

Quick revision on FAQ's

1. Train vs Test vs Validation set:

Imagine that you are a data scientist tasked with building a machine learning model that can predict the price of a house based on its features like size, location, number of rooms, etc. Now, in order to build an accurate model, you need data to train your model on. You have a dataset of 10,000 houses with their respective features and prices.

However, before you can use this dataset to train your model, you need to split it into three different subsets: the train set, the test set, and the validation set.

- The train set is the portion of the data that you will **use to train** your machine learning model. In this example, you would use the features of the houses in the train set to teach your model how to predict the price of a house based on its features.
- The test set is a separate portion of the data that you will **use to evaluate** how well your model performs on data that it has never seen before. You want to make sure that your model is not simply memorizing the features of the houses in the train set, but rather is able to generalize to new, unseen data.
- Finally, the validation set is **used to fine-tune** your model and improve its accuracy. You can use the validation set to test different variations of your model and see which one performs best.

To split the data into these three subsets, you might randomly **assign 70% of the data to the train set, 20% to the test set, and 10% to the validation set**. This split can be adjusted depending on the size of your dataset and the complexity of your model.

For example:

Imagine that you have a **dataset of 1000** images of dogs and cats that you want to use to train a machine learning model to classify images as either "dog" or "cat". You would split this dataset into a train set of **700 images**, a test set of **200 images**, and a validation set of **100 images**.

You would use the 700 images in the train set to teach your model how to recognize features that distinguish dogs from cats. You would then use the test set of 200 images to evaluate how well your model performs on data that it has never seen before. Finally, you would use the validation set of 100 images to fine-tune your model and improve its accuracy.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

2. Data Leakage:

Data leakage occurs when information from outside of the training dataset is used in the model training process, leading to overestimation of the model's performance and generalization error on new data. In other words, **data leakage happens when there is unintentional mixing of training and testing data**, or when the model is trained on data that includes information that would not be available at the time of deployment.

There are **several ways in which data leakage** can occur, including:

- ✓ Using information from the test set during the model training process, such as **using labels from the test set** to improve the model's accuracy. This can result in overfitting and poor generalization on new data.
- ✓ **Using data that is not representative** of the real-world data that the model will be applied to. For example, if you're building a model to predict customer churn, using data from a specific time period when a company had a promotion that kept customers from leaving would lead to overestimating the model's accuracy on new data.
- ✓ **Including features in the model** that are derived from the target variable. For example, if you're building a model to predict the price of a house, including the exact address of the house as a feature would lead to overfitting, as the model would simply memorize the prices of specific houses.
- ✓ **Using data** that is not independent and identically **distributed (IID)**. If the training data is not IID, then the model will not generalize well to new data. For example, if you're building a model to predict credit card fraud, including data from only one region of the world would lead to a model that is biased towards that region and would not generalize well to other regions.

Data leakage can have severe consequences, such as underestimating risk, overestimating the accuracy of the model, and in some cases, legal or ethical issues.

Examples:

- ✓ One example of data leakage is when a credit card company builds a model to predict fraudulent transactions. If the company uses data that includes transaction amounts and timestamps, the model may learn to recognize patterns based on the transaction amounts and times, leading to overfitting and poor generalization to new data. Instead, the company should remove these features before training the model.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

- ✓ Another example of data leakage is in medical research. If researchers use data from the same patients for both the training and testing sets, the model may simply memorize the data of those specific patients rather than learning generalizable patterns. To prevent this, the researchers should use data from different patients for the training and testing sets.

In summary, data leakage can occur in various ways and can have severe consequences on the accuracy and generalization of a model. It is essential to be aware of the potential sources of data leakage and take appropriate measures to prevent it.

