
	2	EN19906	21.0	Female	324	174/99	72	1	0	0	0		9.463426	235262	28.176571	587	4	
	3	JUN3497	84.0	Male	383	163/100	73	1	1	1	0		7.648901	125840	36.464704	378	3	
	4	0F08847	96.0	Male	318	91/88	93	1	1	1	1		1.514821	160555	21.809144	291	1	

```
In [8]: data.tail()
Out[8]:
                 Patient Age
                                                Blood Heart Diabetes Family Smoking Obesity
                                See Cholesterol
                                                                                                                          DMI Triglycerides
          8758 M5V9918 60.0
                                                   9476
                                                                                                    10.806373 235429 19.655895
          8769 QSV6764 20.0 Female
                                            130
                                                 157/102
                                                            73
                                                                                                     3.030030 217001 23.990066
                                                                                                                                       917
          8760 XXXAS925 47.0
                              Male
                                            250
                                                  161/75
                                                           105
                                                                                                     2.375214 36998 35.406146
                                                                                                                                       527
          8761 EPE6801 36.0
                              Male
                                            178
                                                  119/67
                                                            68
                                                                                                     0.029104 209943 27.294020
                                                                                                                                       114
                                                                                                                                                 2
          6762 ZWN9666 25.0 Female
                                            356
                                                  135/57
                                                                                                     9.005234 247338 32.914151
         5 rows × 26 columns
```

```
In [40]: data.shape
Out[40]: (8763, 22)
In [48]: data.dtypes
Out[48]: Age
                                              float64
                                              object
int64
          Chalesteral
          Blood Pressure
                                               object
          Heart Rate
                                                1nt64
          Diabetes
                                                int64
          Family History
                                                int64
          Smoking
                                                Int64
          Obesity
                                                int64
          Alcohol Consumption
                                                int64
          Exercise Hours Per Neek
                                              float64
          Diet
                                              object
          Previous Heart Problems
                                                1nt64
          Medication Use
                                                int64
          Stress Level
                                                int64
                                              float64
          Sedentary Hours Per Day
                                               int64
          Income
                                              float64
          Triglycerides
                                                int64
          Physical Activity Days Per Week
                                                int64
          Sleep Hours Per Day
                                                int64
          Heart Attack Risk
                                                int64
          dtype: object
```

# In [9]: data.info() (class 'pandas.core.frame.DataFrame') RangeIndex: 8763 entries, 0 to 8762 Data columns (total 26 columns): # Column

Non-Null Count Otype Patient ID 8763 non-mull abject Age Sex 8762 non-null float64 8763 non-null abject Cholesterol 8763 non-null int64 Blood Pressure 8763 non-null object Heart Rate Diabetes 8763 non-null 8763 non-null Sat64 Sat64 Family History Smoking 8763 non-null 8763 non-null Set64 Sat64 Obesity Alcohol Consumption 8763 non-null Sat64 30 8763 non-null Set64 22 Exercise Hours Per Week 8761 non-null float64 Dist 0763 non-null 12 ablect. Previous Heart Problems 8763 non-null Medication Use 8763 non-null 24 Sat 64 Stress Level 8763 non-null Sat64 Sedentary Hours Per Day 36 8763 non-null float64 8763 non-null Sat64 8763 non-null float64 18 BACK int64 Triglycerides 8763 non-mull Physical Activity Days Per Heek Sleep Hours Per Day 20 21 8763 non-null Set64 8763 non-null 22 Country 8763 non-null object

```
In [49]: data.count()
Out[49]: Age
           Sex
Cholesterol
                                                     8763
8763
           Blood Pressure
Heart Rate
Diabetes
                                                      8763
                                                     6763
                                                     8763
           Family History
                                                     8763
           Smoking
Obesity
Alcohol Consumption
Exercise Hours Per Week
                                                     6763
                                                     8763
                                                     8763
                                                     6763
           Diet
                                                     8763
           Previous Heart Problems
                                                     8763
           Medication Use
                                                     6763
           Stress Level
                                                     8763
           Sedentary Hours Per Day
                                                     0763
                                                     6763
           Income
           BHI
                                                     8763
           Triglycerides
                                                     0763
           Physical Activity Days Per Week
                                                     6763
           Sleep Hours Per Day
Heart Attack Risk
                                                     8763
                                                     0763
           dtype: int64
```

```
In [41]: data.columns.tolist()
Out[41]: ['Age', 'Sex',
                   'Cholesterol',
                   'Blood Pressure',
'Heart Rate',
                  "Disbetten",
"Pamily History",
'Smoking',
'Obesity',
'Alcohol Consumption',
'Exercise Hours Fer Neek',
                   'Diet',
'Previous Heart Problems',
                    'Medication Use',
                   'Stress Level',
'Sedentary Hours Per Day',
                    'Income',
                   'BHE',
'Triglycerides',
'Physical Activity Days Per Week',
                   'Sleep Hours Per Day',
'Heart Attack Risk']
```

