# Model Evaluation Metrics

Understand how to assess the performance of machine learning models using appropriate evaluation metrics for classification, regression, and ranking tasks.
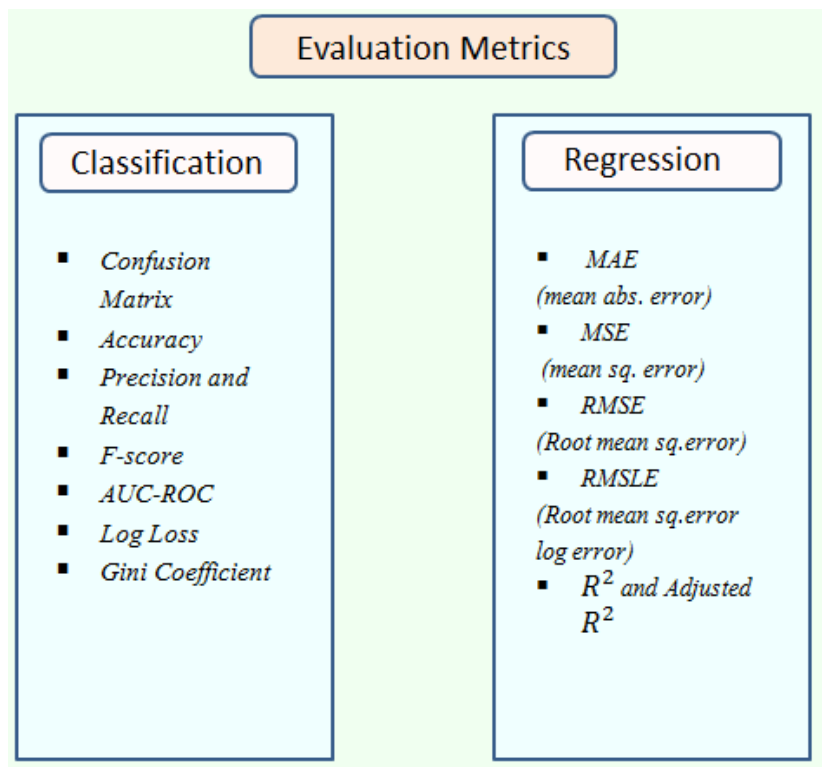
Model evaluation metrics are quantitative measures used to assess the performance of a machine learning model. These metrics provide insights into how well a model generalizes to unseen data and help in comparing different models or algorithms. The choice of metric depends on the specific machine learning task (e.g., classification, regression)
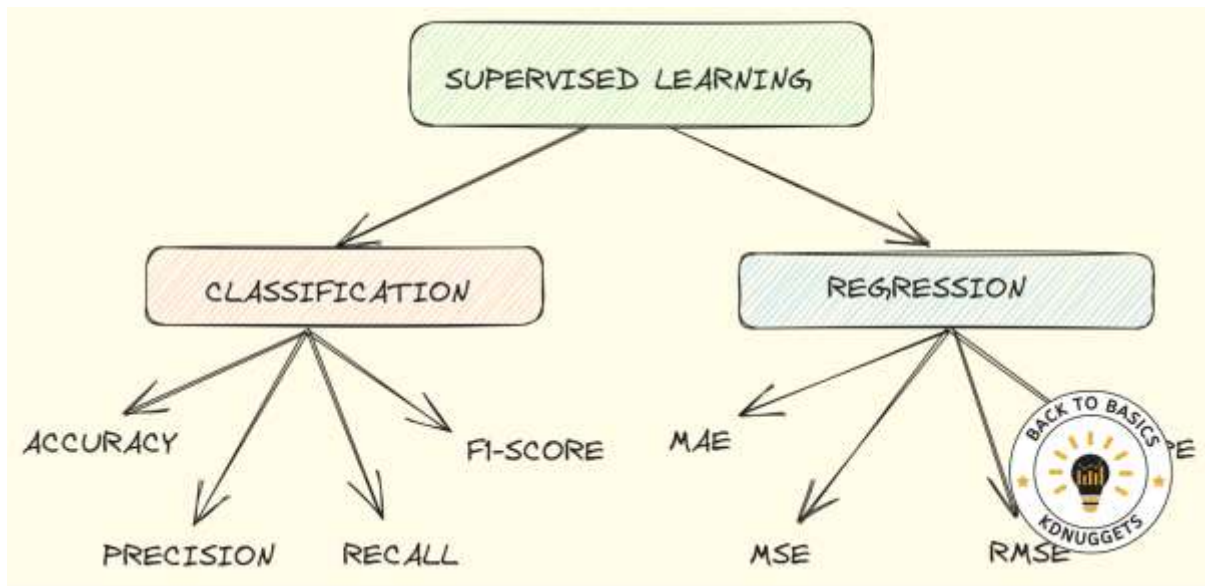
## Why Model Evaluation Is Important

- Helps determine **how well a model generalizes** to unseen data.
- Guides **model selection and tuning**.
- Prevents **overfitting and underfitting**.
- Ensures **reliable performance** in production environments

**Purpose of Model Evaluation:**

- To **assess how well a machine learning model performs** on unseen data.
- To **select the best model** among alternatives.
- To **tune hyperparameters** effectively.
- To **avoid overfitting or underfitting**.



.

## A. For Classification Models:

- **Confusion Matrix:**

  A table summarizing the performance of a classification model, showing true positives, true negatives, false positives, and false negatives.

- **Accuracy:** The proportion of correctly predicted instances out of the total instances.
  Code
  ```
  Accuracy = (True Positives + True Negatives) / Total Predictions
  ```

- **Precision:** The proportion of true positive predictions among all positive predictions. It focuses on minimizing false positives.
  Code
  ```
  Precision = True Positives / (True Positives + False Positives)
  ```

- **Recall (Sensitivity):** The proportion of true positive predictions among all actual positive instances. It focuses on minimizing false negatives.
  Code
  ```
  Recall = True Positives / (True Positives + False Negatives)
  ```

- **F1-Score:** The harmonic mean of precision and recall, providing a balanced measure, especially useful for imbalanced datasets.
  Code
  ```
  F1-Score = 2 * (Precision * Recall) / (Precision + Recall)
  ```

- **ROC Curve and AUC:**

  The Receiver Operating Characteristic (ROC) curve plots the True Positive Rate against the False Positive Rate at various threshold settings. The Area Under the Curve (AUC) provides

a single value summarizing the overall performance of the classifier across all possible thresholds.

- **Log Loss (Cross-Entropy Loss):**
  Measures the performance of a classification model where the prediction input is a probability value between 0 and 1. It penalizes incorrect classifications more heavily when the predicted probability is confident.

## B. For Regression Models:

- **Mean Absolute Error (MAE):**

  The average of the absolute differences between predicted and actual values. It provides a straightforward measure of prediction error.

- **Mean Squared Error (MSE):**

  The average of the squared differences between predicted and actual values. It penalizes larger errors more heavily than MAE.

- **Root Mean Squared Error (RMSE):**

  The square root of MSE, bringing the error back to the original scale of the target variable, making it more interpretable.

- **R-squared (Coefficient of Determination):**

  Represents the proportion of the variance in the dependent variable that is predictable from the independent variables. A higher R-squared indicates a better fit.

- **Adjusted R-squared:**
  A modified version of R-squared that adjusts for the number of predictors in the model, providing a more accurate measure when comparing models with different numbers of features.

## Tips for Metric Selection: Choose metrics based on **task type**, **data distribution**, and **problem context**.

| Scenario | Suggested Metrics |
| --- | --- |
| Binary Classification | F1 Score, ROC-AUC, Precision-Recall |
| Imbalanced Classes | F1 Score, PR Curve |
| Regression | MAE, RMSE, R² |
| Multi-class Classification | Macro F1, Confusion Matrix |
| Recommendation Systems | Precision@k, NDCG, MAP |
| Medical Diagnosis | Recall, Specificity, AUC |

# Classification Metrics

## 1. Confusion Matrix

In the field of <u>machine learning</u> and specifically the problem of <u>statistical classification</u>, a **confusion matrix**, also known as **error matrix**,[1] is a specific <u>table</u> layout that allows visualization of the performance of an algorithm, typically a <u>supervised learning</u> one; in <u>unsupervised learning</u> it is usually called a **matching matrix**.

Each row of the <u>matrix</u> represents the instances in an actual class while each column represents the instances in a predicted class, or vice versa – both variants are found in the literature.[2] The diagonal of the matrix therefore represents all instances that are correctly predicted.[3] The name stems from the fact that it makes it easy to see whether the system is confusing two classes (i.e. commonly mislabeling one as another).

It is a special kind of <u>contingency table</u>, with two dimensions ("actual" and "predicted"), and identical sets of "classes" in both dimensions (each combination of dimension and class is a variable in the contingency table).

## Example

Given a sample of 12 individuals, 8 that have been diagnosed with cancer and 4 that are cancer-free, where individuals with cancer belong to class 1 (positive) and non-cancer individuals belong to class 0 (negative), we can display that data as follows:

| Individual number | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Actual classification** | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |

Assume that we have a classifier that distinguishes between individuals with and without cancer in some way, we can take the 12 individuals and run them through the classifier. The classifier then makes 9 accurate predictions and misses 3: 2 individuals with cancer wrongly predicted as being cancer-free (sample 1 and 2), and 1 person without cancer that is wrongly predicted to have cancer (sample 9).

| Individual number | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Actual classification** | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| **Predicted classification** | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |

Notice, that if we compare the actual classification set to the predicted classification set, there are 4 different outcomes that could result in any particular column. One, if the actual classification is positive and the predicted classification is positive (1,1), this is called a true positive result because the positive sample was correctly identified by the classifier. Two, if the actual classification is positive and the predicted classification is negative (1,0), this is called a false negative result because the positive sample is incorrectly identified by the classifier as being negative. Third, if the actual classification is negative and the predicted

classification is positive (0,1), this is called a false positive result because the negative sample is incorrectly identified by the classifier as being positive. Fourth, if the actual classification is negative and the predicted classification is negative (0,0), this is called a true negative result because the negative sample gets correctly identified by the classifier.

