

Data preparation: Download heart dataset from following link.
<https://www.kaggle.com/zhaoyingzhu/heartcsv> Perform followi
 operation on
 given dataset.
 a) Find Shape of Data
 b) Find Missing Values
 c) Find data type of each column
 d) Finding out Zero's
 e) Find Mean age of patients
 f) Now extract only Age, Sex, ChestPain, RestBP, Chol.
 Randomly divide dataset in training (75%) and testing (25%). Th
 the diagnosis
 test I predicted 100 report as COVID positive, but only 45 of the
 were actually
 positive. Total 50 people in my sample were actually COVID pos
 have total
 500 samples. Create confusion matrix based on above data and
 I. Accuracy
 II. Precision
 III. Recall
 IV. F-1 score
 more theory about this topic
 more information about theory



Data Preparation: Heart Disease Dataset

To begin with, let's walk through the necessary steps involved in data preparation and anal
 the heart disease dataset, followed by a deeper understanding of the confusion matrix met
 (Accuracy, Precision, Recall, and F1 score) in the context of a classification task.

1. Find Shape of Data

In Python (using pandas), the shape of a dataset refers to the number of rows and columns
 shape can be found using:

python



```
import pandas as pd # Load the dataset df = pd.read_csv('heart.csv') # Make sure
dataset is downloaded and in the correct path # Get the shape of the dataset df.s
```

This will return a tuple of the form (number of rows, number of columns) .

2. Find Missing Values

Checking for missing values is essential to understand the quality of the data and handle th
 properly (either by filling or removing them).

python



```
# Check for missing values df.isnull().sum()
```

This will show how many missing values exist in each column. If there are any, they can be
 using techniques like imputation (filling missing values with mean, median, or mode) or by
 rows or columns.

3. Find Data Type of Each Column

To check the data type of each column, which is important for understanding how each feature is represented and whether any type conversion is needed:

```
python
```

```
# Get the data types of each column df.dtypes
```



This will show the data type of each column (e.g., `int64`, `float64`, `object` for strings). For numeric features are of numeric types and categorical features are of type `object` or `category`.

4. Finding Out Zero's

To find columns with zero values (which may not make sense for certain features like age, blood pressure, etc.), you can check if any feature contains zero values:

```
python
```

```
# Check for zero values in the dataset df[df == 0].count()
```



You can also analyze specific columns by doing:

```
python
```

```
# Find zero values in a specific column, e.g., 'Chol' df[df['Chol'] == 0]
```



This could help identify if zero values are valid or if they need to be corrected.

5. Find Mean Age of Patients

The average or mean age of patients can be calculated as follows:

```
python
```

```
# Find the mean of the 'Age' column mean_age = df['Age'].mean() print("Mean Age:", mean_age)
```



This gives the average age of the patients in the dataset.

6. Extract Selected Features (Age, Sex, ChestPain, RestBP, Chol)

You can select only a subset of features to work with:

```
python
```

```
# Extract the selected columns selected_columns = df[['Age', 'Sex', 'ChestPain', 'RestBP', 'Chol']] print(selected_columns.head())
```



This will give you a new dataframe containing just the features you want.

7. Split Dataset into Training and Testing Data

Randomly split the dataset into training (75%) and testing (25%) sets. This is important for the model performance on unseen data.

python



```
from sklearn.model_selection import train_test_split # Split data into features (
target (y) X = df.drop(columns=['Target']) # assuming 'Target' is the diagnosis c
= df['Target'] # Split data into training and testing sets X_train, X_test, y_tra
y_test = train_test_split(X, y, test_size=0.25, random_state=42)
```

8. Confusion Matrix and Performance Metrics

You provided a diagnostic test scenario for COVID detection. Let's interpret this scenario and use a confusion matrix to evaluate the model:

- Predicted COVID Positive = 100
- True COVID Positive = 45
- Total actual COVID Positive = 50
- Total samples = 500

From this information, you can create the confusion matrix and compute various performance metrics.