In [10]: data.nunique()

#### Out[10]: Patient ID 8763 Age Sex 73 2 Cholesterol 281 Blood Pressure Heart Rate 3915 71 Diabetes Family History Smoking Obesity Alcohol Consumption Exercise Hours Per Week 8763 Olet Previous Heart Problems 2 Medication Use 2 Stress Level 10 Sedentary Hours Per Day 8763 8615 Income 8763 Triglycerides Physical Activity Days Per Week 771 8 Sleep Hours Per Day Country Continent 20

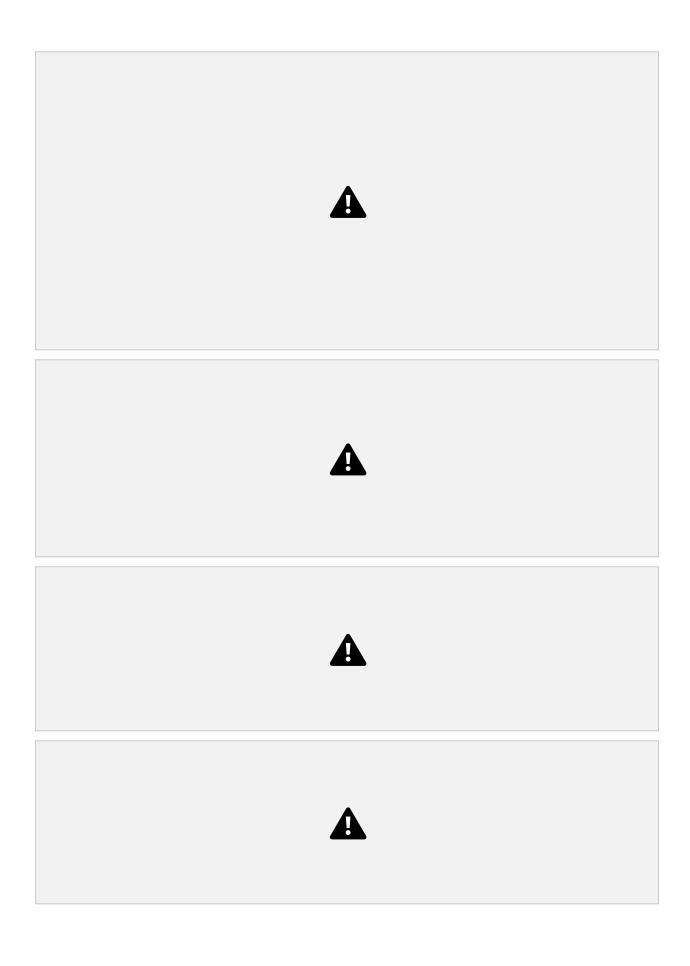
Hemisphere Heart Attack Risk dtype: int64