3. Why do we prefer Random Forest than Decision Tree?

Random Forest is a popular ensemble learning algorithm that is often preferred over decision trees. Here are some reasons why:

- 1) **Reduces overfitting:** Decision trees are prone to overfitting, meaning they can memorize the training data instead of generalizing to new data. Random Forest reduces overfitting by creating multiple trees and averaging their predictions, which reduces the variance and leads to a more stable model.
- 2) **Handles missing data:** Decision trees cannot handle missing data well and can lead to biased results. Random Forest can handle missing data by using mean imputation or using surrogate splits.
- 3) **Reduces bias:** Decision trees can be biased towards features that have more levels or values. Random Forest reduces this bias by randomly selecting a subset of features for each tree.
- 4) **Reduces the effect of outliers:** Decision trees can be sensitive to outliers, leading to incorrect splits. Random Forest reduces the effect of outliers by averaging the predictions of multiple trees.
- 5) **Performs well with large datasets:** Decision trees can become computationally expensive with large datasets. Random Forest can handle large datasets by parallelizing the tree-building process.

Here are some use cases and examples where Random Forest is preferred over Decision Trees:

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

- 1) **Predicting customer churn:** In a telecommunications company, a Random Forest model can be used to predict which customers are likely to churn based on features such as customer demographics, service usage, and billing information.
- 2) **Fraud detection:** In finance, a Random Forest model can be used to identify fraudulent transactions based on features such as transaction amount, location, and time.
- 3) **Medical diagnosis:** In healthcare, a Random Forest model can be used to diagnose a medical condition based on patient symptoms and medical history.
- 4) **Image classification:** In computer vision, a Random Forest model can be used to classify images based on features such as color, texture, and shape.
- 5) **Predicting stock prices:** In finance, a Random Forest model can be used to predict stock prices based on features such as company financials, market trends, and economic indicators.

Overall, Random Forest is a powerful algorithm that can improve upon the limitations of Decision Trees. Its ability to reduce overfitting, handle missing data, reduce bias, reduce the effect of outliers, and perform well with large datasets makes it a popular choice in various industries and use cases.

4. Confusion Matrix:

A confusion matrix is a table that is used to evaluate the performance of a classification model. It shows the number of true positives, true negatives, false positives, and false negatives for each class in the dataset. The four terms in the confusion matrix are defined as follows:

True positives (TP): The number of positive cases that are correctly classified as positive.

True negatives (TN): The number of negative cases that are correctly classified as negative.

False positives (FP): The number of negative cases that are incorrectly classified as positive.

False negatives (FN): The number of positive cases that are incorrectly classified as negative.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

		Predicted Class		
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	Specificity $\frac{TN}{(TN + FP)}$
		Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$

The confusion matrix is often used to calculate various performance metrics for a classification model such as **accuracy, precision, recall, F1-score, and ROC curve**. Here is an example of a confusion matrix for a binary classification problem.

Let's consider a real-world example of a spam email classifier:

Suppose we have a dataset of 1000 emails, out of which 200 are spam and 800 are legitimate. We build a binary classification model to classify these emails as spam or not spam. After training the model, we apply it to a test dataset of 200 emails, where 40 of them are spam and 160 are legitimate.

Suppose the model correctly predicted 30 spam emails as spam (TP), incorrectly predicted 10 legitimate emails as spam (FP), correctly predicted 150 legitimate emails as legitimate (TN), and incorrectly predicted 10 spam emails as legitimate (FN).

Now we can use this confusion matrix to calculate various performance metrics as follows:

- Accuracy:** $(TP + TN) / (TP + TN + FP + FN) = (30 + 150) / (30 + 10 + 150 + 10) = 0.90$ or 90%

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

This means that the model correctly classified 90% of the emails in the test dataset.

- ii. **Precision:** $TP / (TP + FP) = 30 / (30 + 10) = 0.75$ or 75%
This means that out of all the emails that the model predicted as spam, 75% were actually spam.
- iii. **Recall (or sensitivity):** $TP / (TP + FN) = 30 / (30 + 10) = 0.75$ or 75%
This means that out of all the actual spam emails, the model correctly identified 75% of them as spam.
- iv. **F1-score:** $2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall}) = 2 * (0.75 * 0.75) / (0.75 + 0.75) = 0.75$ or 75%
This is the harmonic mean of precision and recall, and it provides a balanced measure of the model's performance.
- v. **Specificity:** $TN / (TN + FP) = 150 / (150 + 10) = 0.94$ or 94%
This means that out of all the actual legitimate emails, the model correctly identified 94% of them as legitimate.
- vi. **False Positive Rate (FPR):** $FP / (FP + TN) = 10 / (10 + 150) = 0.06$ or 6%
This means that out of all the actual legitimate emails, the model incorrectly identified 6% of them as spam.