We can then perform the comparison between actual and predicted classifications and add this information to the table, making correct results appear in green so they are more easily identifiable.

| Individual number | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Actual classification | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| Predicted classification | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| Result | FN | FN | TP | TP | TP | TP | TP | TP | FP | TN | TN | TN |

The template for any binary confusion matrix uses the four kinds of results discussed above (true positives, false negatives, false positives, and true negatives) along with the positive and negative classifications. The four outcomes can be formulated in a 2×2 *confusion matrix*, as follows:

| | Total population = P + N | Predicted condition | |
|---|---|---|---|
| | | **Positive (PP)** | **Negative (PN)** |
| **Actual condition** | **Positive (P)** | **True positive (TP)** | **False negative (FN)** |
| | **Negative (N)** | **False positive (FP)** | **True negative (TN)** |

Sources: [4][5][6][7][8][9][10]

The color convention of the three data tables above were picked to match this confusion matrix, in order to easily differentiate the data.

Now, we can simply total up each type of result, substitute into the template, and create a confusion matrix that will concisely summarize the results of testing the classifier:

| | Total 8 + 4 = 12 | Predicted condition | |
|---|---|---|---|
| | | **Cancer** 7 | **Non-cancer** 5 |
| **Actual condition** | **Cancer** 8 | 6 | 2 |
| | **Non-cancer** 4 | 1 | 3 |

In this confusion matrix, of the 8 samples with cancer, the system judged that 2 were cancer-free, and of the 4 samples without cancer, it predicted that 1 did have cancer. All correct predictions are located in the diagonal of the table (highlighted in green), so it is easy to visually inspect the table for prediction errors, as values outside the diagonal will represent them. By summing up the 2 rows of the confusion matrix, one can also deduce the total number of positive (P) and negative (N) samples in the oiginal dataset,  i.e.
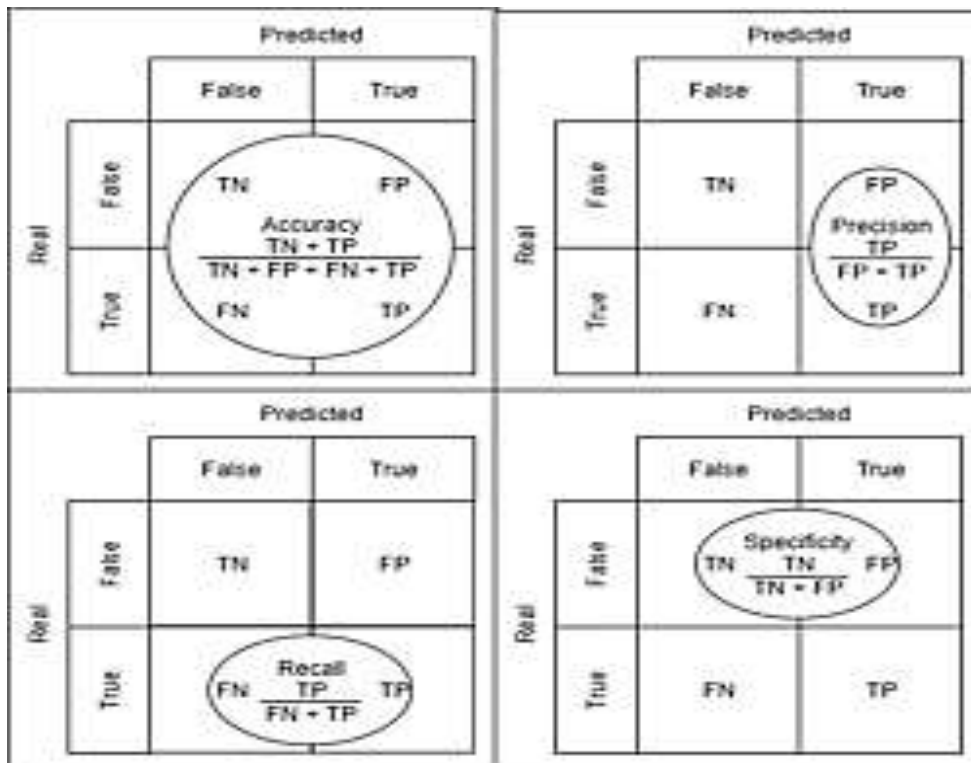
$$P = TP + FN \text{ and } N = FP + TN.$$

.

# Table of confusion

In predictive analytics, a **table of confusion** (sometimes also called a **confusion matrix**) is a table with two rows and two columns that reports the number of *true positives*, *false negatives*, *false positives*, and *true negatives*. This allows more detailed analysis than simply observing the proportion of correct classifications (accuracy). Accuracy will yield misleading results if the data set is unbalanced; that is, when the numbers of observations in different classes vary greatly.

For example, if there were 95 cancer samples and only 5 non-cancer samples in the data, a particular classifier might classify all the observations as having cancer. The overall accuracy would be 95%, but in more detail the classifier would have a 100% recognition rate (sensitivity) for the cancer class but a 0% recognition rate for the non-cancer class. F1 score is even more unreliable in such cases, and here would yield over 97.4%, whereas informedness removes such bias and yields 0 as the probability of an informed decision for any form of guessing (here always guessing cancer).

Other metrics can be included in a confusion matrix, each of them having their significance and use.

Some researchers have argued that the confusion matrix, and the metrics derived from it, do not truly reflect a model's *knowledge*. In particular, the confusion matrix cannot show whether correct predictions were reached through sound reasoning or merely by chance (a problem known in philosophy as epistemic luck). It also does not capture situations where the facts used to make a prediction later change or turn out to be wrong (defeasibility). This means that while the confusion matrix is a useful tool for measuring classification performance, it may give an incomplete picture of a model's true reliability.[20]

Confusion Matrix Metrics. Including accuracy, precision, sensitivity (recall) and specificity

# Confusion matrices with more than two categories

Confusion matrix is not limited to binary classification and can be used in multi-class classifiers as well. The confusion matrices discussed above have only two conditions: positive and negative. For example, the table below summarizes communication of a whistled language between two speakers, with zero values omitted for clarity.[21]

| Perceived vowel Vowel produced | i | e | a | o | u |
|---|---|---|---|---|---|
| i | 15 | | 1 | | |
| e | 1 | | 1 | | |
| a | | | 79 | 5 | |
| o | | | 4 | 15 | 3 |
| u | | | | 2 | 2 |

Multi category confusion matrix. Displaying the positions of false positives, false negatives, true positives and true negatives

## Confusion matrices in multi-label and soft-label classification

Confusion matrices are not limited to underline{single-label classification} (where only one class is present) or hard-label settings (where classes are either fully present, 1, or absent, 0). They can also be extended to underline{Multi-label classification} (where multiple classes can be predicted at once) and soft-label classification (where classes can be partially present).

One such extension is the **Transport-based Confusion Matrix (TCM)**,[22] which builds on the theory of underline{optimal transport} and the underline{principle of maximum entropy}. TCM applies to single-label, multi-label, and soft-label settings. It retains the familiar structure of the standard confusion matrix: a square matrix sized by the number of classes, with diagonal entries indicating correct predictions and off-diagonal entries indicating confusion. In the single-label case, TCM is identical to the standard confusion matrix.

TCM follows the same reasoning as the standard confusion matrix: if class A is overestimated (its predicted value is greater than its label value) and class B is underestimated (its predicted value is less than its label value), A is considered confused with B, and the entry (B, A) is increased. If a class is both predicted and present, it is correctly identified, and the diagonal entry (A, A) increases. Optimal transport and maximum entropy are used to determine the extent to which these entries are updated.

TCM enables clearer comparison between predictions and labels in complex classification tasks, while maintaining a consistent matrix format across settings

## 2. Accuracy

**Accuracy** is the ratio of correctly predicted observations to the total observations.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} = \frac{TP + TN}{TP + TN + FP + FN}$$

Where:

- **TP** = True Positives
- **TN** = True Negatives
- **FP** = False Positives
- **FN** = False Negatives

- **Good for** balanced datasets.
- **Misleading** with imbalanced data.

Example:

**Confusion Matrix:**

|  | Predicted: Yes | Predicted: No |
| --- | --- | --- |
| Actual: Yes | 50 (TP) | 10 (FN) |
| Actual: No | 5 (FP) | 35 (TN) |

- Total predictions = 50 + 10 + 5 + 35 = 100
- Correct predictions = 50 (TP) + 35 (TN) = 85

$$\text{Accuracy} = \frac{85}{100} = 0.85 \text{ or } 85\%$$

# Advantages of Accuracy

1. **Easy to understand and compute**
2. Gives a **quick overview** of model performance.
3. Works well when the **classes are balanced** (roughly equal number of positives and negatives).

# Limitations of Accuracy

1. **Misleading with imbalanced data**
   o Example: A model that always predicts the majority class can have high accuracy but be useless.
   o E.g., in fraud detection where only 1% of transactions are fraud:
      ▪ A model that always predicts "not fraud" will have **99% accuracy**, but will never catch any fraud.
2. **Does not distinguish** between types of errors (FP vs FN).
3. Fails to reflect the **cost of errors**, which might be critical in domains like healthcare or finance.

### 📇 Accuracy in Multi-class Classification

$$\text{Accuracy} = \frac{\text{Number of Correctly Classified Instances}}{\text{Total Instances}}$$

Still measures the proportion of correctly predicted labels but across **multiple classes**.

# When to Use Accuracy

☐ Use Accuracy when:

- Classes are **roughly balanced**.
- The **cost of false positives and false negatives is similar**.

