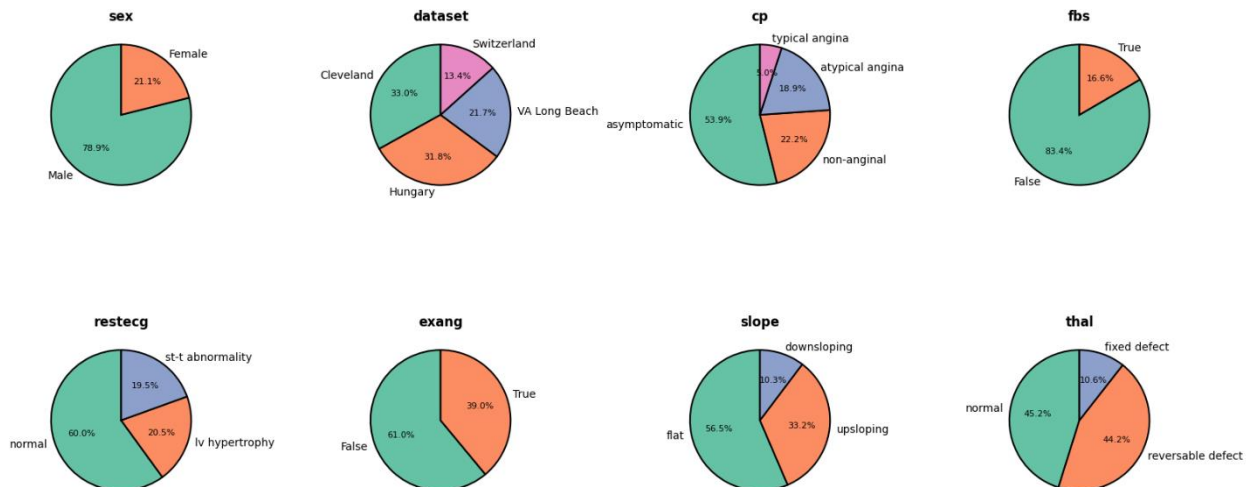# Heart Disease

"This analysis provides insights into a medical dataset focused on cardiovascular characteristics. The data has undergone prior cleaning, with each row representing the attributes of an individual patient. Notably, the 'age' column has been normalized to ensure equal contribution across its range of values in any subsequent modeling or analysis. This transformation provides a standardized basis for comparing the influence of age with other medical variables included in the dataset, such as resting blood pressure ('trestbps'), cholesterol ('chol'), fasting blood sugar levels ('fbs'), and other cardiac-related indicators."

the insights from each individual pie chart:



## 1. Pie Chart for sex

- **Insight:** This chart shows the distribution of gender in the dataset. It indicates that there's a significant class imbalance, with a higher proportion of males compared to females. This imbalance may need to be addressed during model training to avoid potential bias.

## 2. Pie Chart for dataset

- **Insight:** This pie chart illustrates the distribution of patients across different countries. It reveals the proportion of patients belonging to each country . Potential dataset-specific characteristics and biases should be considered during analysis.

## 3. Pie Chart for cp (Chest Pain Type)

- **Insight:** This chart visualizes the distribution of different types of chest pain experienced by the patients. It highlights the prevalence of each chest pain type and indicates that some types are more common than others, suggesting potential patterns in symptom presentation.

**4. Pie Chart for fbs (Fasting Blood Sugar)**

- **Insight:** This chart shows the proportion of patients with fasting blood sugar levels above and below a certain threshold. It reveals the prevalence of high or low blood sugar levels, which could be a risk factor for heart disease.

**5. Pie Chart for restecg (Resting Electrocardiographic Results)**

- **Insight:** This chart depicts the distribution of different resting electrocardiographic results. It indicates the prevalence of various ECG patterns, which could be indicative of underlying heart conditions or abnormalities.

**6. Pie Chart for exang (Exercise Induced Angina)**

- **Insight:** This chart visualizes the proportion of patients who experience angina during exercise. It highlights the prevalence of exercise-induced angina, which could be a significant symptom in diagnosing heart disease.