Few more use-cases for confusion matrix:

1. Medical diagnosis:

A confusion matrix can be used to evaluate the performance of a machine learning model that is designed to diagnose diseases such as cancer.

- **True positives** represent the number of patients correctly diagnosed with the disease.
- **False positives** represent the number of healthy patients incorrectly diagnosed with the disease.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

- **False negatives** represent the number of sick patients who were incorrectly diagnosed as healthy.
- **True negatives** represent the number of healthy patients correctly diagnosed as healthy.

2. Customer churn prediction:

Confusion matrix can be used to evaluate the performance of a machine learning model that is designed to predict whether a customer will churn or not.

- **True positives** represent the number of customers who actually churned and were correctly predicted to churn.
- **False positives** represent the number of customers who did not churn but were incorrectly predicted to churn.
- **False negatives** represent the number of customers who actually churned but were incorrectly predicted to not churn.
- **True negatives** represent the number of customers who did not churn and were correctly predicted to not churn.

In summary, the confusion matrix and various performance metrics calculated from it can give us a better understanding of the strengths and weaknesses of the spam email classifier.

By analyzing the metrics, we can identify areas for improvement and make informed decisions on how to improve the accuracy of the model.

5. Suppose you are given a dataset, & you understood this can be solved by regression, which regression algo you will use & why?

Choosing the right regression algorithm depends on several factors such as:

- ✓ The nature of the dataset.
- ✓ The number and type of features, and the problem statement.

Here are some steps to help you decide which algorithm to use for regression problems:

- 1) **Understand the problem statement:** First, you need to understand the problem statement and the type of regression task you are dealing with. Are you trying to predict a continuous

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

target variable or a categorical one? Is the relationship between the target variable and the features **linear or nonlinear**?

- 2) **Check for linearity and outliers:** If the relationship between the target variable and the features is linear, linear regression may be a good choice. However, if there are nonlinear relationships, polynomial regression or decision tree regression may be more appropriate. If there are outliers in the data, robust regression algorithms such as **Ridge or Lasso regression** may be useful.
- 3) **Check for multicollinearity:** If there are high correlations between the features, **Ridge regression or Lasso regression** may be useful to reduce the variance and improve the model's performance.
- 4) **Consider the number of features:** If there are a large number of features, dimensionality reduction techniques such as **PCA** or feature selection methods such as **Lasso regression** can be used to select the most relevant features.
- 5) **Check for heteroscedasticity:** If the variance of the residuals is not constant across the range of values of the target variable, weighted regression methods such as **WLS or GLS** may be useful.
 - ✓ **GLS (Generalized Least Squares) and WLS (Weighted Least Squares)** are both regression methods used to deal with heteroscedasticity (unequal variances) in the data. However, they differ in their assumptions and implementation.
 - ✓ **WLS** assumes that the variance of the errors is not constant and can vary across the range of the target variable. It assigns weights to each observation based on the variance of the residuals
 - ✓ **GLS** assumes that the variance of the errors follows a specific pattern, which can be modeled using a covariance matrix. It takes into account the correlation between the errors and assigns weights to each observation based on the covariance matrix.
 - ✓ WLS and GLS are both useful regression methods for dealing with heteroscedasticity in the data.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

- 6) **Consider the size of the dataset:** If the dataset is small, simpler models such as linear regression or decision tree regression may be preferred. If the dataset is large, more **complex models** such as **SVR or Random Forest regression** may be used.