☐ Avoid using Accuracy when:

- You have **imbalanced datasets**.
- You care more about specific types of errors (e.g., missing cancer diagnoses).

| Feature | Description |
|---|---|
| Metric Type | Classification |
| Formula | (TP + TN) / (TP + TN + FP + FN) |
| Best Use Case | Balanced datasets |
| Not Ideal When | Classes are imbalanced |
| Related Metrics | Precision, Recall, F1 Score |

## 3. Precision (Positive Predictive Value)

- Measures correctness among predicted positives.
- **High Precision** = Low false positive rate.

Precision (PPV) measures **how many of the items the model predicted as positive are actually positive**. In other words: when the model says "positive", how likely is it to be correct?

Mathematically:

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

Here:

- TP = number of instances correctly predicted as positive
- FP = number of instances wrongly predicted as positive (i.e., predicted positive but actually negative)

  The Data Comp... +2

- So the denominator (TP + FP) is *all instances predicted as positive*.

## Example Scenario

Suppose we're building a binary classification model to detect **fraudulent transactions**.

- "Positive" = transaction is fraudulent
- "Negative" = transaction is non-fraudulent

After testing the model, we obtain the following counts:

|  | Predicted Fraudulent (Positive) | Predicted Non-Fraudulent (Negative) |
|---|---|---|
| **Actual Fraudulent** | True Positives (TP) = 50 | False Negatives (FN) = 10 |
| **Actual Non-Fraudulent** | False Positives (FP) = 40 | True Negatives (TN) = 900 |

So in summary:

- TP = 50
- FP = 40
- FN = 10
- TN = 900

### Calculating Precision

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{50}{50 + 40} = \frac{50}{90} \approx 0.5556 \approx 55.56\%$$

**Interpretation**: Of all the transactions the model predicted as "fraudulent" (50 + 40 = 90 transactions), only 50 were actually fraudulent. So when the model raises a fraud alert, it's correct about **55.56% of the time**.

## Why Precision Matters in This Scenario

- If a transaction is flagged as fraudulent (positive), you might block it or require extra verification — this has cost (inconvenience, customer experience, manual reviews).
- A **low precision** means many false alarms (FPs) — non-fraudulent transactions flagged as fraudulent. That leads to unnecessary cost.
- Therefore, a higher precision is **important** because you want predictions of "fraud" to be *very reliable*.

### Things to note

- Precision does **not** account for the transactions you missed (false negatives, FN = 10). That's the focus of **Recall**.
- Even though Precision ~ 55.6% here, you could still have high recall (catching many frauds) or low recall depending on FN and TP.
- If you tried to **increase** precision (make the model more conservative in flagging frauds), you might reduce FP (good) but you might also **increase** FN (miss more frauds) → trade-off between precision and recall.

## Variation / Additional Example

Another simpler example:

Imagine a spam filter that predicted 200 emails as spam. Out of those 200 flagged, 150 were actually spam (TP = 150) and 50 were legitimate emails wrongly flagged as spam (FP = 50). Then:

$$\text{Precision} = \frac{150}{150 + 50} = \frac{150}{200} = 0.75 = 75\%$$

So when this filter flags something as spam, it is correct 75% of the time.

# Additional Considerations

- **Threshold tuning**: For models that output probability scores (e.g., logistic regression), you can adjust the threshold above which you predict "positive". Raising the threshold typically increases precision (fewer predictions as positive → fewer FP) but may reduce recall. Lowering threshold does opposite.
- **Precision-Recall Curve**: Plotting precision vs recall for different thresholds gives insight into model trade-offs. Especially useful for imbalanced datasets. Medium
- **Dependence on prevalence**: In diagnostic testing, the PPV (precision) depends not only on the test's sensitivity / specificity but also on the prevalence of the condition in the population. Wikipedia
- **Interpretation in multi-class/multi-label**: For multi-class classification, you can compute precision per class (treating each class as "positive" in turn) and then average (macro/micro/weighted) depending on your focus.

| Metric | Definition | Focus | Use Case |
|---|---|---|---|
| Precision | TP / (TP + FP) | Accuracy of positive predictions | Minimizing false positives, ensuring positive predictions are reliable |
| Recall | TP / (TP + FN) | Capturing all positives | Minimizing false negatives, catching as many positives as possible |
| F1 Score | Harmonic mean of precision & recall | Balance between precision & recall | When both FP & FN matter and you need a single summary metric |

## 4. . Recall (Sensitivity, True Positive Rate)

- Measures ability to find all positive samples.
- **High Recall** = Low false negative rate.

Recall (Sensitivity / True Positive Rate) measures **how many of the actual positive class instances the model correctly identifies as positive**.
In other words, when something *is* positive, how likely is the model to pick it up?

Mathematically:

$$Recall = \frac{True\ Positives\ (TP)}{True\ Positives\ (TP) + False\ Negatives\ (FN)}$$

Where:

- TP = number of instances correctly predicted as positive
- FN = number of instances that are positive but the model predicted negative
- So the denominator is *all actual positives* = TP + FN.  gao-hongnan.gi... +1
- From a probability viewpoint:

$$Recall = P(\text{model predicts positive} \mid \text{actual is positive})$$

Thus recall is a measure of **completeness** of the positive predictions: of all the positives there are in truth, how many did you manage to catch?

### 🔴 Intuition & When It Matters

- Recall asks the question: *"Of all the positives, how many did I predict as positive?"*
- It doesn't care about how many were predicted positive wrongly (that is what precision handles). It focuses on missing positives.
- Use recall when **missing a positive is costly** or undesirable.

### Example contexts where high recall is important:

- A medical test for a disease: if a person actually has the disease (positive), you want the model/test to *catch* it. A false negative (the model says "no disease" when actually disease) is high cost.
- Fraud detection: if a transaction *is* fraudulent (positive), you want to catch it rather than letting it slip through (false negative).
- Safety / security applications: if an item/event *is* dangerous (positive), you want detection rather than missing it.

In all these, the priority is ensuring actual positives aren't "missed".

**Example**

Imagine a binary classification model to detect a rare disease (positive class = "disease present"). Suppose:

- TP = 90 (90 patients are diseased and correctly predicted positive)
- FN = 10 (10 patients are diseased but predicted negative)
- And – say – many who don't have disease (but we focus here on positives)

Then:

$$\text{Recall} = \frac{90}{90 + 10} = \frac{90}{100} = 0.90 \quad \text{(or 90\%)}$$

Meaning: the model catches 90% of all diseased patients.

**Example Scenario**

Let's imagine a model that diagnoses whether people have a certain disease ("Positive") or not ("Negative").
After testing on 200 patients, the confusion matrix comes out as:

|  | **Predicted: Disease (Positive)** | **Predicted: No Disease (Negative)** |
|---|---|---|
| **Actual: Disease** | True Positive (TP) = 45 | False Negative (FN) = 15 |
| **Actual: No Disease** | False Positive (FP) = 20 | True Negative (TN) = 120 |

So:

- TP = 45
- FN = 15
- FP = 20
- TN = 120
- Total actual positives = TP + FN = 45 + 15 = 60
- Total actual negatives = FP + TN = 20 + 120 = 140
- Total patients = 200

**Calculating Recall**

Recall = TP / (TP + FN)
= 45 / (45 + 15)
= 45 / 60
= 0.75 or **75%**

**Interpretation**: Out of all the patients who actually have the disease (60 people), the model correctly identified 45 of them. So the model catches 75% of the actual positive cases.

**Why This Matters**

- A recall of 75% means **25% of the diseased patients** were missed (FN = 15).
- If missing a diseased patient is very serious (e.g., a disease with major consequences), then a 75% recall might be **insufficient**.
- If the priority is to catch as many diseased cases as possible (minimize FN), you might want to **improve recall**, even if it means accepting more false positives (FP) and a lower precision.

**Summary**

- Recall tells you: **Of all the actual positives**, how many did we correctly identify?
- In this example: 75%.
- Good recall means few false negatives; high recall is **essential** in domains where missing a positive is very costly.
- But recall alone doesn't tell you about false positives—so you also need to look at precision (and other metrics) to get the full picture.

If recall were, say, 0.50, then half of the diseased people would go undetected.

**Additional Views**

- In this case, precision could also be computed:
  Precision = TP / (TP + FP) = 45 / (45 + 20) = 45 / 65 ≈ 0.692 or ~69.2%
  So when the model predicts "disease", it's correct about ~69.2% of the time.
- Balanced view: The model has a trade-off: recall = 75% (good, but could be better) while precision ~ 69% (so some false alarms).
- Depending on the domain, you might adjust the decision threshold to boost recall (label more cases as positive) but that will likely reduce precision (more FP).

**Key Properties**

- Range: 0 to 1 (or 0% to 100%). The maximum recall is 1 when FN = 0 (i.e., no actual positive is missed). [Google for Developers](#)
- Does *not* account for false positives (FP). So recall might be high even if the model wrongly flags many negatives as positive.
- Especially important when the **cost of false negatives** (missed positives) is high.
- Often there's a **trade-off** between recall and precision: raising recall (catching more positives) may increase false positives, reducing precision. [Medium+1](#)
- Recall is independent of negatives count in the denominator — so if the negative class is huge, recall isn't directly affected by those.

---

**When to Use Recall**

- When the main concern is **not missing positives** (i.e., false negatives must be minimized).
- When you have a domain where a missed positive is more harmful than a false alarm (false positive).
- In imbalanced datasets where the positive class is rare and you care about catching as many of them as possible.
- As part of a balanced metric set along with precision (and F1).

### When Recall Alone is Not Enough

- If you also care about false positives (i.e., you need positive predictions to be reliable), then precision must also be considered.
- If the cost of false positives is high and you only maximize recall, you might generate many false positives.
- Using recall only may ignore the model's performance on negatives (or how many false positives it generates).

