# What is Exploratory Data Analysis (EDA)?

**Exploratory Data Analysis (EDA)** is the process of analyzing and summarizing the main characteristics of a dataset, often using visual and statistical methods. It helps identify patterns, relationships, anomalies, and insights in data, which can guide further analysis or model building.

EDA is a crucial step in the data analysis pipeline and is used to:

1. **Understand the Dataset**: Identify the structure, size, types of data, and distributions.
2. **Identify Data Quality Issues**: Detect missing values, outliers, and inconsistencies.
3. **Discover Relationships**: Explore correlations and patterns between variables.
4. **Feature Engineering**: Inform decisions about creating, combining, or removing features.
5. **Prepare for Modeling**: Ensure data is clean and ready for machine learning or statistical analysis.

---

## Key Points in EDA

Here are the critical aspects to focus on during EDA:

---

## 1. Understand the Data Structure

- **Key Actions**:
    - Look at the shape of the dataset (rows and columns).
    - Inspect data types (`int`, `float`, `object`, etc.).
    - View sample data (e.g., `head()` and `tail()`).
- **Code Example**:
- `print(data.shape)`
- `print(data.info())`
- `print(data.head())`

---

## 2. Handle Missing Values

Missing values can cause issues in analysis or modeling and need proper handling.

- **Key Questions**:
    - How many missing values are there per column?
    - Are they missing completely at random or due to patterns?
- **Handling Options**:
    - **Imputation**: Fill missing values with mean, median, mode, or other strategies.
    - **Dropping**: Remove rows or columns with excessive missing values.

o **Placeholder**: Use a placeholder value (e.g., `-1`, `Unknown`).
- **Code Example**:
- `print(data.isnull().sum())`
- `data.fillna(data.mean(), inplace=True)  # Example: Replace missing values with the column mean`

---

## 3. Check for Duplicates

Duplicate rows may introduce bias and should be handled.

- **Key Action**:
- `print(data.duplicated().sum())  # Check duplicates`
- `data.drop_duplicates(inplace=True)  # Remove duplicates`

---

## 4. Analyze Data Distribution

Understand how data is distributed for numerical and categorical features.

- **Numerical Features**:
    o Use histograms, boxplots, or density plots.
    o Check for skewness and outliers.
- **Categorical Features**:
    o Use bar plots or pie charts.
- **Code Example**:
- `data.hist(figsize=(10, 8))`
- `sns.boxplot(data['feature_name'])`
- `sns.countplot(data['categorical_feature'])`

---

## 5. Identify and Handle Outliers

Outliers can distort analysis and affect model performance.

- **Methods to Detect Outliers**:
    o Boxplot (IQR rule).
    o Z-score or Standard Deviation.
    o Visual inspection.
- **Code Example**:
- `sns.boxplot(data['feature'])`

---

## 6. Explore Relationships Between Features

- **Correlations**:
    o Find relationships between numerical features.

- o   Use a heatmap for visualization.
- **Feature Relationships**:
  - o   Compare features with the target variable (e.g., scatterplots, bar plots).
- **Code Example**:
- `corr_matrix = data.corr()`
- `sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')`

---

## 7. Summarize Descriptive Statistics

Understand the central tendency, dispersion, and shape of the data distribution.

- **Key Metrics**:
  - o   Mean, median, mode.
  - o   Standard deviation, variance.
  - o   Minimum, maximum, and quartiles.
- **Code Example**:
- `print(data.describe())`

---

## 8. Feature Engineering Opportunities

- **Transformations**:
  - o   Normalize or standardize numerical features.
  - o   Encode categorical features (e.g., one-hot encoding).
- **Creation**:
  - o   Combine features or create new ones (e.g., BMI from height and weight).

---

## Summary

EDA is about asking questions like:

1. Are there missing values or duplicates?
2. What are the distributions of numerical features? Any skewness or outliers?
3. How are categorical variables distributed?
4. Which features are correlated with each other or the target variable?

