Assignment

Q1. Describe the decision tree classifier algorithm and how it works to make predictions. Ans:A decision tree classifier is a popular machine learning algorithm used for both classification and regression tasks. It works by recursively partitioning the input space into regions and assigning a specific class label (for classification) or predicting a continuous value (for regression) within each region. Decision trees are constructed based on a set of rules that are learned from the training data.

Here are the main steps in the decision tree classifier algorithm:

1. Decision Tree Construction:

a. Root Node:

- Selection: Choose the best feature to split the data based on a criterion (discussed later).
- Split: Partition the dataset into subsets based on the selected feature.

b. Internal Nodes:

- Repeat: Recursively apply the splitting process to each subset (child node).
- Selection: At each internal node, choose the best feature to split the data among its child nodes.

c. Leaf Nodes:

 Assignment: Assign a class label to each leaf node based on the majority class in the corresponding subset.

2. Feature Selection:

a. Information Gain (for Entropy-based Trees):

- Concept: Measures the reduction in entropy (uncertainty) achieved by splitting the data based on a particular feature.
- Calculation:
- Information

Gain=Entropy(parent) $-\sum$ (Weight(child)×Entropy(child)Weight(parent))

- Information Gain=Entropy(parent)−∑(
 - Weight(parent)
 - Weight(child)×Entropy(child)

•

b. Gini Impurity (for Gini-based Trees):

- Concept: Measures the probability of misclassifying an instance.
- Calculation:
- Gini Impurity= $1-\sum$ (Probability of Class2)
- Gini Impurity= $1-\sum$ (Probability of Class
- 2
-)

c. Other Criteria (e.g., Gain Ratio, Chi-square):

 Depending on the specific implementation or library, other criteria may be used for feature selection.

3. Stopping Criteria:

- Max Depth: Limit the depth of the tree to prevent overfitting.
- Min Samples Split: Specify the minimum number of samples required to split a node.
- Min Samples Leaf: Specify the minimum number of samples required to be in a leaf node.
- Max Features: Limit the number of features considered for a split.

4. Prediction:

- Traversal: Given a new instance, traverse the decision tree from the root node to a leaf node based on the feature values of the instance.
- Class Assignment: Assign the class label associated with the leaf node.

Advantages of Decision Trees:

Interpretability:

 Decision trees are easy to understand and interpret, making them suitable for explaining model decisions.

Handling Non-Linearity:

• Decision trees can capture non-linear relationships in the data without explicit feature engineering.

Variable Importance:

• Feature importance can be easily determined, providing insights into which features contribute most to predictions.

Robust to Outliers:

• Decision trees are less sensitive to outliers compared to some other algorithms.

Challenges of Decision Trees:

Overfitting:

• Decision trees can easily overfit the training data, capturing noise rather than the underlying patterns.

Instability:

 Small changes in the data can lead to different tree structures, making decision trees less stable.

Bias Towards Dominant Classes:

- In imbalanced datasets, decision trees may be biased towards dominant classes. Global Optimization:
 - Decision trees use a greedy approach for local optimization, which may not lead to the best global tree structure.

Despite these challenges, decision trees are widely used and form the basis for ensemble methods like Random Forests and Gradient Boosting, which address some of the limitations of individual decision trees.

Q2. Provide a step-by-step explanation of the mathematical intuition behind decision tree classification.

Ans:The mathematical intuition behind decision tree classification involves concepts such as entropy, information gain, and Gini impurity. Here, I'll provide a step-by-step explanation of these concepts and how they are used in the decision tree classification algorithm:

1. Entropy:

- Entropy measures the uncertainty or disorder in a set of data.
- For a binary classification problem, the entropy (
- �
- *H*) of a set
- •
- S with respect to the class labels
- 🛭
- *p* and
- •
- q is calculated as follows:
- $\lozenge(\lozenge) = -\lozenge \cdot \log 2(\lozenge) \lozenge \cdot \log 2(\lozenge)$

- $H(S) = -p \cdot \log$
- 2
- •
- $(p)-q \cdot \log$
- 2
- •
- (q)
 - where
- �
- *p* and
- •
- q represent the probabilities of the positive and negative classes in the set.
- The goal is to minimize entropy by finding the best split that reduces uncertainty.