**7. Pie Chart for slope (Slope of the Peak Exercise ST Segment)**

- **Insight:** This chart shows the distribution of different slope patterns during exercise. It indicates the prevalence of various slope categories, which could reflect the severity or type of heart condition.
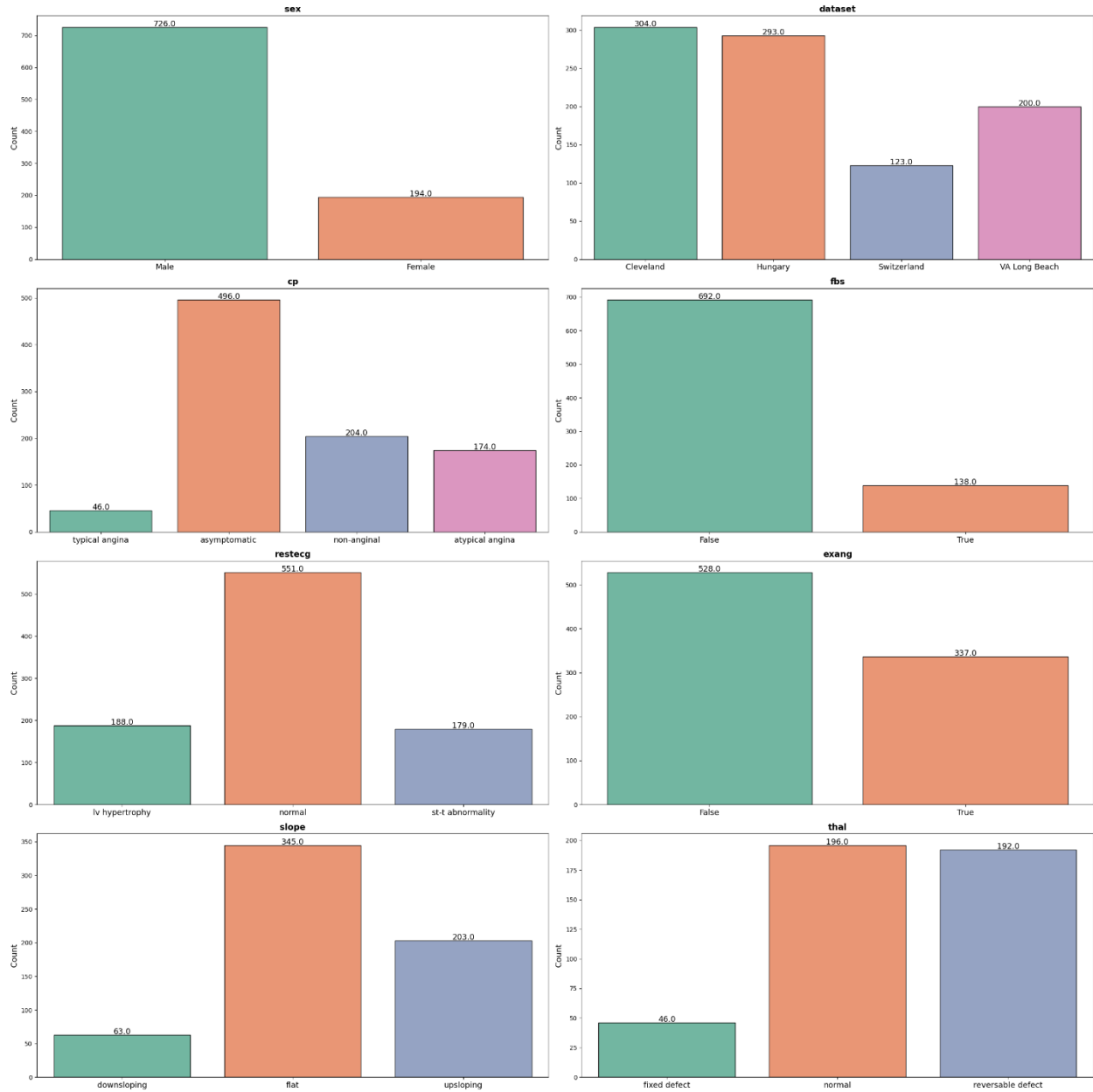
**8. Pie Chart for thal (Thalassemia)**

- **Insight:** This chart depicts the distribution of different thalassemia types or blood disorders. It highlights the prevalence of different thalassemia categories, which could be a potential risk factor for heart disease.