Confusion Matrix

A confusion matrix has the following structure:

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

Given the scenario:

- True Positives (TP) = 45 (COVID positive and predicted positive)
- False Negatives (FN) = 5 (COVID positive but predicted negative, i.e., 50 total actual positives minus 45 TP)
- False Positives (FP) = 55 (COVID negative but predicted positive, i.e., 100 predicted positives minus 45 TP)
- True Negatives (TN) = 400 (the remaining samples, as there are 500 samples in total)

You can create the confusion matrix and calculate the metrics as follows:

python



```
from sklearn.metrics import confusion_matrix # Assuming the model prediction for
test set is stored in y_pred y_pred = [1] * 100 + [0] * 400 # Predicted positives
100, the rest negatives y_true = [1] * 50 + [0] * 450 # Actual 50 positives, rest
negatives # Create confusion matrix cm = confusion_matrix(y_true, y_pred)
print("Confusion Matrix:\n", cm)
```

Metrics:

1. **Accuracy:** The accuracy measures how often the model is correct. It is calculated as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

2. **Precision:** Precision measures the proportion of positive predictions that are actually correct. It is calculated as:

$$\text{Precision} = \frac{TP}{TP + FP}$$

3. **Recall (Sensitivity):** Recall measures the proportion of actual positives that are correctly identified. It is calculated as:

$$\text{Recall} = \frac{TP}{TP + FN}$$

4. **F1 Score:** The F1 Score is the harmonic mean of Precision and Recall, and is used when a balance between Precision and Recall. It is calculated as:

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Using Python, you can calculate these metrics using `sklearn.metrics` :

python



```
from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score
# Calculate accuracy
accuracy = accuracy_score(y_true, y_pred) # Calculate precision
precision = precision_score(y_true, y_pred) # Calculate recall
recall = recall_score(y_true, y_pred) # Calculate F1 score
f1 = f1_score(y_true, y_pred)
print(f"Accuracy: {accuracy}") print(f"Precision: {precision}") print(f"Recall: {recall}") print(f"F1 Score: {f1}")
```

Theory Behind the Confusion Matrix Metrics:

1. **Accuracy:** Accuracy is the simplest evaluation metric. It works well when the dataset is balanced but may not be reliable in cases where the dataset is imbalanced (e.g., 95% negatives, 5% positives).
2. **Precision:** Precision answers the question: *Of all the instances where the model predicted positive, how many were actually positive?* It is particularly useful when the cost of a false positive is high (e.g., falsely diagnosing a patient with a disease).
3. **Recall (Sensitivity):** Recall answers the question: *Of all the actual positive instances, how many were correctly identified by the model?* It is important when the cost of a false negative is high (e.g., failing to identify a disease in a patient who actually has it).
4. **F1 Score:** The F1 score is the harmonic mean of precision and recall. It is useful when you need a balance between precision and recall, especially in imbalanced datasets. A high F1 score indicates both high precision and high recall.

Conclusion

By following these steps, you'll be able to prepare the data, perform basic exploratory data analysis (EDA), and calculate key performance metrics using the confusion matrix. These steps are fundamental when assessing the performance of machine learning classification models.

[more theory](#)



Deep Dive: Theory Behind Data Preparation and Evaluation Metrics

To gain a clearer understanding of the theory behind data preparation and evaluation metrics like accuracy, precision, recall, and F1 score, let's break down each step and concept in more detail.

Data Preparation and Preprocessing

Data preparation is a crucial step in any data analysis or machine learning task, especially in real-world scenarios where data is often messy and incomplete. Proper preprocessing is essential for making meaningful predictions and ensuring that the model can learn effectively from the data. Here are some key concepts:

1. Shape of Data

The shape of a dataset refers to its dimensions, specifically the number of rows (observations) and columns (features). For example:

- Rows represent individual samples or records (e.g., individual patient data).
- Columns represent different features (e.g., age, sex, blood pressure, cholesterol, etc.).