2. Information Gain:

- Information gain measures the reduction in entropy achieved by splitting the data based on a particular feature.
- Given a feature
- �
- ullet A and its possible values
- {**\dagge**1,**\dagge**2,...,**\dagge**\dagge}
- {a
- 1
- •
- ,a
- 2
- •
- ,...,a
- *k*
- •
- }, the information gain (
- ��
- *IG*) is calculated as:
- $\bullet \quad \diamondsuit \diamondsuit (\diamondsuit, \diamondsuit) = \diamondsuit (\diamondsuit) \sum \diamondsuit = 1 \diamondsuit |\diamondsuit \diamondsuit| |\diamondsuit| \cdot \diamondsuit (\diamondsuit \diamondsuit)$
- $IG(S,A)=H(S)-\sum$
- *i*=1
- *k*

- | S
- *i*
- •
- · *H*(*S*
- , ;
- •
-)

where

- ��
- 5
- *i*
- •
- is the subset of
- •
- ullet S when feature
- �
- A takes the value
- ��
- a
- i
- _
- , and
- | ��|
- |S
- ;
- •
- is the size of
- ��
- S
- i
- •
- .
- Select the feature that maximizes information gain as the best feature for splitting.

3. Gini Impurity:

- Gini impurity measures the probability of misclassifying an instance.
- For a set
- �
- *S*, the Gini impurity (
- �
- G) is calculated as:
- $\Diamond(\Diamond)=1-\sum \Diamond=1 \Diamond \Diamond \Diamond 2$
- $G(S)=1-\sum$
- *i*=1
- k
- •
- p
- *i*
- 2
- •
- where
- ��
- n
- i
- _
- is the probability of an instance belonging to class
- �
- *i* in the set
- �
- S.
- Similar to entropy, the goal is to minimize Gini impurity by finding the best split.

4. Splitting:

- The decision tree algorithm iteratively selects the best feature to split the data based on the criterion (entropy or Gini impurity).
- For each split, the dataset is partitioned into subsets, and the process is repeated recursively for each subset.

5. Stopping Criteria:

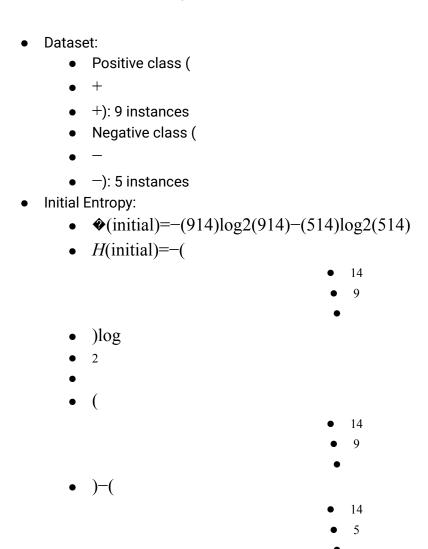
- Stopping criteria, such as maximum depth, minimum samples per split, or minimum samples per leaf, are used to prevent overfitting.
- If a stopping criterion is met, the algorithm stops recursion and assigns a class label to the leaf node.

6. Prediction:

- Given a new instance, it traverses the decision tree from the root to a leaf based on the feature values.
- The class label associated with the leaf node is assigned to the instance.

Example:

Let's consider a simple binary classification problem:



-)log
- •
- (

145

-)
- Information Gain:
 - Calculate information gain for each feature, select the one with the highest gain.
- Splitting:
 - Split the dataset based on the selected feature.
- Gini Impurity:
 - Calculate Gini impurity for each split.
- Recursive Process:
 - Repeat the process for each subset until stopping criteria are met.
- Prediction:
 - Assign class labels based on the leaf nodes.

By iteratively applying these steps, the decision tree algorithm constructs a tree that optimally partitions the data based on the selected features, reducing uncertainty and making predictions for new instances.

Q3. Explain how a decision tree classifier can be used to solve a binary classification problem. Ans:A decision tree classifier can be used to solve a binary classification problem by creating a tree structure that recursively splits the dataset based on different features, eventually assigning a class label to each region of the input space. Here's a step-by-step explanation of how a decision tree classifier works for binary classification:

1. Initial Dataset:

一).

Assume you have a dataset with instances belonging to one of two classes: positive	(
+	
+) or negative (
_	

2. Entropy or Gini Impurity Calculation:

- Calculate the initial entropy or Gini impurity of the entire dataset.
 - For entropy:
 - ϕ (initial)= $-\phi \cdot \log 2(\phi) \phi \cdot \log 2(\phi)$
 - $H(initial) = -p \cdot log$
 - 2
 - •
 - $(p)-q \cdot \log$
 - 2
 - •
 - (q) where
 - . 0
 - \bullet p and
 - �
 - q are the probabilities of positive and negative instances.