**Overall Insight:**

The pie charts collectively provide a visual representation of the distribution of various categorical features within the dataset. They highlight potential class imbalances, prevalence of different categories, and patterns in symptom presentation or risk factors, offering valuable insights for further analysis and model development.

the insights provided by each individual bar chart (countplot):

## 1. Bar Chart for sex

- **Insight:** This chart visually represents the distribution of gender in the dataset. It clearly shows that there's a significant class imbalance, with a higher number of males compared to females. This imbalance could potentially introduce bias during model training, and strategies like data augmentation or resampling might be needed to address it.

## 2. Bar Chart for dataset

- **Insight:** This bar chart illustrates the distribution of patients across different datasets. It reveals the number of patients belonging to each dataset category, highlighting potential variations or biases associated with specific datasets. Considering dataset-specific characteristics is crucial during analysis to avoid skewed interpretations.

### 3. Bar Chart for cp (Chest Pain Type)

- **Insight:** This chart visualizes the frequency of different types of chest pain experienced by the patients. It clearly shows that some types of chest pain are more prevalent than others, suggesting potential patterns in symptom presentation and their association with heart disease.

### 4. Bar Chart for fbs (Fasting Blood Sugar)

- **Insight:** This chart depicts the distribution of patients with fasting blood sugar levels above and below a certain threshold. It reveals the number of patients with high or low blood sugar levels, indicating the prevalence of this potential risk factor for heart disease within the dataset.

### 5. Bar Chart for restecg (Resting Electrocardiographic Results)

- **Insight:** This chart illustrates the frequency of different resting electrocardiographic results. It shows the number of patients with various ECG patterns, highlighting the prevalence of specific ECG abnormalities or indicators that could be associated with heart conditions.

### 6. Bar Chart for exang (Exercise Induced Angina)

- **Insight:** This chart visualizes the distribution of patients who experience angina during exercise. It clearly shows the number of patients with and without exercise-induced angina, indicating the prevalence of this significant symptom in the dataset and its potential role in diagnosing heart disease.