### Relationship with Other Metrics

- Recall is the complement of the **False Negative Rate (FNR)**:

$$FNR = \frac{FN}{TP + FN} = 1 - Recall$$

aryaxai.com +1

- In binary classification, recall = sensitivity = true positive rate (TPR).

$$TPR = \frac{TP}{TP + FN}$$

Wikipedia +1

- Paired with precision:

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

(so precision focuses on *predicted positives*, recall on *actual positives*)   Wikipedia +1

One common summary metric is the **F1 Score**, the harmonic mean of precision and recall, used when both are important.

## Additional Considerations

- **Threshold tuning**: For probabilistic classifiers, you can adjust the decision threshold to increase recall (e.g., label more cases as positive). This often increases FP, so precision may drop.
- **Imbalanced data**: On rare positive classes, recall becomes critical because you might otherwise miss most positives.
- **Domain cost trade-offs**: In e.g. cancer screening, false negatives (missing a sick person) are often far worse than false positives (unnecessary further tests). So recall tends to be prioritized.
- **Multi-class / multi-label**: For more than two classes, recall can be computed per-class (i.e., for each class treat it as positive vs rest) and then averaged (macro/micro).
- **Interpretation caution**: A model can have *perfect recall (1.0)* by predicting *all* instances positive. That catches all actual positives (FN=0) but likely produces many false positives. So recall alone doesn't tell you everything.

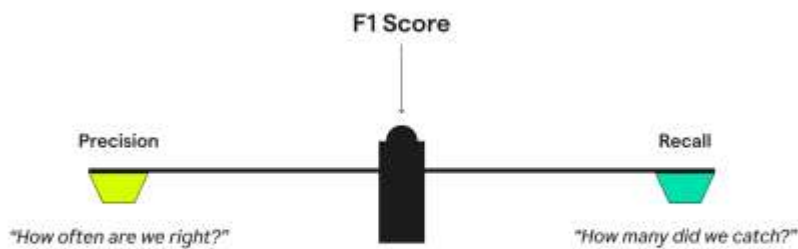| Metric | Definition | Focus | Use-case condition |
|---|---|---|---|
| Recall (Sensitivity, TPR) | $TP / (TP + FN)$ | How many actual positives are caught | When missing positives is costly |
| Precision | $TP / (TP + FP)$ | How correct predicted positives are | When false positives are costly |
| F1 Score | Harmonic mean of precision & recall | Balance between precision & recall | When both FP & FN matter |

# 5. F1 Score

**Best for** imbalanced datasets.

ML practitioners face a common challenge when building classification models: training the model to catch all cases while avoiding false alarms. This is particularly important in critical applications like financial fraud detection and medical diagnosis, where false alarms and missing important classifications have serious consequences. Achieving the right balance is particularly important when dealing with imbalanced datasets, where a category like fraudulent transactions is much rarer than the other category (legitimate transactions).

**Precision and recall**

To measure model performance quality, the F1 score combines two related metrics:

- **Precision**, which answers, "When the model predicts a positive case, how often is it correct?"
- **Recall**, which answers, "Of all actual positive cases, how many did the model correctly identify?"

A model with high precision but low recall is overly cautious, missing many true positives, while one with high recall but low precision is overly aggressive, generating many false positives. The F1 score strikes a balance by taking the harmonic mean of precision and recall, which gives more weight to lower values and ensures that a model performs well on both metrics rather than excelling in just one.

F1 Score

Precision | Recall

"How often are we right?" | "How many did we catch?"

## Precision and recall example

To better understand precision and recall, consider a spam detection system. If the system has a high rate of correctly flagging emails as spam, this means it has high precision. For example, if the system flags 100 emails as spam, and 90 of them are actually spam, the precision is 90%. High recall, on the other hand, means the system catches most actual spam emails. For example, if there are 200 actual spam emails and our system catches 90 of them, the recall is 45%.

## Variants of the F1 score

In multiclass classification systems or scenarios with specific needs, the F1 score can be calculated in different ways, depending on what factors are important:

- **Macro-F1:** Calculates the F1 score separately for each class and takes the average
- **Micro-F1:** Calculates recall and precision over all predictions
- **Weighted-F1**: Similar to Macro-F1, but classes are weighted based on frequency

## Beyond the F1 score: The F-score family

The F1 score is part of a larger family of metrics called the F-scores. These scores offer different ways to weight precision and recall:

- **F2:** Places greater emphasis on recall, which is useful when false negatives are costly
- **F0.5:** Places greater emphasis on precision, which is useful when false positives are costly

## How to calculate an F1 score

The F1 score is mathematically defined as the harmonic mean of precision and recall. While this might sound complex, the calculation process is straightforward when broken down into clear steps.

**The formula for the F1 score:**

$$F_1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

Before diving into the steps to calculate F1, it's important to understand the key components of what's called a *confusion matrix*, which is used to organize classification results:

- **True positives (TP):** The number of cases correctly identified as positive
- **False positives (FP):** The number of cases incorrectly identified as positive
- **False negatives (FN):** The number of cases missed (actual positives that were not identified)

The general process involves training the model, testing predictions and organizing results, calculating precision and recall, and calculating the F1 score.

**Step 1: Train a classification model**

First, a model must be trained to make binary or multiclass classifications. This means that the model needs to be able to classify cases as belonging to one of two categories. Examples include "spam/not spam" and "fraud/not fraud."

**Step 2: Test predictions and organize results**

Next, use the model to perform classifications on a separate dataset that wasn't used as part of the training. Organize the results into the confusion matrix. This matrix shows:

- TP: How many predictions were actually correct
- FP: How many positive predictions were incorrect
- FN: How many positive cases were missed

The confusion matrix provides an overview of how the model is performing.

**Step 3: Calculate precision**

Using the confusion matrix, precision is calculated with this formula:

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

For example, if a spam detection model correctly identified 90 spam emails (TP) but incorrectly flagged 10 nonspam emails (FP), the precision is 0.90.

$$Precision = \frac{90}{90 + 10} = 0.90$$

**Step 4: Calculate recall**

Next, calculate recall using the formula:

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

Using the spam detection example, if there were 200 total spam emails, and the model caught 90 of them (TP) while missing 110 (FN), the recall is 0.45.

$$Recall = \frac{90}{90 + 110} = 0.45$$

**Step 5: Calculate the F1 score**

With the precision and recall values in hand, the F1 score can be calculated.

The F1 score ranges from 0 to 1. When interpreting the score, consider these general benchmarks:

- **0.9 or higher:** The model is performing great, but should be checked for overfitting.
- **0.7 to 0.9:** Good performance for most applications
- **0.5 to 0.7:** Performance is OK, but the model could use improvement.
- **0.5 or less:** The model is performing poorly and needs serious improvement.

Using the spam detection example calculations for precision and recall, the F1 score would be 0.60 or 60%.

$$F_1 = 2 \times \frac{0.90 \times 0.45}{0.90 + 0.45}$$

In this case, the F1 score indicates that, even with high precision, the lower recall is affecting overall performance. This suggests that there's room for improvement in catching more spam emails.

**F1 score vs. accuracy**

While both F1 and *accuracy* quantify model performance, the F1 score provides a more nuanced measure. Accuracy simply calculates the percentage of correct predictions. However, just relying on accuracy to measure model performance can be problematic when the number of instances of one category in a dataset significantly outnumbers the other category. This problem is referred to as the *accuracy paradox*.

To understand this problem, consider the example of the spam detection system. Suppose an email system receives 1,000 emails every day, but only 10 of those are actually spam. If spam detection simply classifies every email as not spam, it will still achieve 99% accuracy. This is because 990 predictions out of 1,000 were correct, even though the model is actually useless when it comes to spam detection. Clearly, accuracy does not give an accurate picture of the quality of the model.

The F1 score avoids this problem by combining the precision and recall measurements. Therefore, F1 should be used instead of accuracy in the following cases:

- **The dataset is imbalanced.** This is common in fields like diagnosis of obscure medical conditions or spam detection, where one category is relatively rare.
- **FN and FP are both important.** For example, medical screening tests seek to balance catching actual issues with not raising false alarms.
- **The model needs to strike a balance between being too aggressive and too cautious.** For example, in spam filtering, an overly cautious filter might let through too much spam (low recall) but rarely make mistakes (high precision). On the other hand, an overly aggressive filter might block real emails (low precision) even if it does catch all spam (high recall).

**Applications of the F1 score**

The F1 score has a wide range of applications across various industries where balanced classification is critical. These applications include financial fraud detection, medical diagnosis, and content moderation.

**Financial fraud detection**

Models designed to detect financial fraud are a category of systems well suited for measurement using the F1 score. Financial firms often process millions or billions of transactions daily, with actual cases of fraud being relatively rare. For this reason, a fraud detection system needs to catch as many fraudulent transactions as possible while simultaneously minimizing the number of false alarms and resulting inconvenience to customers. Measuring the F1 score can help financial institutions determine how well their systems balance the twin pillars of fraud prevention and a good customer experience.

**Medical diagnosis**

In medical diagnosis and testing, FN and FP both have serious consequences. Consider the example of a model designed to detect rare forms of cancer. Incorrectly diagnosing a healthy patient could lead to unnecessary stress and treatment, while missing an actual cancer case will have dire consequences for the patient. In other words, the model needs to have both high precision and high recall, which is something that the F1 score can measure.

### Content moderation

Moderating content is a common challenge in online forums, social media platforms, and online marketplaces. To achieve platform safety without overcensoring, these systems must balance precision and recall. The F1 score can help platforms determine how well their system balances these two factors.

### Benefits of the F1 score

In addition to generally providing a more nuanced view of model performance than accuracy, the F1 score provides several key advantages when evaluating classification model performance. These benefits include faster model training and optimization, reduced training costs, and catching [overfitting](overfitting) early.

### Faster model training and optimization

The F1 score can help speed up model training by providing a clear reference metric that can be used to guide optimization. Instead of tuning recall and precision separately, which generally involves complex trade-offs, ML practitioners can focus on increasing the F1 score. With this streamlined approach, optimal model parameters can be identified quickly.