In summary, choosing the right regression algorithm depends on the characteristics of the dataset and the problem statement. It is important to consider factors such as **linearity, multicollinearity, outliers, heteroscedasticity, number of features, and dataset size** when selecting an appropriate algorithm.

6) What is pre-pruning & post-pruning? Why it is used?

Decision trees are prone to overfitting, which occurs when the model captures the noise in the training data instead of the underlying patterns. Pre-pruning and post-pruning are two techniques used to prevent overfitting in decision trees.

- ✓ **Pre-pruning** refers to stopping the tree construction process before it reaches the point where the model becomes overfitted.
 - In pre-pruning, **the tree is stopped** from growing beyond a certain level or depth, or if the number of instances in a node falls below a certain threshold. This ensures that the tree is not too complex and captures only the most important patterns in the data. Pre-pruning can reduce the computational cost of building a decision tree and can result in faster and more accurate predictions.
- ✓ **Post-pruning**, on the other hand, involves **growing the tree to its full depth** and then removing or collapsing nodes that do not provide any significant improvement in accuracy.
 - In post-pruning, the **tree is grown using an algorithm such as ID3 or C4.5**, and then the nodes that result in overfitting are pruned. The algorithm uses statistical tests or validation data to determine which nodes to prune. **Post-pruning can result in more accurate and robust models and can be useful when pre-pruning results in underfitting.**

The choice of pre-pruning or post-pruning depends on the nature of the dataset and the problem statement.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

- ✓ **Pre-pruning is preferred when the dataset is large and complex**, and when there is a high risk of overfitting.
- ✓ **Post-pruning is preferred when the dataset is small** or when pre-pruning results in underfitting.

In summary, pre-pruning and post-pruning are two techniques used to prevent overfitting in decision trees. Pre-pruning involves stopping the tree construction process before it becomes overfitted, while post-pruning involves growing the tree to its full depth and then removing nodes that do not improve accuracy. The choice of pre-pruning or post-pruning depends on the nature of the dataset and the problem statement.

7) What are Ensemble Methods? How it is used?

Ensemble methods are a type of machine learning technique that combines the predictions of multiple models to improve the overall accuracy and robustness of the final model.

The basic idea behind ensemble methods is that by combining the predictions of multiple models, the resulting model will be more accurate and less prone to overfitting than any of the individual models.

There are two main types of ensemble methods:

- i. Bagging
- ii. Boosting.

- ✓ **In bagging**, multiple **models are trained independently on random subsets of the data**, and their predictions are averaged or combined in some way to produce the final prediction. Bagging is **useful** when the **base models are unstable or have high variance**, as it reduces the variance and produces a more stable and robust model.
- ✓ **In boosting**, on the other hand, **multiple weak models are trained sequentially**, with each model trying to improve on the errors of the previous model. Boosting is useful when the base models are relatively simple and stable, but have high bias. Boosting can help to **reduce the bias** and produce a more accurate and robust model.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

Some of the popular ensemble methods include **Random Forest, Gradient Boosting Machines (GBM), XG-Boost, and AdaBoost.**

- Ensemble methods are used in a wide range of applications, including classification, regression, and anomaly detection.
- They are particularly useful when dealing with complex and high-dimensional data, and when the goal is to achieve high accuracy and robustness in the final model.

Some use cases of ensemble methods include:

- ✓ Predicting customer churn in telecom industry
- ✓ Fraud detection in finance industry
- ✓ Identifying cancer patients based on gene expression data
- ✓ Predicting stock prices and market trends in finance industry

In summary, ensemble methods are a powerful machine learning technique that can help to improve the accuracy and robustness of the final model.

Bagging and boosting are the two main types of ensemble methods, and popular algorithms include Random Forest, GBM, XGBoost, and AdaBoost. Ensemble methods are used in a wide range of applications and are particularly useful when dealing with complex and high-dimensional data.

8) How ensemble methods are used in Random Forest, Gradient Boosting Machines (GBM), XG-Boost, and AdaBoost?