### 7. Bar Chart for slope (Slope of the Peak Exercise ST Segment)

- **Insight:** This chart depicts the frequency of different slope patterns during exercise. It shows the number of patients belonging to various slope categories, revealing the prevalence of specific slope patterns and their potential association with the severity or type of heart condition.
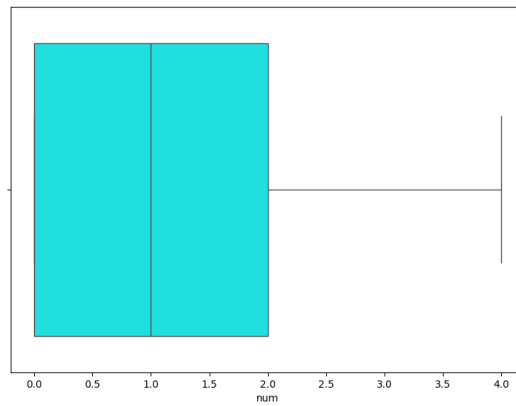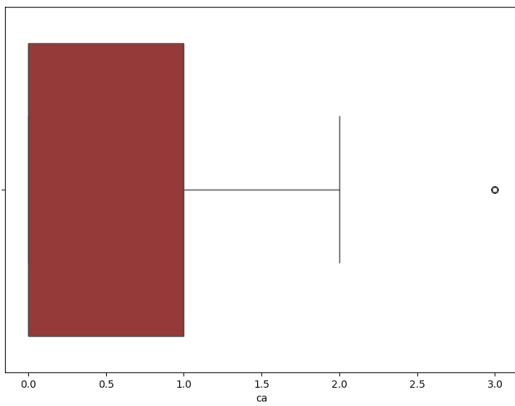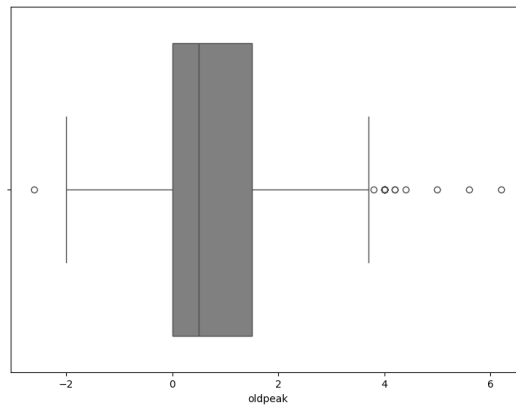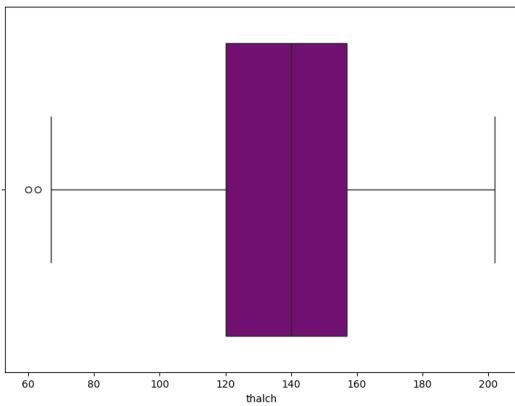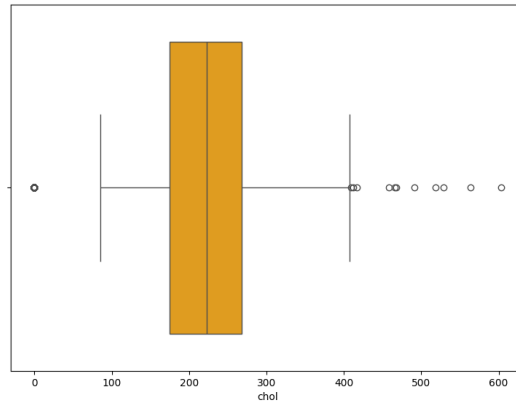
### 8. Bar Chart for thal (Thalassemia)

- **Insight:** This chart illustrates the distribution of different thalassemia types or blood disorders. It shows the number of patients in each thalassemia category, highlighting the prevalence of specific thalassemia types and their potential link to heart disease risk.

**Overall Insight:**

The bar charts (countplots) provide a clear and comparative view of the distribution of categorical features in the dataset. They reveal class imbalances, highlight the frequency of different categories, and offer insights into potential patterns, risk factors, and symptom presentations associated with heart disease. These visualizations are valuable for understanding the data's characteristics and informing further analysis and model development.

the insights gleaned from each individual box plot:

**1. Box Plot for age**

- **Insight:** This box plot visualizes the distribution of age in the dataset. It shows the median age, the interquartile range (IQR), and potential outliers. It reveals the typical age range of patients in the dataset and identifies any unusually young or old individuals. It also gives you an idea about the spread and skewness of the age distribution.

**2. Box Plot for trestbps (Resting Blood Pressure)**

- **Insight:** This box plot depicts the distribution of resting blood pressure values. It highlights the median resting blood pressure, the IQR, and any outliers. It shows the typical range of resting blood pressure and identifies any individuals with unusually high or low blood pressure, which could be indicative of health issues.

**3. Box Plot for chol (Serum Cholesterol)**

- **Insight:** This box plot visualizes the distribution of serum cholesterol levels. It shows the median cholesterol level, the IQR, and potential outliers. It helps understand the typical cholesterol range and identify any individuals with unusually high or low cholesterol, which could be a risk factor for heart disease.