3. Feature Selection:

- Identify the feature that provides the maximum information gain or minimizes the Gini impurity when used for splitting.
 - Information Gain:
 - $\Diamond \Diamond (\Diamond, \Diamond) = \Diamond (\Diamond) \sum \Diamond = 1 \Diamond |\Diamond \Diamond | |\Diamond | \cdot \Diamond (\Diamond \Diamond)$
 - $IG(S,A)=H(S)-\sum$
 - *i*=1
 - *k*
 - •

• |S|

- |S
- i
- •

• |

- $\cdot H(S)$
- i
-)
- Gini Impurity:
- $\Diamond(\Diamond)=1-\sum \Diamond=1 \Diamond \Diamond \Diamond 2$

- $G(S)=1-\sum$
- i=1
- k
- •
- p
- *i*
- 2
- •
- •

4. Splitting the Dataset:

- Split the dataset into subsets based on the chosen feature.
 - Each subset corresponds to a branch in the decision tree.

5. Recursive Splitting:

- Repeat the process recursively for each subset, selecting the best feature for splitting at each internal node.
 - Continue until a stopping criterion is met (e.g., maximum depth, minimum samples per split).

6. Leaf Nodes and Class Assignment:

- When the recursive splitting process reaches a leaf node, assign a class label based on the majority class in that subset.
 - For example, if most instances in a leaf node are positive, assign the positive class label.

7. Prediction:

- Given a new instance, traverse the decision tree from the root to a leaf node based on the feature values of the instance.
- Assign the class label associated with the leaf node to the instance.

Example:

Let's consider a simple binary classification problem where we want to predict whether a person will buy a product (

```
+
+) or not (
-
-) based on two features: Age and Income.
```

```
Initial Dataset: Instances with labeled classes (
+
+ or
-
-).
```

Entropy Calculation: Calculate the entropy of the entire dataset.

Feature Selection: Choose the feature that provides the maximum information gain or minimizes Gini impurity.

Splitting the Dataset: Divide the dataset into subsets based on the chosen feature. Recursive Splitting: Repeat the process for each subset until a stopping criterion is met. Leaf Nodes and Class Assignment: Assign class labels to leaf nodes based on majority class in each subset.

Prediction: Given a new instance, traverse the tree to reach a leaf node and assign the associated class label.

The resulting decision tree serves as a predictive model that can be used to classify new instances into one of the two classes. The decision tree captures the decision boundaries based on the selected features, making it a powerful and interpretable tool for binary classification problems.

Q4. Discuss the geometric intuition behind decision tree classification and how it can be used to make predictions.

Ans:The geometric intuition behind decision tree classification involves the partitioning of the input space into regions, each associated with a specific class label. The decision boundaries are axis-aligned, and the tree structure can be visualized as a series of splits along the feature axes. Let's explore this geometric intuition and how it facilitates predictions:

Geometric Intuition:

Decision Boundaries:

 Decision tree classification creates regions in the input space, each associated with a particular class label. These regions are separated by decision boundaries that are perpendicular to the feature axes.

Axis-Aligned Splits:

 Decision trees make splits along individual features, creating decision boundaries that are axis-aligned. Each split divides the input space into two regions based on a threshold value for a specific feature.

Recursive Partitioning:

• The recursive nature of decision tree construction involves further partitioning each region into subregions. This process continues until a stopping criterion is met (e.g., maximum depth, minimum samples per leaf).

Leaf Nodes:

• The terminal nodes of the tree, known as leaf nodes, represent the final decision regions. Each leaf node is associated with a majority class within its region.

Making Predictions:

Traversal through the Tree:

 To make predictions for a new instance, traverse the decision tree from the root node to a leaf node based on the feature values of the instance.

Decision at Each Node:

 At each internal node, make a decision based on the feature value. If the feature value satisfies the condition, move to the left child; otherwise, move to the right child.

Reaching a Leaf Node:

• Continue traversing until reaching a leaf node. The class label associated with that leaf node is the prediction for the new instance.

Example:

Consider a simple binary classification problem with two features, Age and Income, predicting whether a person will buy a product or not. The decision tree may have splits like:

- Root Node: Split based on Age <= 30.
 - If true, move to the left child; if false, move to the right child.
- Left Child: Further split based on Income <= \$50,000.
 - Continue the process until reaching leaf nodes.

Visually, the decision tree creates rectangular regions in the Age-Income space, each associated with a class label. The splits along the Age and Income axes form axis-aligned decision boundaries.