Ensemble methods are used differently in different algorithms. Here's a brief overview of how ensemble methods are used in some popular algorithms:

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

i. Random Forest:

- ✓ Random Forest is an ensemble learning method that uses bagging as the base method.
- ✓ In Random Forest, multiple decision trees are trained on different subsets of the data, and their predictions are combined by taking the average (for regression) or majority vote (for classification) to produce the final prediction.
- ✓ Random Forest also uses feature bagging, where a random subset of features is selected for each tree to further improve the robustness and reduce overfitting.

ii. Gradient Boosting Machines (GBM):

- ✓ GBM is a boosting algorithm that trains a sequence of weak models (usually decision trees) sequentially, with each model trying to improve on the errors of the previous model.
- ✓ In GBM, the objective function is minimized by adding new trees to the model, with each new tree fitting the residual errors of the previous trees. The final prediction is the sum of the predictions of all the trees.

iii. XGBoost:

- ✓ XGBoost is a variant of GBM that uses a more regularized model to prevent overfitting. It also uses a technique called "weighted quantile sketch" to handle large-scale data efficiently.
- ✓ XGBoost is known for its speed and performance and has been used to win several machine learning competitions.

iv. AdaBoost:

- ✓ AdaBoost is a boosting algorithm that assigns higher weights to the misclassified samples in each iteration, which makes the subsequent models focus on the harder examples.
- ✓ In AdaBoost, a series of weak models (usually decision stumps) are trained sequentially, with each model trying to improve on the errors of the previous model.
- ✓ The final prediction is a weighted sum of the predictions of all the models.

In summary, ensemble methods are used in different ways in different algorithms, but the general idea is to combine the predictions of multiple models to improve the overall accuracy and robustness of the final model.

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

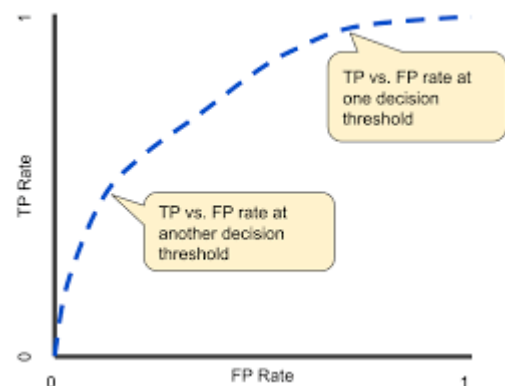
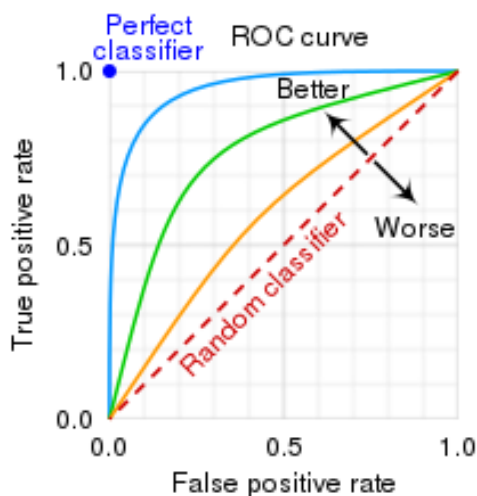
9) What is ROC-AUC curve? What does it signify? How to study ROC-AUC curve to get the insights?

ROC (Receiver Operating Characteristic) curve is a graphical representation of the performance of a binary classification model at different classification thresholds. It plots

- ✓ The **True Positive Rate (TPR)** against the **False Positive Rate (FPR)** for different threshold values. The area under the ROC curve (AUC) is a metric that measures the performance of the model across all possible classification thresholds.
- ✓ The **ROC-AUC** curve can be used to evaluate the overall performance of a binary classification model. A model with an **AUC of 1.0** is a perfect classifier, while a model with an **AUC of 0.5** is no better than **random guessing**.