### Reduced training costs

The F1 score can help ML practitioners make informed decisions about when a model is ready for deployment by providing a nuanced, single measure of model performance. With this information, practitioners can avoid unnecessary training cycles, investments in computational resources, and having to acquire or create additional training data. Overall, this can lead to substantial cost reductions when training classification models.

### Catching overfitting early

Since the F1 score considers both precision and recall, it can help ML practitioners identify when a model is becoming too specialized in the training data. This problem, called overfitting, is a common issue with classification models. The F1 score gives practitioners an early warning that they need to adjust training before the model reaches a point where it is unable to generalize on real-world data.

## Limitations of the F1 score

Despite its many benefits, the F1 score has several important limitations that practitioners should consider. These limitations include a lack of sensitivity to true negatives, not being suited for some datasets, and being harder to interpret for multiclass problems.

**Lack of sensitivity to true negatives**

The F1 score doesn't account for true negatives, which means that it is not well suited for applications where measuring this is important. For example, consider a system designed to identify safe driving conditions. In this case, correctly identifying when conditions are genuinely safe (true negatives) is just as important as identifying dangerous conditions. Because it doesn't track FN, the F1 score wouldn't accurately capture this aspect of overall model performance.

**Not suited for some datasets**

The F1 score may not be suited for datasets where the impact of FP and FN are significantly different. Consider the example of a cancer screening model. In such a situation, missing a positive case (FN) could be life-threatening, while wrongly finding a positive case (FP) only leads to additional testing. So, using a metric that can be weighted to account for this cost is a better choice than the F1 score.

**Harder to interpret for multiclass problems**

While variations like micro-F1 and macro-F1 scores mean that the F1 score can be used to evaluate multiclass classification systems, interpreting these aggregated metrics is often more complex than the binary F1 score. For example, the micro-F1 score might hide poor performance in classifying less frequent classes, while the macro-F1 score might overweight rare classes. Given this, businesses need to consider whether equal treatment of classes or overall instance-level performance is more important when choosing the right F1 variant for multiclass classification models.
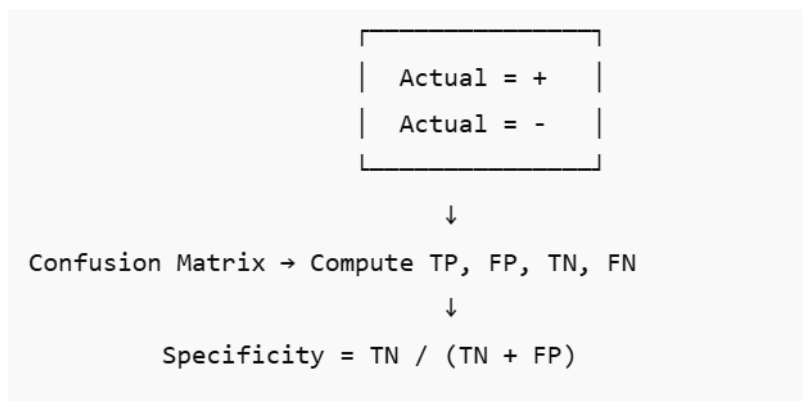
---

# 6. Specificity (True Negative Rate)

In machine learning—especially in **classification problems**—we evaluate model performance using various metrics derived from the **confusion matrix**.

- Specificity is a **core diagnostic metric** in medical, fraud detection, and quality control applications.
- It should **never be evaluated in isolation**—always alongside sensitivity or other balanced metrics (AUC, F1, etc.).
- Understanding the **cost of errors** (FP vs FN) guides whether to prioritize specificity or sensitivity.

The **confusion matrix** summarizes how many instances were correctly or incorrectly classified by the model.

For a **binary classification problem**, the confusion matrix has four outcomes:

| Actual \ Predicted | Positive | Negative |
|---|---|---|
| **Positive** (True class = 1) | True Positive (TP) | False Negative (FN) |
| **Negative** (True class = 0) | False Positive (FP) | True Negative (TN) |

```
        ┌─────────────────┐
        │   Actual = +    │
        │   Actual = -    │
        └─────────────────┘
                 ↓
Confusion Matrix → Compute TP, FP, TN, FN
                 ↓
        Specificity = TN / (TN + FP)
```

**Specificity** measures the **proportion of actual negatives that are correctly identified as negatives**.

Specificity=TN/TN+FP

**Intuition:**

Specificity answers the question:

*"Out of all the truly negative cases, how many did the model correctly classify as negative?"*

**Relation to Other Metrics**

| Metric | Formula | Focuses On |
|---|---|---|
| Sensitivity (Recall / True Positive Rate) | $\frac{TP}{TP+FN}$ | Correctly identifying positives |
| Specificity (True Negative Rate) | $\frac{TN}{TN+FP}$ | Correctly identifying negatives |
| Precision (Positive Predictive Value) | $\frac{TP}{TP+FP}$ | Correctness of positive predictions |
| Accuracy | $\frac{TP+TN}{TP+TN+FP+FN}$ | Overall correctness |

Specificity is the **complement** of the **False Positive Rate (FPR)**:

$$FPR = 1 - \text{Specificity}$$

$$FPR = \frac{FP}{FP + TN}$$

**When is Specificity Important?**

Specificity is crucial when **false positives** are costly or dangerous.

**Example Scenarios:**

| Domain | High Specificity is Needed Because… |
|---|---|
| Medical diagnostics (e.g., HIV test) | We don't want to wrongly tell healthy people they are infected. |
| Spam detection | We don't want to classify legitimate emails as spam. |
| Fraud detection | Wrongly flagging normal transactions as fraudulent harms users. |

**Interpreting Specificity**

- **High Specificity (~1.0):** Model rarely classifies negatives as positives (low false alarm rate).
- **Low Specificity (~0.0):** Model frequently misclassifies negatives as positives (many false alarms).

**Example:**

| TN | FP | Specificity |
|---|---|---|
| 90 | 10 | 0.90 |
| 90 | 30 | 0.75 |
| 90 | 0 | 1.00 |

**Trade-Off with Sensitivity**

There is often a **trade-off** between **Sensitivity** and **Specificity**.

- Increasing sensitivity may lower specificity (and vice versa).
- The trade-off is typically visualized via the **Receiver Operating Characteristic (ROC) Curve**.

**ROC Context:**

- X-axis: **False Positive Rate (1 - Specificity)**
- Y-axis: **True Positive Rate (Sensitivity)**
- **AUC (Area Under Curve)** measures how well the model balances both.

**Balanced Metrics: Combining Sensitivity & Specificity**

To evaluate models considering both aspects:

- Balanced Accuracy:

$$\text{Balanced Accuracy} = \frac{\text{Sensitivity} + \text{Specificity}}{2}$$

- Youden's Index (J-statistic):

$$J = \text{Sensitivity} + \text{Specificity} - 1$$

(A higher J indicates better test performance.)

**Example**

Let's consider a **disease test** on 1000 people:

| Actual / Predicted | Positive | Negative | Total |
|---|---|---|---|
| Positive (diseased) | 80 (TP) | 20 (FN) | 100 |
| Negative (healthy) | 30 (FP) | 870 (TN) | 900 |

**Specificity Calculation:**

$$\text{Specificity} = \frac{TN}{TN + FP} = \frac{870}{870 + 30} = 0.9667$$

✅ Interpretation:

The test correctly identifies 96.67% of healthy individuals as disease-free.

| Concept | Summary |
|---|---|
| Meaning | Proportion of true negatives correctly classified |
| Formula | $\frac{TN}{TN+FP}$ |
| Range | 0 to 1 |
| High Value Means | Few false positives |
| Useful When | False positives are more critical |
| Complement | False Positive Rate (FPR = 1 - Specificity) |
| Trade-off | With Sensitivity (True Positive Rate) |

---

# 7. ROC Curve (Receiver Operating Characteristic)

In classification tasks—especially **binary classification**—we often need to measure **how well a model distinguishes between two classes** (positive vs. negative).

The **ROC Curve (Receiver Operating Characteristic Curve)** is a **graphical tool** used to evaluate a model's **discriminative ability** at various classification thresholds.
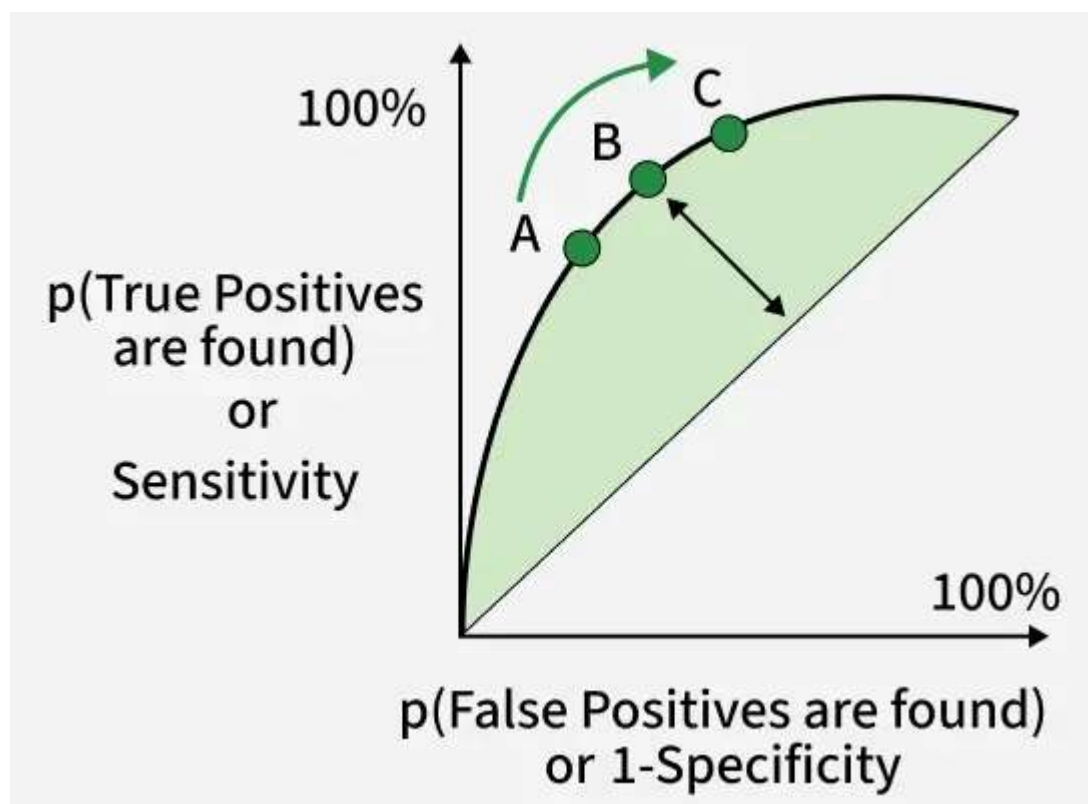
- ROC curve shows **model performance across thresholds**.
- **AUC** summarizes this performance in a single value.
- **Higher AUC ⇒ better ability** to distinguish between classes.
- Use **ROC + AUC** to compare models' ranking quality.
- For **imbalanced datasets**, also consider the **Precision–Recall Curve**.
  - Plots **True Positive Rate (Recall)** vs. **False Positive Rate**.
  - Area Under Curve (**AUC**) shows performance:
    - **AUC = 1.0:** Perfect model
    - **AUC = 0.5:** Random guessing