**4. Box Plot for thalch (Maximum Heart Rate Achieved)**

- **Insight:** This box plot depicts the distribution of maximum heart rates achieved during exercise. It highlights the median maximum heart rate, the IQR, and any outliers. It shows the typical range of maximum heart rates and identifies any individuals with unusually high or low values, which could reflect their fitness levels or underlying heart conditions.

**5. Box Plot for oldpeak (ST Depression Induced by Exercise Relative to Rest)**

- **Insight:** This box plot visualizes the distribution of ST depression induced by exercise. It shows the median ST depression, the IQR, and potential outliers. It helps understand the typical range of ST depression and identify individuals with unusually high or low values, which could indicate the severity of heart disease.

**6. Box Plot for ca (Number of Major Vessels Colored by Fluoroscopy)**

- **Insight:** This box plot depicts the distribution of the number of major vessels colored by fluoroscopy. It highlights the median number of vessels, the IQR, and any outliers. It shows the typical range of vessel coloration and identifies any individuals with unusually high or low values, which could reflect the extent of coronary artery disease.
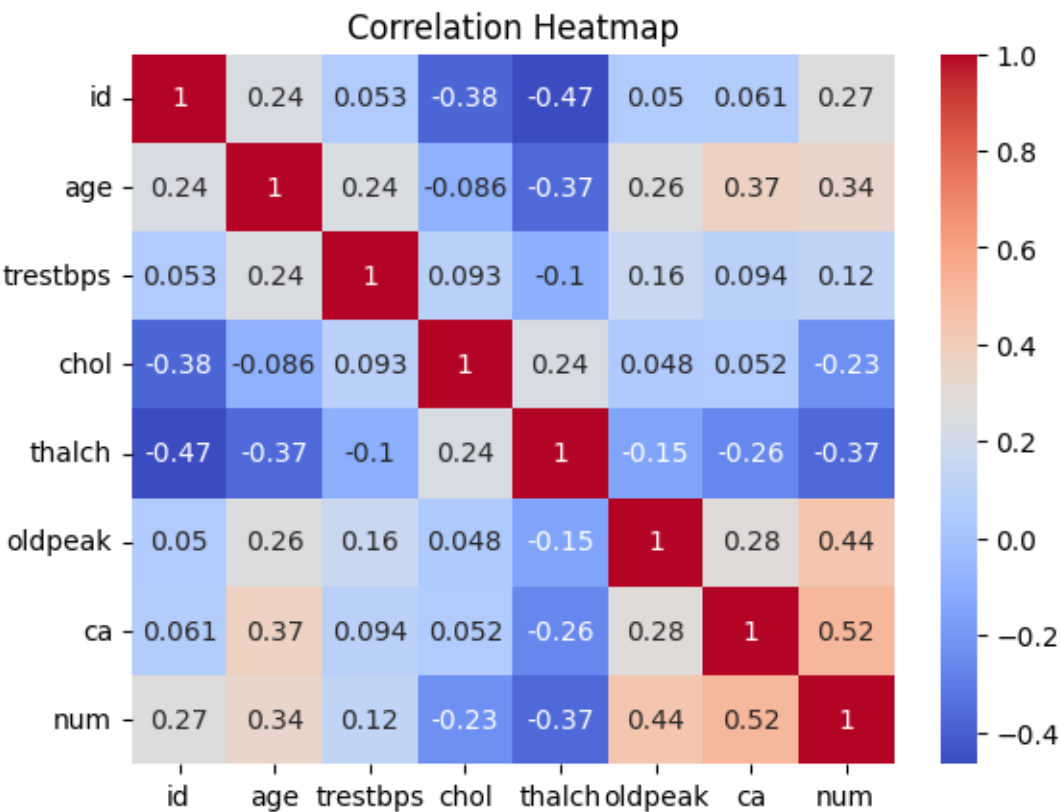
### 7. Box Plot for num (Diagnosis of Heart Disease)

- **Insight:** This box plot visualizes the distribution of the diagnosis of heart disease. Since this is likely a categorical variable (0 or 1), the box plot might not be the most informative visualization. However, it could still show the median and IQR, giving a basic overview of the distribution of heart disease diagnoses in the dataset.