EDA lays the foundation for robust analysis or modeling. Would you like help performing EDA on a specific dataset?

Performing **Exploratory Data Analysis (EDA)** on the Heart Disease dataset helps to understand its structure, patterns, and key insights. Here's a step-by-step guide with code snippets:

---

# 1. Import Necessary Libraries

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Settings for better visuals
sns.set(style='whitegrid')
```

---

# 2. Load the Dataset

```python
# Load the dataset
file_path = 'processed.cleveland.data'  # Replace with your file path
columns = [
    "age", "sex", "cp", "trestbps", "chol", "fbs", "restecg",
    "thalach", "exang", "oldpeak", "slope", "ca", "thal", "target"
]  # Replace with actual column names

data = pd.read_csv(file_path, header=None, names=columns)

# Preview the dataset
print(data.head())
```

---

# 3. Data Cleaning

### a. Handle Missing Values

```python
# Replace '?' with NaN
data.replace('?', np.nan, inplace=True)

# Convert numeric columns to appropriate data types
data = data.apply(pd.to_numeric, errors='coerce')

# Check missing values
print(data.isnull().sum())

# Impute or drop missing values
data.fillna(data.mean(), inplace=True)
```

### b. Check for Duplicates

```python
# Check for duplicate rows
print("Duplicate Rows:", data.duplicated().sum())

# Drop duplicates if any
data.drop_duplicates(inplace=True)
```

---

## 4. Summary Statistics

### a. General Overview

```python
# Dataset overview
print(data.info())

# Summary statistics
print(data.describe())
```

### b. Target Distribution

```python
# Check target variable distribution
print(data['target'].value_counts())

# Visualize target distribution
sns.countplot(data['target'])
plt.title('Distribution of Target Variable')
plt.show()
```

---

## 5. Visualizations

### a. Correlation Heatmap

```python
# Correlation matrix
corr_matrix = data.corr()

# Heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(corr_matrix, annot=True, fmt=".2f", cmap="coolwarm")
plt.title("Correlation Matrix")
plt.show()
```

### b. Numerical Feature Distributions

```python
# Histograms for numerical features
data.hist(figsize=(12, 10), bins=20, color='skyblue', edgecolor='black')
plt.suptitle("Feature Distributions")
plt.show()
```

### c. Relationship with Target

```python
# Boxplot for numerical features vs target
plt.figure(figsize=(12, 6))
sns.boxplot(x='target', y='age', data=data)
plt.title('Age vs Target')
plt.show()

# Pairplot for selected features
sns.pairplot(data, vars=['age', 'trestbps', 'chol', 'thalach', 'oldpeak'],
hue='target', diag_kind='kde')
```

```
plt.suptitle('Pairplot of Features (by Target)', y=1.02)
plt.show()
```

### d. Categorical Feature Analysis

```
# Bar plots for categorical features
categorical_features = ['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope',
'thal']

for feature in categorical_features:
    plt.figure(figsize=(6, 4))
    sns.countplot(x=feature, hue='target', data=data)
    plt.title(f'{feature.capitalize()} vs Target')
    plt.show()
```

---

## 6. Key Insights

1. **Target Distribution**: Check for class imbalance in the target variable.
2. **Correlations**: Analyze how strongly features like `thalach` (maximum heart rate achieved) or `oldpeak` correlate with heart disease.
3. **Feature Relationships**: Identify trends like whether higher cholesterol (`chol`) or resting blood pressure (`trestbps`) are associated with heart disease.
4. **Categorical Feature Impact**: Understand how `sex` (gender) or `cp` (chest pain type) affects heart disease outcomes.

---

## 7. Next Steps

- Use findings to engineer new features (e.g., `age_group`, `chol_level`).
- Prepare the dataset for machine learning by scaling numerical features and encoding categorical variables.

Would you like to focus on a specific part of the EDA, such as deeper insights or more advanced visualizations?