Benefits of Geometric Intuition:

Interpretability:

• The geometric interpretation of decision trees is intuitive and easy to understand. Decision boundaries are visualized as splits along features.

Visualization:

 Decision trees can be visually represented, allowing practitioners and stakeholders to interpret and explain the model easily.

Feature Importance:

• The splits in decision trees inherently indicate the importance of different features in making predictions.

Insights into Decision Logic:

• The path from the root to a leaf node provides insights into the decision logic followed by the model for a specific instance.

Handling Non-Linearity:

• Decision trees can capture complex, non-linear decision boundaries without requiring explicit feature engineering.

While decision trees have interpretability and simplicity advantages, they may struggle with capturing certain complex relationships and might be prone to overfitting. Ensemble methods like Random Forests or Gradient Boosted Trees are often used to address these limitations while retaining the geometric intuition of decision trees.

Q5. Define the confusion matrix and describe how it can be used to evaluate the performance of a

classification model.

Ans:A confusion matrix is a table that is used to evaluate the performance of a classification model. It provides a comprehensive view of the model's predictions by breaking down the outcomes into different categories. In a binary classification scenario, the confusion matrix has four entries: true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). These elements are used to calculate various performance metrics. Here's a breakdown of the confusion matrix elements:

Confusion Matrix Elements:

True Positives (TP):

• Instances that belong to the positive class and are correctly predicted as positive by the model.

True Negatives (TN):

• Instances that belong to the negative class and are correctly predicted as negative by the model.

False Positives (FP):

• Instances that belong to the negative class but are incorrectly predicted as positive by the model (Type I error).

False Negatives (FN):

• Instances that belong to the positive class but are incorrectly predicted as negative by the model (Type II error).

Structure of a Confusion Matrix:

mathematica

Copy code

Actual Positive Actual Negative Predicted Positive TP FP Predicted Negative FN TN

Performance Metrics Derived from Confusion Matrix:

Accuracy:

- The proportion of correctly classified instances out of the total instances.
- Accuracy=��+���+��+��+��
- Accuracy=

- $\bullet \quad \textit{TP+FP+FN+TN}$
 - TP+TN

•

Precision (Positive Predictive Value):

- The proportion of true positive predictions among all instances predicted as positive.
- Precision=���+��
- Precision=

- TP+FP
- TP

_

Recall (Sensitivity, True Positive Rate):

- The proportion of true positive predictions among all actual positive instances.
- Recall=���+��
- Recall=

TP+FNTP

•

Specificity (True Negative Rate):

- The proportion of true negative predictions among all actual negative instances.
- Specificity=���+��
- Specificity=

TN+FPTN

•

F1 Score:

- The harmonic mean of precision and recall. It balances precision and recall.
- • 1 Score=2×Precision×RecallPrecision+Recall
- F1 Score=

Precision+Recall

2×Precision×Recall

•

•

Use of Confusion Matrix:

- Model Evaluation:
 - The confusion matrix provides a detailed breakdown of a model's performance, allowing practitioners to understand how well the model is performing on different types of predictions.
- Decision-Making:
 - It helps in making informed decisions based on the balance between false positives and false negatives, depending on the specific application.
- Threshold Tuning:
 - Adjusting the classification threshold can impact the balance between precision and recall, and the confusion matrix helps in understanding the trade-offs.
- Identifying Model Limitations:
 - By examining the confusion matrix, one can identify patterns of misclassifications and gain insights into areas where the model may struggle.

• Performance Monitoring:

Periodically updating and analyzing the confusion matrix helps in monitoring

model performance over time and making necessary adjustments.

In summary, the confusion matrix is a valuable tool for assessing the performance of a

classification model. It provides a detailed breakdown of prediction outcomes, allowing

practitioners to choose appropriate performance metrics based on their specific goals and

requirements.

Q6. Provide an example of a confusion matrix and explain how precision, recall, and F1 score

can be

calculated from it.

Ans:Let's consider an example of a confusion matrix for a binary classification problem.

Suppose we have a model that predicts whether an email is spam (+) or not spam (-). The

confusion matrix might look like this:

mathematica

Copy code

Actual Spam Actual Not Spam

Predicted Spam 85 15

Predicted Not Spam 10 290

In this confusion matrix:

• True Positives (TP) = 85: Emails correctly predicted as spam.

• True Negatives (TN) = 290: Emails correctly predicted as not spam.

• False Positives (FP) = 15: Emails predicted as spam but actually not spam.