Knowing the shape helps to understand the scale of the data and how much information you are working with.

2. Missing Values

Missing data is common in many datasets, and it can arise due to a variety of reasons:

- Data not collected or recorded
- Errors in data entry
- Non-response in surveys or studies

Dealing with missing values is important because many machine learning algorithms can't handle missing data directly. Common strategies include:

- **Imputation:** Filling missing values with a calculated value, such as the mean, median, or mode of the column.
- **Removal:** Deleting rows or columns that contain too many missing values (if they are not critical).
- **Prediction:** Using models to predict missing values based on other features.

3. Data Types

Different types of data require different handling:

- **Numerical data** (e.g., age, blood pressure, cholesterol levels) can be used directly in machine learning models, though they may require scaling or normalization.
- **Categorical data** (e.g., sex, chest pain type) may need to be encoded into numerical form using techniques like one-hot encoding or label encoding to be usable in most machine learning algorithms.
- **Ordinal data** (e.g., levels of severity) has a natural order and may require special handling, such as assigning numerical values based on the order.

4. Zero Values

In some contexts, zeros might be valid values (e.g., a patient with zero cholesterol). However, in other contexts, zeros may indicate missing data or errors (e.g., age or blood pressure values of zero might be unrealistic). Identifying zeros is important to understand whether they are meaningful or need to be corrected.

5. Feature Selection

Extracting relevant features (columns) is crucial for the efficiency and accuracy of your model. By selecting only the most important features (e.g., age, sex, chest pain type), you can reduce the complexity of your model and improve its interpretability.

Confusion Matrix and Performance Metrics

Once you have a trained model, the next step is to evaluate its performance. The confusion matrix provides a detailed view of how well the model is performing, and from that, we can derive important metrics that help in model evaluation.

Confusion Matrix

A **confusion matrix** is a table that shows the performance of a classification algorithm by comparing the actual versus predicted values. It breaks down the predictions into four categories:

- **True Positives (TP):** These are cases where the model correctly predicts a positive class (e.g., predicting a patient has heart disease when they actually do).
- **False Positives (FP):** These are cases where the model incorrectly predicts a positive class (e.g., predicting a patient has heart disease when they do not).
- **True Negatives (TN):** These are cases where the model correctly predicts a negative class (e.g., predicting a patient does not have heart disease when they really don't).
- **False Negatives (FN):** These are cases where the model incorrectly predicts a negative class (e.g., predicting a patient does not have heart disease when they actually do).

The confusion matrix helps to break down the performance of the model into more specific categories, rather than just providing an overall accuracy.

Performance Metrics

1. Accuracy

Accuracy is the most straightforward evaluation metric. It is the proportion of correct predictions (both positive and negative) out of all predictions made.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Pros:** Simple to calculate and understand.
- **Cons:** It may be misleading when dealing with imbalanced datasets (e.g., in fraud detection or rare disease detection, where the majority of cases are negative).

2. Precision (Positive Predictive Value)

Precision measures how many of the instances predicted as positive are actually positive. It focuses on the performance of the model in predicting the positive class.

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **High precision** means that when the model predicts a positive class, it is usually correct.
- **Use case:** Precision is particularly important when false positives are costly. For example, in medical diagnoses, a false positive may lead to unnecessary treatments or tests.

3. Recall (Sensitivity or True Positive Rate)

Recall measures how many of the actual positive instances were correctly predicted by the model. It focuses on the model's ability to identify positive cases.

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **High recall** means that the model correctly identifies most of the positive cases.
- **Use case:** Recall is especially important when false negatives are costly. For example, in cancer detection, failing to detect cancer (false negative) could have severe consequences.

4. F1 Score

The F1 score is the harmonic mean of precision and recall. It provides a balance between the two metrics, especially when there is an uneven class distribution (e.g., when one class is much rarer than the other).

prevalent than the other).

$$F1 \text{ Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- **Use case:** The F1 score is useful when you need a balance between precision and recall, especially useful when the dataset has imbalanced classes (e.g., a rare disease with only a small percentage of positive cases).