-----------------

**AUC-ROC curve** is a graph used to check how well a binary classification model works. It helps us to understand how well the model separates the positive cases like people with a disease from the negative cases like people without the disease at different threshold level. It shows how good the model is at telling the difference between the two classes by plotting:

- **True Positive Rate (TPR):** how often the model correctly predicts the positive cases also known as **Sensitivity or Recall**.
- **False Positive Rate (FPR):** how often the model incorrectly predicts a negative case as positive.
- **Specificity:** measures the proportion of actual negatives that the model correctly identifies. It is calculated as 1 - FPR.

The higher the curve the better the model is at making correct predictions.



Sensitivity versus False Positive Rate plot
These terms are derived from the **confusion matrix** which provides the following values:
- **True Positive (TP)**: Correctly predicted positive instances

- **True Negative (TN)**: Correctly predicted negative instances
- **False Positive (FP)**: Incorrectly predicted as positive
- **False Negative (FN)**: Incorrectly predicted as negative

Confusion Matrix for a Classification Task

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

- **ROC Curve** :

  From this matrix, we define two critical rates used in the ROC curve:

$$\text{True Positive Rate (TPR)} = \frac{TP}{TP + FN}$$

$$\text{False Positive Rate (FPR)} = \frac{FP}{FP + TN}$$

- It plots TPR vs. FPR at different thresholds. It represents the trade-off between the sensitivity and specificity of a classifier.

◆ The ROC Curve plots:

$$\text{TPR (Sensitivity)} \quad \text{vs.} \quad \text{FPR (1 - Specificity)}$$

- **X-axis:** False Positive Rate (FPR)
- **Y-axis:** True Positive Rate (TPR)

Each point on the curve corresponds to a **different classification threshold**.

**Steps to Construct an ROC Curve**
1. **Obtain predicted probabilities** for the positive class.
2. **Sort predictions** from highest to lowest.
3. For each unique threshold:
   - Predict **positive** if probability ≥ threshold.
   - Compute **TPR** and **FPR**.
4. **Plot (FPR, TPR)** for each threshold.

# Interpreting the ROC Curve

**Ideal Cases:**

| Type of Model | Description | ROC Curve Shape | AUC Value |
|---|---|---|---|
| Perfect classifier | Always correct | Passes through (0,1) | 1.0 |
| Random classifier | No discriminative power | Diagonal line (0,0) → (1,1) | 0.5 |
| Poor classifier | Worse than random (inverted) | Below diagonal | < 0.5 |

The ROC curve can guide **threshold selection**:

- Choose the point **closest to the top-left corner (0,1)** — minimizes both FPR and FN.
- Can also use **Youden's Index (J)**:

$$J = \text{Sensitivity} + \text{Specificity} - 1$$

Threshold with maximum $J$ gives the best balance.

## Area Under the ROC Curve (AUC or AUROC)

**AUC(Area Under the Curve)**: measures the area under the ROC curve. A higher AUC value indicates better model performance as it suggests a greater ability to distinguish between classes. An AUC value of 1.0 indicates perfect performance while 0.5 suggests it is random guessing.

The **AUC** quantifies the **overall ability of a classifier to discriminate** between positive and negative classes

$$\text{AUC} = \int_0^1 TPR(FPR^{-1}(x)) \, dx$$

or simply the **area under the ROC curve.**

- Interpretation:
- AUC = 1 → Perfect separation
- AUC = 0.5 → No discrimination (random guessing)
- AUC < 0.5 → Worse than random (model may be flipped)

AUC represents the **probability that a randomly chosen positive instance** is ranked higher than a randomly chosen negative instance.
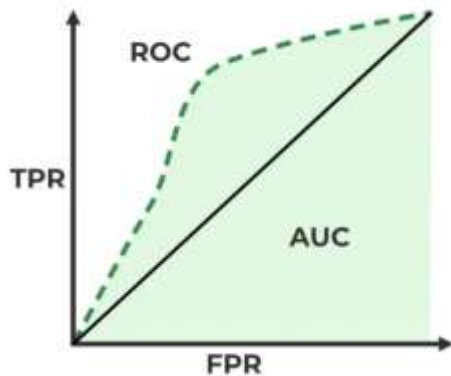
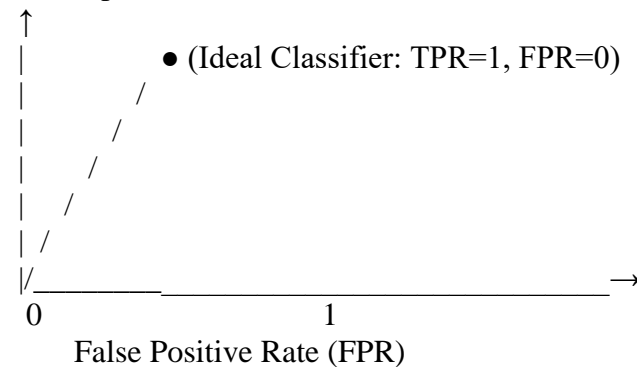$$P(\text{score}(x^+) > \text{score}(x^-))$$

**How AUC-ROC Works**
AUC-ROC curve helps us understand how well a classification model distinguishes between the two classes. Imagine we have 6 data points and out of these:
- **3 belong to the positive class:** Class 1 for people who have a disease.

- **3 belong to the negative class:** Class 0 for people who don't have disease.



ROC Space



AUC represents the **probability that a randomly chosen positive instance** is ranked higher than a randomly chosen negative instance.

$$P(\text{score}(x^+) > \text{score}(x^-))$$

## ROC-AUC Classification Evaluation Metric

Now the model will give each data point a predicted probability of belonging to Class 1. The AUC measures the model's ability to assign higher predicted probabilities to the positive class than to the negative class. Here's how it work:

1. **Randomly choose a pair**: Pick one data point from the positive class (Class 1) and one from the negative class (Class 0).
2. **Check if the positive point has a higher predicted probability**: If the model assigns a higher probability to the positive data point than to the negative one for correct ranking.
3. **Repeat for all pairs**: We do this for all possible pairs of positive and negative examples.

**When to Use AUC-ROC**
AUC-ROC is effective when:
- The dataset is balanced and the model needs to be evaluated across all thresholds.
- False positives and false negatives are of similar importance.

*In cases of highly imbalanced datasets AUC-ROC might give overly optimistic results. In such cases the Precision-Recall Curve is more suitable focusing on the positive class.*
Model Performance with AUC-ROC:

- **High AUC (close to 1)**: The model effectively distinguishes between positive and negative instances.
- **Low AUC (close to 0)**: The model struggles to differentiate between the two classes.
- **AUC around 0.5**: The model doesn't learn any meaningful patterns i.e it is doing random guessing.

In short AUC gives you an overall idea of how well your model is doing at sorting positives and negatives, without being affected by the threshold you set for classification. A higher AUC means your model is doing good.

## Advantages of the ROC Curve

**Threshold-independent** – Evaluates model performance across all thresholds
**Comprehensive** – Shows trade-off between TPR and FPR
**Useful for imbalanced datasets** – Focuses on relative rates, not absolute counts
**Visual diagnostic tool** – Easy to compare multiple models

## Limitations of ROC Curve

☐ **Can be misleading** when data is highly imbalanced
→ Even high TPR and low FPR may not mean good precision
☐ **Does not account for misclassification costs**
→ Sometimes, false positives and false negatives have very different impacts
☐ **Only compares ranking quality**, not calibrated probabilities.
In such cases, **Precision-Recall (PR) curves** can be more informative.

## Selecting an Optimal Threshold

The ROC curve can guide **threshold selection**:

- Choose the point **closest to the top-left corner (0,1)** — minimizes both FPR and FN.
- Can also use **Youden's Index (J)**:

$$J = \text{Sensitivity} + \text{Specificity} - 1$$

Threshold with maximum $J$ gives the best balance.

## Summary

- ROC curve shows **model performance across thresholds**.

- **AUC** summarizes this performance in a single value.

- **Higher AUC ⇒ better ability** to distinguish between classes.

- Use **ROC + AUC** to compare models' ranking quality.

- For **imbalanced datasets**, also consider the **Precision–Recall Curve**.