• False Negatives (FN) = 10: Emails predicted as not spam but actually spam.

Now, let's calculate precision, recall, and F1 score based on these values:

Precision:

Precision=True PositivesTrue Positives+False Positives

Precision=

True Positives+False Positives

True Positives

Precision=8585+15=85100=0.85

Precision=

85+15

85

=

100

85

=0.85

Recall (Sensitivity, True Positive Rate):

Recall=True PositivesTrue Positives+False Negatives

Recall=

True Positives+False Negatives

True Positives

Recall=8585+10=8595≈0.8947

Recall=

85+10

85

≈0.8947

F1 Score:

F1 Score=2×Precision×RecallPrecision+Recall

F1 Score=

Precision+Recall

2×Precision×Recall

F1 Score=2×0.85×0.89470.85+0.8947≈0.8727

F1 Score=

0.85+0.8947

 $2 \times 0.85 \times 0.8947$

 ≈ 0.8727

So, in this example:

- Precision is 0.85, indicating that out of all instances predicted as spam, 85% were actually spam.
- Recall is approximately 0.8947, indicating that out of all actual spam instances, the model correctly identified around 89.47%.
- F1 Score is approximately 0.8727, providing a balance between precision and recall.

These metrics help in assessing different aspects of the model's performance, and the choice between precision and recall depends on the specific goals and requirements of the application.

Q7. Discuss the importance of choosing an appropriate evaluation metric for a classification problem and

explain how this can be done.

Ans:Choosing an appropriate evaluation metric for a classification problem is crucial because it determines how the performance of the model is assessed and compared to the desired goals. Different metrics emphasize different aspects of model performance, and the choice depends on the specific characteristics and requirements of the problem at hand. Here are some key considerations and steps for choosing an appropriate evaluation metric:

Importance of Choosing the Right Metric:

Alignment with Business Goals:

 The choice of metric should align with the ultimate goals of the business or application. For example, in a medical diagnosis task, the cost of false positives and false negatives may differ, and the metric should reflect these considerations.

Trade-offs between Precision and Recall:

Precision and recall are often inversely related. Choosing a metric that balances
these trade-offs is important, especially when there are costs associated with
false positives and false negatives.

Understanding Class Imbalance:

• In imbalanced datasets, where one class is significantly more prevalent than the other, accuracy alone may be misleading. Metrics like precision, recall, F1 score, or area under the ROC curve (AUC-ROC) may be more informative.

Application-Specific Considerations:

• Different applications may prioritize different aspects of model performance. For instance, in fraud detection, recall (identifying as many fraudulent cases as possible) might be more critical than precision.

Interpretability:

 Some metrics, like accuracy and precision, are easy to interpret and communicate to stakeholders. However, more complex metrics may be necessary for a nuanced evaluation, especially in cases with multiple classes or imbalanced datasets.

Steps for Choosing an Appropriate Metric:

Understand the Problem Domain:

 Gain a deep understanding of the specific characteristics of the problem, including the nature of the data, class distribution, and any domain-specific considerations.

Define Success Criteria:

Clearly define what success looks like in the context of the application. Identify
the outcomes that have the most significant impact and align with the business
goals.

Consider Class Imbalance:

 If the dataset is imbalanced, consider metrics that account for this, such as precision, recall, F1 score, or area under the ROC curve.

Evaluate Trade-offs:

 Evaluate the trade-offs between precision and recall. Some metrics, like the F1 score, provide a balanced view of both. Consider whether false positives or false negatives have a more significant impact on the problem.

Review Domain-Specific Requirements:

 Some applications may have specific requirements or regulations that dictate the importance of certain types of errors. Ensure that the chosen metric reflects these considerations.

Explore Multiple Metrics:

 Use multiple metrics to get a holistic view of model performance. For instance, ROC curves and AUC-ROC provide insights into the model's ability to discriminate between classes.

Iterative Process:

• Choosing an evaluation metric is often an iterative process. As the project progresses and more insights are gained, revisit the choice of metric to ensure it aligns with the evolving understanding of the problem.

Examples of Evaluation Metrics:

Accuracy:

 Appropriate for balanced datasets but may be misleading in imbalanced scenarios.

Precision:

 Suitable when false positives are costly or when there's a need for a high positive predictive value.

Recall:

• Suitable when false negatives are costly or when there's a need for a high true positive rate.

F1 Score:

• Balances precision and recall, providing a single metric that considers both false positives and false negatives.