Other Evaluation Metrics

Besides the standard metrics discussed above, there are other evaluation metrics that can be valuable depending on the task and context:

1. **Specificity (True Negative Rate):** Measures how well the model identifies negative cases. It is the complement of recall:

$$\text{Specificity} = \frac{TN}{TN + FP}$$

2. **Area Under the ROC Curve (AUC-ROC):** The ROC curve is a graphical representation of a classifier's performance, plotting the true positive rate (recall) against the false positive rate. The area under the ROC curve (AUC) provides an aggregate measure of the classifier's performance across all possible thresholds.
3. **Log Loss (Cross-Entropy Loss):** This metric evaluates the probability output of a classification model rather than just the final predicted class. It is used primarily in probabilistic models (e.g., logistic regression).
4. **Matthews Correlation Coefficient (MCC):** MCC is a balanced measure that takes into account true positives, true negatives, false positives, and false negatives. It ranges from -1 (perfect inverse prediction) to 1 (perfect prediction), with 0 indicating random prediction.

Conclusion

- **Data Preparation:** The quality of your data largely determines the effectiveness of your model. Proper preprocessing, handling of missing values, feature selection, and understanding the shape and types of data are all foundational to building robust models.
- **Confusion Matrix and Evaluation Metrics:** The confusion matrix offers a detailed breakdown of a model's performance, and metrics like accuracy, precision, recall, and F1 score help evaluate different aspects of the model's effectiveness. Precision and recall are especially useful for imbalanced datasets, where accuracy alone may not provide an accurate picture of model performance.

Ultimately, choosing the right evaluation metrics depends on the problem at hand and the consequences of different types of errors. For example, in medical applications, you may prioritize recall to ensure that you detect as many positive cases as possible, even at the cost of some false positives. In other cases, you might prioritize precision to minimize the cost of false positives.



Download temperature data from below link.
<https://www.kaggle.com/venky73/temperaturesofindia?select=temperatures.csv> This data consists of temperatures of all places in India, averaging the temperatures of all places month wise. Temperature values are recorded in CELSIUS.
 a. Apply Linear Regression using suitable library function and plot the Monthwise temperature.

- b. Assess the performance of regression models using MSE, MAE, R-Square metrics
- c. Visualize simple regression model



Temperature Prediction Using Linear Regression

In this task, we will work with the temperature data of India, which includes monthly average temperatures recorded in Celsius for various places across India. Our objective is to apply **Linear Regression** to predict the monthly temperature, assess the model's performance using various evaluation metrics, and visualize the results.

Step-by-Step Breakdown

1. Download the Dataset

- First, download the temperature dataset from the given Kaggle link.
- The dataset likely contains columns like "Month" and "Temperature" (in Celsius), where temperature represents the average temperature for each month across various locations in India.

2. Data Loading and Preprocessing

- Load the dataset into a pandas DataFrame.
- Preprocess the data (check for missing values, ensure correct data types, etc.).
- Split the data into features (X) and target (Y), where the feature will be the month, and the target will be the temperature.

3. Apply Linear Regression

- Use a linear regression model to predict the monthly temperature based on the features.
- Fit the model using a suitable library function (like `sklearn.linear_model.LinearRegression`).

4. Assess Model Performance

- Evaluate the model performance using the following metrics:
 - Mean Squared Error (MSE)
 - Mean Absolute Error (MAE)
 - R-Squared (R^2)

5. Visualize the Regression Model

- Plot the predicted values against the actual values.
- Visualize the linear regression line over the data points to better understand how well the model fits the data.

Code Implementation

Let's break it down into code.