To study the ROC-AUC curve, you can follow these steps:

- i. Train a binary classification model and obtain the predicted probabilities for the test set.
- ii. Plot the ROC curve by plotting TPR against FPR for different classification thresholds. The diagonal line represents a random classifier.
- iii. Compute the AUC score by calculating the area under the ROC curve.
- iv. Interpret the ROC-AUC curve to gain insights into the performance of the model. For example, a higher AUC score indicates better performance, and the shape of the curve can provide information on the balance between sensitivity (TPR) and specificity (1 - FPR).



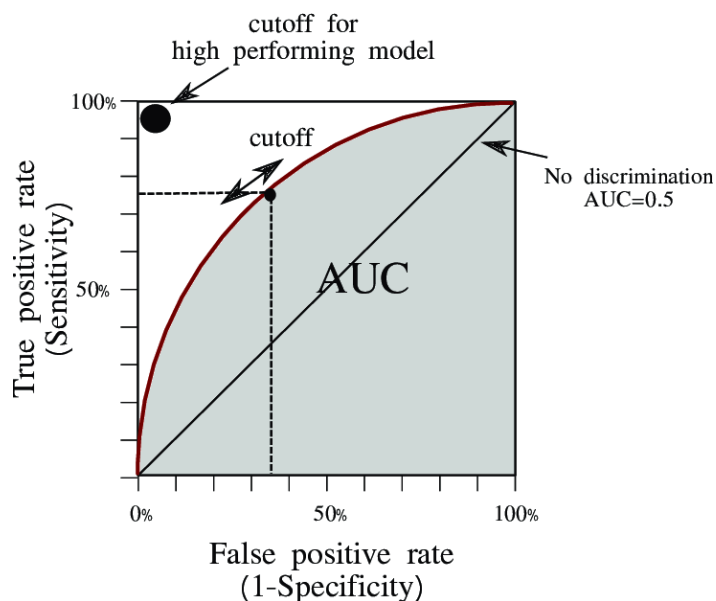
Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

Here are some key insights you can gain from studying the ROC-AUC curve:

- i. **Discrimination threshold:** The ROC curve can help you choose an appropriate classification threshold that balances sensitivity and specificity based on the problem's specific requirements.
- ii. **Imbalanced data:** The ROC-AUC curve can help you evaluate the performance of the model on imbalanced data by providing a summary measure of the overall model performance across different classification thresholds.
- iii. **Comparison between models:** You can use the ROC-AUC curve to compare the performance of different models on the same dataset.
- iv. **Interpretation of model performance:**
The ROC-AUC curve can provide insights into how well the model is distinguishing between positive and negative samples and how sensitive it is to false positives.



In summary, the ROC-AUC curve is a useful tool for evaluating and **interpreting the performance of binary classification models**. It provides a summary measure of the model's performance across different classification thresholds and can help you choose an appropriate threshold for the problem at hand.

10) When to use which performance metrics should be used, & why?

- a) **Binary classification problems:** When the problem involves classifying data into two categories, common performance metrics include **accuracy, precision, recall, F1-score, and ROC-AUC curve**. The choice of metric depends on the specific problem requirements, such as the importance of false positives versus false negatives.
- b) **Multi-class classification problems:** When the problem involves classifying data into more than two categories, common performance metrics include **accuracy, macro-**

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>

averaged precision, macro-averaged recall, macro-averaged F1-score, and confusion matrix.

- c) **Regression problems:** When the problem involves predicting a continuous outcome, common performance metrics include mean squared error (**MSE**), root mean squared error (**RMSE**), mean absolute error (**MAE**), **R-squared**, and coefficient of determination.
- d) **Imbalanced data:** When the data is imbalanced, it may be more appropriate to use metrics such as **precision, recall, and F1-score instead of accuracy, as accuracy can be misleading when the classes are imbalanced.**
- e) **Cost-sensitive learning:** When the cost of different types of errors is different, it may be necessary to use metrics that take into account the cost of different types of errors, such as **weighted precision and recall.**

Name: Shobhandeb Paul

Linkedin: <https://www.linkedin.com/in/shobhandeb-paul-b6914a168/>

Github: <https://github.com/herbert0419>