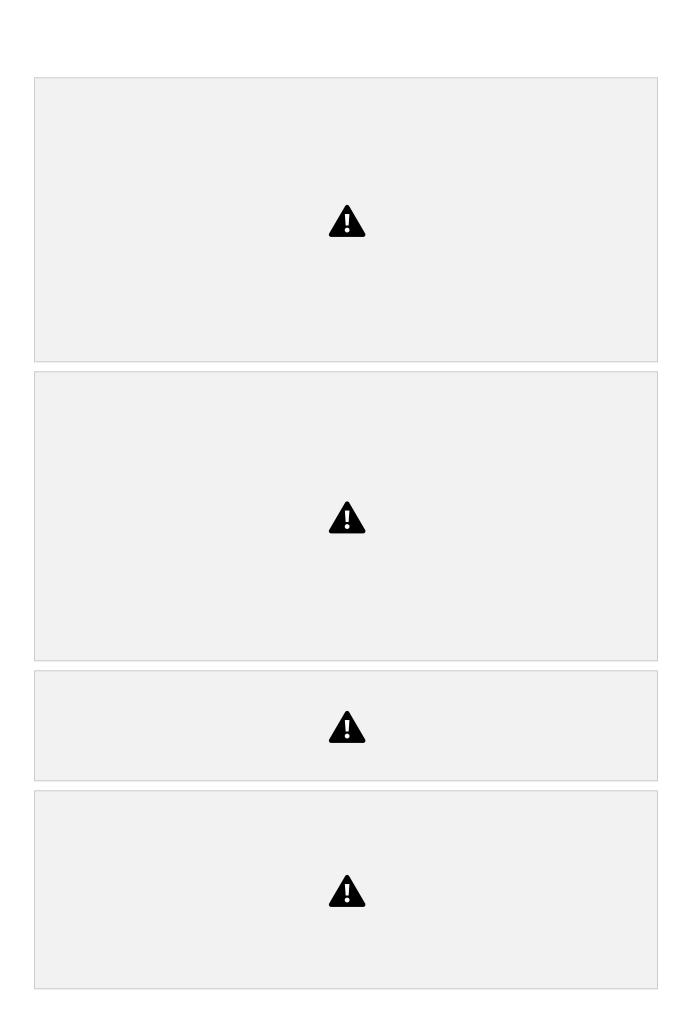
6

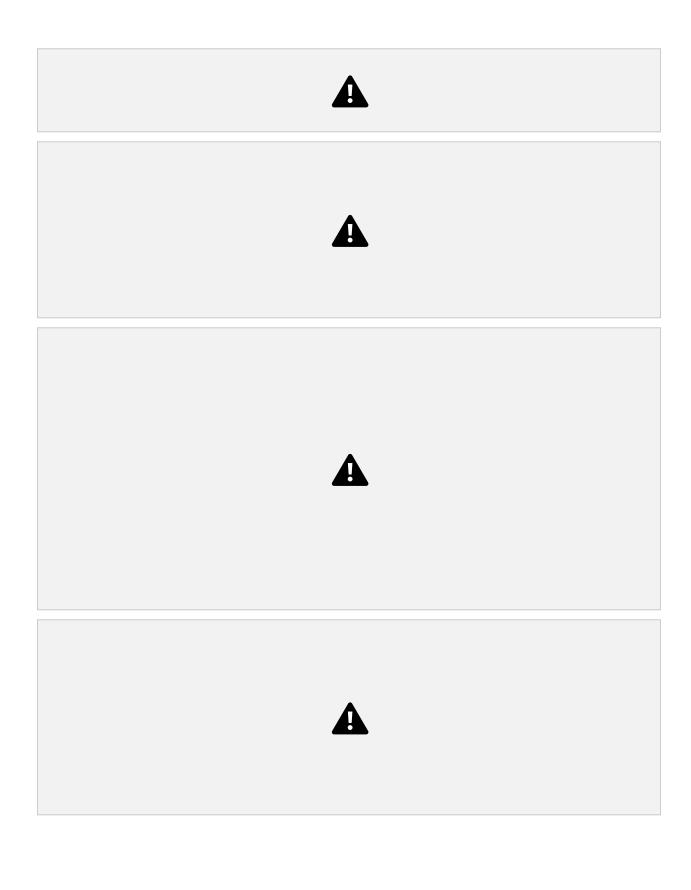
```
In [11]: data.isnull().sum()
Out[11]: Patient ID
            Age
Sex
                                                         0 0 0 0 0
             Cholesterol
            Blood Pressure
Heart Rate
            Diabetes
Family History
             Smoking
            Obesity
Alcohol Consumption
             Exercise Hours Per Week
            Diet
Previous Heart Problems
             Medication Use
            Stress Level
Sedentary Hours Per Day
             Income
             BME
             Triglycerides
             Physical Activity Days Per Week
             Sleep Hours Per Day
            Country
Continent
            Hemisphere
Heart Attack Risk
dtype: int64
In [68]: age_mean = data['Age'].mean()
data['Age'].fillna(value = age_mean, inplace = True)
data['Age']
Out[68]: 0
                      67.0
21.0
                       21.0
                      84.0
66.0
            8758
8759
                       60.0
                      28.0
             8760
                      47.0
             8761
                      36.0
             8762
                      25.0
             Name: Age, Length: 8763, dtype: float64
 In [69]: data.isnull().sum()
Out[69]| Age
Sex
Cholesterol
                                                         0
            Blood Pressure
Heart Rate
                                                         00000
            Diabetes
Family History
Smoking
            Obesity
Alcohol Consumption
             Exercise Hours Per Week
            Diet
Previous Heart Problems
             Medication Use
            Stress Level
Sedentary Hours Per Day
             Income
            BMI
Triglycerides
             Physical Activity Days Per Week
            Sleep Hours Per Day
Heart Attack Risk
             dtype: int64
```