**Overall Insight:**

The box plots provide a visual representation of the distribution of numerical features in the dataset. They highlight central tendency, spread, and potential outliers, offering insights into the typical ranges of values and identifying individuals with unusual characteristics. These visualizations are crucial for understanding data distribution, identifying potential data quality issues, and guiding further analysis.

the insights we can derive from the correlation analysis presented , specifically focusing on the two correlation heatmaps:



Correlation Heatmap

|           | id     | age    | trestbps | chol   | thalch | oldpeak | ca    | num   |
|-----------|--------|--------|----------|--------|--------|---------|-------|-------|
| id        | 1      | 0.24   | 0.053    | -0.38  | -0.47  | 0.05    | 0.061 | 0.27  |
| age       | 0.24   | 1      | 0.24     | -0.086 | -0.37  | 0.26    | 0.37  | 0.34  |
| trestbps  | 0.053  | 0.24   | 1        | 0.093  | -0.1   | 0.16    | 0.094 | 0.12  |
| chol      | -0.38  | -0.086 | 0.093    | 1      | 0.24   | 0.048   | 0.052 | -0.23 |
| thalch    | -0.47  | -0.37  | -0.1     | 0.24   | 1      | -0.15   | -0.26 | -0.37 |
| oldpeak   | 0.05   | 0.26   | 0.16     | 0.048  | -0.15  | 1       | 0.28  | 0.44  |
| ca        | 0.061  | 0.37   | 0.094    | 0.052  | -0.26  | 0.28    | 1     | 0.52  |
| num       | 0.27   | 0.34   | 0.12     | -0.23  | -0.37  | 0.44    | 0.52  | 1     |

### 1. Correlation Heatmap between Numerical Features:

- **Insight:** This heatmap visualizes the correlation coefficients between pairs of numerical features in your dataset. It provides insights into the relationships between variables:

    - **Strong Positive Correlation:** Features with correlation coefficients close to +1 indicate a strong positive linear relationship. For example, "cp" (chest pain type) and "thalach" (maximum heart rate achieved) may show a positive correlation, suggesting that certain types of chest pain are associated with higher heart rates.

    - **Strong Negative Correlation:** Features with correlation coefficients close to -1 indicate a strong negative linear relationship. For example, "age" and "thalach" may show a negative correlation, suggesting that older individuals tend to have lower maximum heart rates.

    - **Weak or No Correlation:** Features with correlation coefficients close to 0 indicate a weak or no linear relationship. This means that changes in one variable are not strongly associated with changes in the other variable.

- **Overall:** This heatmap helps identify potential multicollinearity (high correlation between independent variables) and understand which features might be more influential in predicting the target variable.

Correlation of Features with Target

| Feature | num |
|---|---|
| num | 1 |
| ca | 0.52 |
| oldpeak | 0.44 |
| dataset_Hungary | -0.38 |
| thalch | -0.37 |
| exang_True | 0.35 |
| cp_atypical angina | -0.34 |
| age | 0.34 |
| dataset_Switzerland | 0.28 |
| thal_reversable defect | 0.28 |
| id | 0.27 |
| sex_Male | 0.26 |
| slope_flat | 0.24 |
| dataset_VA Long Beach | 0.24 |
| chol | -0.23 |
| thal_normal | -0.22 |
| restecg_normal | -0.17 |
| cp_non-anginal | -0.16 |
| slope_upsloping | -0.14 |
| restecg_st-t abnormality | 0.13 |
| fbs_True | 0.13 |
| trestbps | 0.12 |
| cp_typical angina | -0.056 |

**2. Correlation Heatmap with Target Variable (num)**

- **Insight:** This heatmap focuses on the correlation between each feature and the target variable (num), which represents the diagnosis of heart disease. It helps identify features that are most strongly associated with the presence or absence of heart disease:

    - **Positive Correlation:** Features with positive correlation coefficients indicate that higher values of the feature are associated with a higher likelihood of heart disease. For example, "oldpeak" (ST depression) may show a positive correlation with num, suggesting that higher ST depression values are linked to a higher risk of heart disease.