Step 1: Data Loading and Preprocessing

First, make sure you've downloaded the dataset and have the necessary libraries installed:

```
bash
```

```
pip install pandas matplotlib scikit-learn
```

Now, let's load the data and inspect it:

```
python
```

```
import pandas as pd # Load the dataset
df = pd.read_csv('temperatures.csv') # Inspect the first few rows of the dataset
print(df.head()) # Check for missing values and
```



```
types print(df.info()) # Drop any rows with missing values (if present)
df.dropna(inplace=True)
```

Assuming the dataset has a "Month" and "Temperature" column, let's move on to applying regression model.

Step 2: Apply Linear Regression

We will use `scikit-learn` to create a linear regression model. The target variable (`y`) is the **Temperature**, and the feature (`x`) is the **Month**.

python



```
from sklearn.linear_model import LinearRegression from sklearn.model_selection import train_test_split # Extract feature (X) and target variable (y) X = df[['Month']]
'Month' as the independent variable y = df['Temperature'] # 'Temperature' as the dependent variable # Split the data into training and testing sets (80% training, testing) X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) # Initialize the linear regression model model = LinearRegression() # Fit the model on the training data model.fit(X_train, y_train) # Predict the temperatures on the test set y_pred = model.predict(X_test)
```

Step 3: Model Evaluation

We will use **Mean Squared Error (MSE)**, **Mean Absolute Error (MAE)**, and **R-Squared (R^2)** to evaluate the model's performance.

python



```
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score # Calculate Mean Squared Error (MSE) mse = mean_squared_error(y_test, y_pred) # Calculate Mean Absolute Error (MAE) mae = mean_absolute_error(y_test, y_pred) # Calculate R-Squared ( $R^2$ ) r_squared = r2_score(y_test, y_pred) # Print the performance metrics print(f"Mean Squared Error (MSE): {mse}") print(f"Mean Absolute Error (MAE): {mae}") print(f"R-Squared ( $R^2$ ): {r_squared}")
```

- **MSE:** Measures the average of the squares of the errors—that is, the average squared between the actual and predicted values. Lower values are better.
- **MAE:** Measures the average magnitude of errors in a set of predictions, without considering their direction. Like MSE, lower values are better.
- **R^2 (R-Squared):** This is the proportion of the variance in the dependent variable that is predictable from the independent variable. R^2 values range from 0 to 1, where 1 indicates a perfect fit.

Step 4: Visualize the Simple Regression Model

Now, let's visualize the fitted linear regression line over the actual data points. We'll use **Matplotlib** to plot the data.

python



```
import matplotlib.pyplot as plt # Plot the actual data points plt.scatter(X_test, y_test, color='blue', label='Actual Data') # Plot the regression line (Predicted) plt.plot(X_test, y_pred, color='red', linewidth=2, label='Regression Line') # Set title and labels plt.title('Linear Regression: Temperature Prediction') plt.xlabel('Month')
```

```
plt.ylabel('Temperature (Celsius)') # Show the legend plt.legend() # Show the plot  
plt.show()
```

In this plot:

- **Blue dots** represent the actual data points (i.e., the real temperatures for each month).
- **Red line** represents the linear regression model's predictions.

Summary of Results

After running the code:

1. **Model Performance:** You'll get values for MSE, MAE, and R^2 , which help assess the performance of the model.
 - **MSE and MAE** indicate the average error in the model's predictions.
 - **R^2** tells us how well the model explains the variance in the temperature data. A high value suggests the model fits the data well.
2. **Visualizing the Model:** The regression plot will show how well the linear model fits the temperature data. A straight line indicates a linear relationship between the month and temperature.

Additional Insights

- **Data Trend:** In the real-world, temperature data typically shows some seasonality or cyclical behavior (e.g., higher temperatures in summer months and lower temperatures in winter months). A simple linear regression might not fully capture these seasonal patterns, which may require more advanced models like **Polynomial Regression** or **Time Series Models** (e.g., SARIMA).
- **Model Evaluation:** If the model's performance metrics (MSE, MAE) are high, or the R^2 is low, this suggests that a simple linear model may not be the best choice for this dataset. You could explore other regression techniques or even feature engineering (e.g., encoding months as sine/cosine values to capture seasonal behavior).