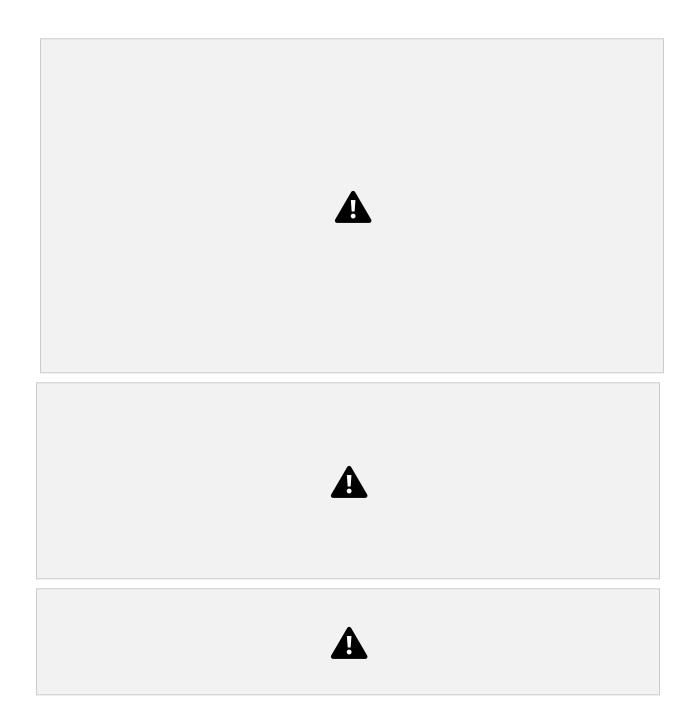
```
In [12]: (data.ismull().sum()/(len(data)))*100
Out[12]: Patient ID
                                                     0.000000
           Age
Sex
Cholesterol
                                                      0.011412
                                                     0.000000
           Blood Pressure
Heart Rate
                                                      0.000000
                                                     0.000000
           Oiabetes
Family History
                                                      0.000000
            Smoking
                                                     0.000000
           Obesity
Alcohol Consumption
                                                     0.000000
            Exercise Hours Fer Week
                                                     0.000000
           Diet
Previous Heart Problems
                                                      0.000000
           Medication Use
Stress Level
Sedentary Hours Per Day
                                                     0.000000
                                                      0.000000
            Income
                                                      0.000000
            5910
           oni
Triglycerides
Mhysical Activity Days Per Week
Sleep Hours Per Day
                                                      0.000000
                                                     0.000000
                                                      0.000000
           Country
Continent
                                                      0.000000
                                                      0.000000
            Hendsphere
                                                      0.000000
           Heart Attack Risk
dtype: float64
                                                     0.000000
In [70]: (deta--0).sum()
Out[78]: Age
                                                      2652
            Sex
            Cholesterol
                                                         0
            Blood Pressure
            Heart Rate
           Diabetes
Family History
                                                      3047
                                                      4443
            Smoking
                                                      984
           Obesity
Alcohol Consumption
                                                      4369
                                                      3522
            Exercise Hours Per Week
            Diet
                                                         0
            Previous Heart Problems
                                                      4418
           Medication Use
Stress Level
                                                      4396
            Sedentary Hours Per Day
            Income
            BMI
            Triglycerides
           Physical Activity Days Per Week
Sleep Hours Per Day
                                                     1865
           Heart Attack Risk
dtype: int64
                                                      5624
```











### **CONCLUSION**

THE EXPLORATORY DATA ANALYSIS FOR THE DATASET IS PERFORMED TO UNDERSTAND THE DATASET, BASED ON THE HIGH FACTORS AFFECTING THE HEART ATTACK RISK, THE INDEPENDENT VARIABLES AND DEPENDENT VARIABLES, AND

# CORRELATION BETWEEN THE VARIABLES.