Area under the ROC Curve (AUC-ROC):

• Measures the model's ability to discriminate between positive and negative instances across different probability thresholds.

Area under the Precision-Recall Curve (AUC-PR):

 Focuses on the trade-off between precision and recall, often used in imbalanced datasets.

Ultimately, the choice of an evaluation metric should be guided by a comprehensive understanding of the problem and its specific requirements. Regular communication with stakeholders and domain experts is essential to ensure that the chosen metric aligns with the goals of the project.

Q8. Provide an example of a classification problem where precision is the most important metric, and explain why.

Ans:Let's consider a medical diagnosis scenario where the classification problem involves identifying whether a patient has a rare and potentially life-threatening disease (positive class) or not (negative class). In this context, precision becomes a crucial metric. Here's why:

Example: Medical Diagnosis for a Rare Disease

Class Definitions:

- Positive Class (+): Patients with the rare disease (e.g., only 1% of the population is affected).
- Negative Class (-): Patients without the disease.

Importance of Precision:

Consequences of False Positives:

False positives in this context mean classifying a patient as having the disease
when they do not actually have it. This could lead to unnecessary stress, anxiety,
additional medical tests, and potentially harmful treatments.

Low Prevalence of the Disease:

 Since the disease is rare (1% prevalence), a classifier that predicts everyone as negative would achieve high accuracy but would be useless for detecting the actual cases.

Emphasis on Avoiding False Alarms:

 The medical community and patients would likely prioritize avoiding false alarms (false positives). Identifying a few true positive cases is important, but it's more critical to minimize the number of healthy individuals incorrectly classified as positive.

Precision Definition:

• Precision is the ratio of true positives to the total number of instances predicted as positive. In this case, it's the proportion of correctly identified patients with the disease among all instances predicted as positive.

Precision=True PositivesTrue Positives+False Positives
Precision=

True Positives+False Positives
True Positives

Scenario:

- Suppose a model achieves high precision (e.g., 0.90).
- Out of 100 instances predicted as positive, 90 are true positives, and 10 are false positives.
- Precision = 90 / (90 + 10) = 0.90 (90%).

Conclusion:

In this medical diagnosis example, precision is the most important metric because it reflects the accuracy of positive predictions among all instances predicted as positive. High precision ensures that when the model predicts a patient has the rare disease, there is a high likelihood that the prediction is correct, minimizing the number of false positives and the associated negative consequences for patients.

Q9. Provide an example of a classification problem where recall is the most important metric, and explain

why.

Ans:Let's consider a scenario in the context of spam email detection, where the classification problem involves identifying whether an email is spam or not. In this case, recall becomes a crucial metric. Here's why:

Example: Spam Email Detection

Class Definitions:

- Positive Class (+): Spam emails.
- Negative Class (-): Non-spam (legitimate) emails.

Importance of Recall:

Consequences of False Negatives:

 False negatives in this context mean failing to classify a spam email correctly, allowing it to reach the user's inbox. The consequences of missing a potentially harmful spam message, such as phishing attempts or malware, can be severe.

High Volume of Non-Spam Emails:

Email inboxes often receive a large volume of non-spam emails. Even if the spam
prevalence is relatively low, missing a small percentage of spam emails could still
result in a significant number of false negatives.

User Experience and Trust:

 Users generally expect their spam filters to be effective in catching spam. If the system allows too many spam emails to slip through (false negatives), users may lose trust in the filtering system, and their experience with the email service could be compromised.

Emphasis on Identifying All Spam:

 In spam detection, there is often a preference for being overly cautious and capturing as much spam as possible, even if it means some false positives. It's better to have a few legitimate emails marked as spam (false positives) than to let actual spam reach the inbox (false negatives).

Recall Definition:

 Recall is the ratio of true positives to the total number of actual positive instances. In this case, it's the proportion of correctly identified spam emails among all actual spam instances.

Recall=True PositivesTrue Positives+False Negatives
Recall=

True Positives+False Negatives
True Positives

Scenario:

- Suppose a spam filter achieves high recall (e.g., 0.95).
- Out of 100 actual spam emails, the system correctly identifies 95, and 5 are missed (false negatives).
- Recall = 95 / (95 + 5) = 0.95 (95%).

Conclusion:

In the context of spam email detection, recall is the most important metric because it reflects the ability of the system to capture a high percentage of actual spam emails. High recall ensures that the spam filter is effective in identifying and blocking the majority of spam, even if

it means some false positives. This emphasis on capturing all spam emails contributes to a better user experience and maintains the trustworthiness of the email filtering system.