---

# 8. Precision-Recall Curve
# Work well with imbalance data

The **Precision–Recall Curve (PR Curve)** is a graphical representation used to evaluate the performance of a **binary classification model**, particularly when dealing with **imbalanced datasets** (where one class occurs far more frequently than the other).

While the **ROC Curve (Receiver Operating Characteristic)** is commonly used for general performance evaluation, the **PR Curve** provides **better insight** when the **positive class is rare** or when the **cost of false positives/negatives** differs significantly.

**Terminology**

| Term | Definition | Formula |
|------|-----------|---------|
| True Positive (TP) | Correctly predicted positive samples | — |
| False Positive (FP) | Incorrectly predicted as positive | — |
| True Negative (TN) | Correctly predicted negative samples | — |
| False Negative (FN) | Incorrectly predicted as negative | — |
| Precision (P) | Fraction of predicted positives that are actually positive | $\text{Precision} = \frac{TP}{TP+FP}$ |
| Recall (R) | Fraction of actual positives that are correctly predicted | $\text{Recall} = \frac{TP}{TP+FN}$ |

## Understanding Precision and Recall

### 3.1. Precision

- Also called **Positive Predictive Value (PPV)**.
- Measures how *precise* your positive predictions are.
- High precision means: few false positives.

**Example:**
If a model identifies 100 emails as spam, and 90 are actually spam →
$\text{Precision} = 90/100 = 0.9$

---

### 3.2. Recall

- Also called **Sensitivity** or **True Positive Rate (TPR)**.
- Measures how many of the actual positives your model caught.
- High recall means: few false negatives.

**Example:**
If there are 120 spam emails total, and your model correctly identifies 90 →
$\text{Recall} = 90/120 = 0.75$

---

### 3.3. The Precision–Recall Tradeoff

- There is an **inverse relationship** between Precision and Recall.
- If you increase the model's **threshold**, you predict fewer positives → **Precision increases**, **Recall decreases**.
- If you decrease the **threshold**, you predict more positives → **Recall increases**, **Precision decreases**.

---

# 4. Constructing the Precision–Recall Curve

The PR curve is constructed by **varying the classification threshold** of a model and computing **Precision** and **Recall** for each threshold.

## Steps:

1. Get model scores/probabilities for the positive class.
2. Sort predictions by descending probability.
3. For each threshold:
   - Compute TP, FP, FN.
   - Calculate Precision and Recall.
4. Plot **Recall (x-axis)** vs **Precision (y-axis)**.

---

# 5. Interpreting the PR Curve

- Each point on the curve represents a (Precision, Recall) pair at a specific threshold.
- The **ideal PR curve** reaches the top-right corner (Precision = 1, Recall = 1).

- A **model with random predictions** will have a flat PR curve close to the **baseline**, which equals the proportion of positive samples in the dataset.

---

# 6. The Baseline in PR Curve

The **baseline** (reference line) in the PR curve equals the **ratio of positive instances** in the dataset.

$$\text{Baseline Precision} = \frac{\text{Number of Positives}}{\text{Total Samples}}$$

For example:

- If 10% of the samples are positive → baseline = 0.1.
- A good model should maintain precision significantly **above** this line.
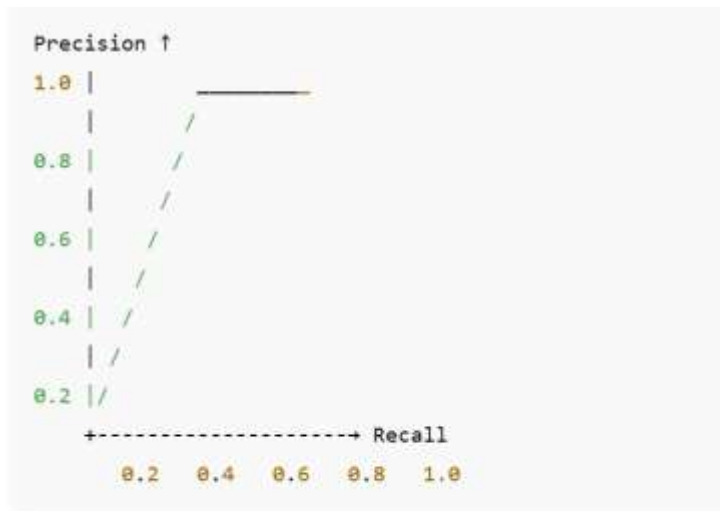
---

# 7. Area Under the PR Curve (AUC–PR)

The **Area Under the Precision–Recall Curve (AUC–PR)** is a single-value summary of the model's performance.

- High AUC–PR → better model.
- AUC–PR = 1 → perfect classifier.
- AUC–PR ≈ baseline → poor classifier.

**Mathematically:**

$$\text{AUC–PR} = \int_0^1 P(R)\, dR$$

Computed numerically using interpolation between points.

```
Precision ↑
1.0 |            _____
    |        /
0.8 |      /
    |    /
0.6 |   /
    |  /
0.4 | /
    |/
0.2 |/
    +-------------------→ Recall
      0.2  0.4  0.6  0.8  1.0
```

The curve shows how Precision decreases as Recall increases.

## Related Metrix

### F1-Score

- Harmonic mean of Precision and Recall.

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

- Useful when you need a **single metric** to balance both.

### Average Precision (AP)

- A smoothed version of AUC–PR, used especially in **object detection (e.g., mAP)**.

# Key Points

- The PR curve highlights **performance on the positive class**.
- It is **robust for imbalanced data** where ROC can be misleading.
- **High AUC–PR** means the model maintains good precision without sacrificing too much recall.
- Use **Average Precision (AP)** or **F1-score** as compact performance summaries.

- Use **ROC curves** for balanced data and general diagnostics.
- Use **Precision–Recall curves** for **imbalanced datasets** or **rare event detection** tasks where **precision** and **recall** are more meaningful than overall accuracy.

# ROC Curve vs. Precision–Recall Curve



PR Curve vs ROC Curve

## Comparison

| Aspect | ROC Curve (Receiver Operating Characteristic) | Precision–Recall (PR) Curve |
|---|---|---|
| Purpose | Evaluates overall classifier performance. | Focuses on performance for the positive class. |
| Use Case | Best for **balanced** datasets. | Best for **imbalanced** datasets. |
| Axes | X-axis: **False Positive Rate (FPR)** Y-axis: **True Positive Rate (TPR)** | X-axis: **Recall (TPR)** Y-axis: **Precision (PPV)** |
| What It Shows | Trade-off between **TPR** and **FPR**. | Trade-off between **Precision** and **Recall**. |

## Summary

| Concept | Formula | Notes |
|---|---|---|
| Precision | $TP/(TP + FP)$ | "Of predicted positives, how many are correct?" |
| Recall | $TP/(TP + FN)$ | "Of actual positives, how many are found?" |
| F1-score | $2PR/(P + R)$ | Harmonic mean of Precision and Recall |
| Baseline | Positives/Total Samples | Reference line in PR curve |
| AUC–PR | Area under Precision–Recall curve | Overall performance measure |

# B. Regression Metrics

Regression models are used to **predict continuous values** — such as prices, temperatures, or scores — rather than categories.
To evaluate how well a regression model performs, we use **error metrics** that quantify the difference between the predicted values and the actual values.

Common regression metrics include:

- **Mean Absolute Error (MAE)**
- **Mean Squared Error (MSE)**
- **Root Mean Squared Error (RMSE)**
- **R² (Coefficient of Determination)**
- **Adjusted R²**

## 1. Mean Absolute Error (MAE)

The **Mean Absolute Error** is the **average of the absolute differences** between the predicted and the actual (true) values.

It measures the **magnitude of errors** in a set of predictions — without considering their direction (positive or negative).

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

where:

- $n$ = total number of observations
- $y_i$ = actual (true) value
- $\hat{y}_i$ = predicted value
- $|y_i - \hat{y}_i|$ = absolute error for each observation

- The **absolute error** for each data point measures **how far** the prediction is from the true value.

- The **mean** of all absolute errors gives an **overall measure** of model accuracy.

- The **smaller** the MAE, the **better** the model fits the data.

Example:

```
| Observation | True Value (y_i) | Predicted (ŷ_i) | Absolute Error |y_i − ŷ_i| |
|--------------|---------------------|--------------------------|---------------------------------
----|
|1|5|4|1|
|2|6|8|2|
|3|3|2|1|
|4|7|7|0|
```

$$\text{MAE} = \frac{1 + 2 + 1 + 0}{4} = 1$$

So, the model's **average absolute error is 1 unit**.

**Characteristics of MAE**

| Property | Explanation |
|---|---|
| **Range** | $[0,\infty)$[0, \infty)$[0,\infty)$ — cannot be negative |
| **Interpretation** | Average magnitude of prediction error |
| **Sensitivity** | Linear — every error contributes proportionally |
| **Robustness** | Less sensitive to outliers compared to MSE |
| **Unit** | Same as the target variable (e.g., ₹, °C, km, etc.) |

---

**6. Geometric Interpretation**

MAE represents the **average vertical distance** between the actual data points and the predicted regression line (or curve).

Visually:

- The distance between each actual point and the regression line is measured.
- All distances are taken as **absolute values** (no cancellation of positive/negative errors).
- The mean of these distances gives MAE.

MAE treats all errors **equally** — whether big or small — unlike MSE or RMSE which **emphasize large errors** more.

Imagine plotting actual vs predicted values:

```pgsql
          |
Actual ↑  |                  *
          |             *    *
          |          *          *
          |       *                *
          |    *                 *
          +-------------------------→ Predicted
          (Perfect line y = x)
```

- The **vertical distances** between points and the 45° line represent the **absolute errors**.
- The **average** of those distances is the **MAE**.

# When to Use MAE

Use **Mean Absolute Error** when:

- You want a **simple, interpretable** measure of prediction accuracy.
- Outliers are not very influential.
- Equal penalty for over- and under-predictions is desired.