    - **Negative Correlation:** Features with negative correlation coefficients indicate that lower values of the feature are associated with a higher likelihood of heart disease. For example, "thalach" (maximum heart rate achieved) may show a negative correlation with num, suggesting that lower maximum heart rates are linked to a higher risk of heart disease.

- **Overall:** This heatmap provides valuable insights into feature importance and can guide feature selection for building predictive models. Features with stronger

correlations to the target variable are likely to be more informative in predicting heart disease.

**General Considerations:**

- **Correlation does not imply causation:** While correlation reveals relationships between variables, it does not necessarily mean that one variable causes the other. Further analysis and domain expertise are required to establish causal relationships.

- **Linearity:** Correlation coefficients measure linear relationships between variables. Non-linear relationships may exist but might not be captured by correlation analysis.

- **Context is key:** The insights from correlation analysis should be interpreted in the context of the specific dataset and the research question. Domain expertise is crucial for understanding the practical implications of the observed correlations.

the insights we can gain from the scatter plots (specifically the pairplot) :

**Pairplot Insights**

The sns.pairplot function creates a matrix of scatter plots, allowing you to visualize the relationships between pairs of numerical features in your dataset. Here's a breakdown of the key insights you can gather from these scatter plots:

**1. Relationships between Features:**

- **Linear Relationship:** If the points in a scatter plot roughly form a straight line, it suggests a linear relationship between the two features. This could be either a positive linear relationship (as one feature increases, the other also tends to

increase) or a negative linear relationship (as one feature increases, the other tends to decrease).

- **Non-linear Relationship:** If the points form a curve or other non-linear pattern, it suggests a more complex relationship between the features. This could indicate a quadratic relationship, an exponential relationship, or other types of non-linear dependencies.

- **No Relationship:** If the points are scattered randomly without any discernible pattern, it suggests that there is no clear relationship between the two features.

## 2. Clusters and Outliers:

- **Clusters:** Scatter plots can reveal clusters of data points, indicating groups of observations with similar characteristics. These clusters can provide insights into underlying patterns or subgroups within the dataset.

- **Outliers:** Outliers are data points that are significantly different from the majority of the data. They can be identified as points that are far away from the main cluster of points in a scatter plot. Outliers can be indicative of data errors or unusual observations that require further investigation.

## 3. Distribution of Individual Features:

- **Histograms:** The diagonal plots in a pairplot typically show histograms, which display the distribution of individual features. Histograms can reveal the shape of the distribution (e.g., normal, skewed), the central tendency, and the spread of the data.

**Specific Insights from your Pairplot:**

Based on the features included in your pairplot
(age, trestbps, chol, thalch, oldpeak, ca, num), you can gather insights such as:

- **Relationship between age and thalch:** You might observe a negative linear relationship, indicating that as age increases, the maximum heart rate achieved tends to decrease.

- **Relationship between oldpeak and num:** You might observe a positive relationship, suggesting that higher values of ST depression are associated with a higher likelihood of heart disease.

- **Clusters based on ca and num:** You might observe clusters of data points based on the number of major vessels colored by fluoroscopy and the diagnosis of heart

disease. This could indicate distinct subgroups of patients with different characteristics.

- **Outliers:** You might identify outliers in any of the features, which could warrant further investigation to determine if they are data errors or unusual observations.

**Overall:**

The scatter plots in your pairplot provide a visual exploration of the relationships between numerical features in your dataset. By examining the patterns, clusters, and outliers in these plots, you can gain valuable insights into the underlying structure of your data and identify potential factors that contribute to heart disease. Remember to interpret these insights in the context of your specific research question and domain knowledge.