- MAE provides a **clear, intuitive measure** of average prediction error.

- It's **robust**, easy to explain to non-technical audiences, and valuable for most regression tasks.

- However, if **large errors are particularly undesirable**, consider **MSE** or **RMSE** instead.

- You want to evaluate models where **robustness** is important.

# 2. Mean Squared Error (MSE)

The **Mean Squared Error (MSE)** measures the **average of the squares** of the errors — that is, the average squared difference between the actual values and the predicted values.

It tells us **how close** a regression line (or model) fits the actual data points.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

where:

- $n$ = total number of observations
- $y_i$ = actual (true) value
- $\hat{y}_i$ = predicted value
- $(y_i - \hat{y}_i)$ = residual or prediction error

- MSE squares each error term — so **larger errors have a disproportionately higher impact**.

- The squaring ensures **all errors are positive** (no cancellation of positive and negative errors).

- The lower the MSE, the **better** the model performance.

**Example**

| Observation | Actual ($y_i$) | Predicted ($\hat{y}_i$) | Error ($y_i - \hat{y}_i$) | Squared Error |
|---|---|---|---|---|
| 1 | 5 | 4 | 1 | 1 |
| 2 | 6 | 8 | -2 | 4 |
| 3 | 3 | 2 | 1 | 1 |
| 4 | 7 | 7 | 0 | 0 |

$$\text{MSE} = \frac{1 + 4 + 1 + 0}{4} = 1.5$$

So the model's **mean squared error = 1.5**.

## Characteristics of MSE

| Property | Explanation | |
|---|---|---|
| Range | $[0, \infty)$ | |
| Best Value | 0 (perfect predictions) | |
| Units | Square of the target variable (e.g., if target = ₹, MSE = ₹²) | |
| Sensitivity | Highly sensitive to large errors (due to squaring) | |
| Error Penalization | Quadratic — large errors are penalized more severely | |

## Advantages of MSE

**Mathematically convenient** — differentiable everywhere, useful for optimization (e.g., Gradient Descent).
**Emphasizes large errors** — good if big mistakes are costly (e.g., medical or financial predictions).
**Smooth and convex loss function** — ensures a unique global minimum.

## Limitations of MSE

☐ **Not interpretable in original units** — results are in squared units (e.g., ₹²).
☐ **Very sensitive to outliers** — a few large errors can inflate MSE dramatically.
☐ **Not robust** — if data has noisy points or outliers, MAE may be a better alternative.

MSE is used as a **loss function** in regression models such as:

- Linear Regression
- Ridge and Lasso Regression
- Neural Networks (for continuous output)
- Support Vector Regression (SVR) with squared loss

Optimization goal:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i; \theta))^2$$

where $f(x_i; \theta)$ is the model prediction.
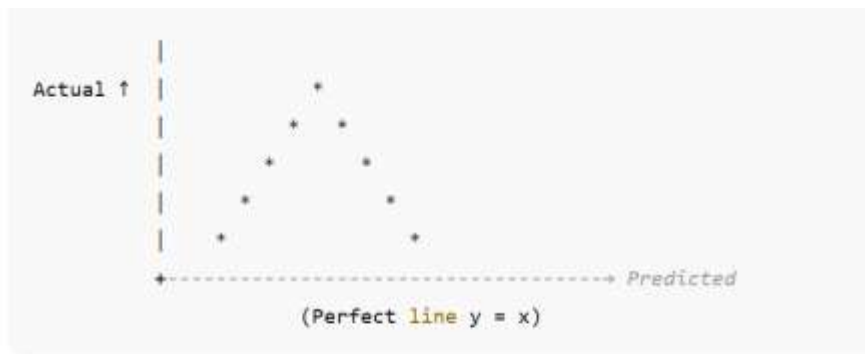
**Mathematical Properties**

1. **Always non-negative:**

$$(y_i - \hat{y}_i)^2 \geq 0$$

2. **Zero only for perfect prediction.**

3. **Convex function** — guarantees a single global minimum (good for gradient descent).

4. **Differentiable** — makes optimization efficient:

$$\frac{\partial \text{MSE}}{\partial \hat{y}_i} = -\frac{2}{n}(y_i - \hat{y}_i)$$

5. **Quadratic growth** — errors increase rapidly as deviations grow.

```
          |
          |
Actual ↑  |            .
          |         .     .
          |      .      .
          |    .      .
          |  .       .
          | .          .
          +- - - - - - - - - - - - - - - - - -→ Predicted
              (Perfect line y = x)
```

- Each data point's **vertical deviation** from the line is squared.
- The **average of these squared distances** = MSE.

Larger deviations contribute **disproportionately more** to the total error.

# When to Use MSE

Use **MSE** when:

- You want to **penalize large errors heavily**.
- Your data **does not have significant outliers**.
- You're training models using **gradient-based optimization**.
- You care more about **mathematical convenience and differentiability**.

Avoid MSE when:

- The dataset contains **outliers** or **noise** — use **MAE** instead.

# Key Points

- **MSE** measures the **average squared deviation** between predicted and true values.
- It is **smooth, differentiable, and convex**, making it ideal for training models via gradient descent.
- However, it **magnifies large errors**, making it less robust to outliers.
- Always pair it with **RMSE** for better interpretability.

-------------------------------------------------------------------

# 3. Root Mean Squared Error (RMSE)

The **Root Mean Squared Error (RMSE)** is defined as the **square root of the average of squared differences** between predicted values and actual (true) values.

It provides a single measure that summarizes how far, on average, the model's predictions are from the true outcomes.

**RMSE** measures how much predictions deviate from actual values, on average, in the **same units as the output**, giving a clear picture of model accuracy.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$$

where:

- $n$ = total number of observations
- $y_i$ = actual value
- $\hat{y}_i$ = predicted value
- $(y_i - \hat{y}_i)$ = residual or prediction error

- Each prediction error $(y_i - \hat{y}_i)$ is squared → larger errors have more influence.
- The mean of these squared errors gives the MSE.
- Taking the square root **brings the metric back to the same scale** as the target variable.

Hence, RMSE tells you:

> "On average, how much do predictions deviate from the actual values?"

RMSE is directly derived from MSE:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

$$\text{RMSE} = \sqrt{\text{MSE}}$$

Thus, RMSE can be viewed as a **scaled (unit-corrected)** version of MSE.

While MSE provides mathematical convenience, **RMSE is easier to interpret** because it has the **same units as the original data**.

| Observation | Actual ($y_i$) | Predicted ($\hat{y}_i$) | Error ($y_i - \hat{y}_i$) | Squared Error |
|---|---|---|---|---|
| 1 | 5 | 4 | 1 | 1 |
| 2 | 6 | 8 | -2 | 4 |
| 3 | 3 | 2 | 1 | 1 |
| 4 | 7 | 7 | 0 | 0 |

$$\text{MSE} = \frac{1 + 4 + 1 + 0}{4} = 1.5$$

$$\text{RMSE} = \sqrt{1.5} \approx 1.225$$

☑ Interpretation:

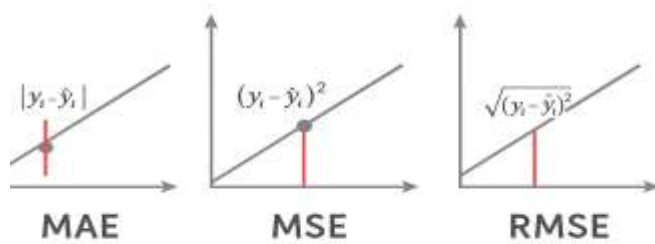On average, the predictions deviate from the true values by **1.225 units**.

## Characteristics of  RMSE

| Property | Explanation |
|---|---|
| Range | $[0, \infty)$ |
| Best Value | 0 (perfect predictions) |
| Units | Same as target variable |
| Error Penalization | Quadratic — larger errors penalized more heavily |
| Sensitivity | Highly sensitive to outliers |
| Interpretation | Average magnitude of error, with higher emphasis on large errors |

## Comparison: RMSE vs MAE vs MSE

| Metric | Formula | Penalty Type | Outlier Sensitivity | Units | Interpretation |
|---|---|---|---|---|---|
| MAE | ( \frac{1}{n}\sum | y_i - \hat{y}_i ) | | Linear | Low |
| MSE | $\frac{1}{n} \sum (y_i - \hat{y}_i)^2$ | Quadratic | High | Squared units | Average squared error |
| RMSE | $\sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$ | Quadratic | High | Same as target | Standard deviation of prediction errors |

• **MAE** is more robust and interpretable.
• **RMSE** is more sensitive and highlights large deviations.

$|y_i - \hat{y}_i|$     $(y_i - \hat{y}_i)^2$     $\sqrt{(y_i - \hat{y}_i)^2}$

MAE     MSE     RMSE

# Advantages of RMSE

☐ **Interpretable** — in same units as target variable.
☐ **Differentiable and convex**, ideal for optimization in ML models.
☐ **Strongly penalizes large errors**, making it useful when big mistakes are undesirable (e.g., forecasting, finance, healthcare).
☐ **Standard metric** — used in competitions, literature, and industry benchmarks.

## Limitations of RMSE

☐ **Highly sensitive to outliers** — a single large error can dominate the metric.
☐ **Non-linear error contribution** — not all errors are treated equally.
☐ **Difficult to decompose** — unlike MAE, RMSE can't easily show direction or magnitude separately.
☐ **Not robust** in datasets with noise or extreme values.

## Key Points

- **RMSE** is a powerful and interpretable regression metric.
- It combines the **mathematical properties** of MSE with **practical interpretability**.
- It penalizes large errors more severely, which can improve model robustness in precise prediction tasks.
- However, it can be **distorted by outliers**, so context matters.

# When to Use RMSE

**Use RMSE when:**

- You want a **sensitive metric** that highlights large errors.
- All errors have the **same cost functionally**.
- You need a **standardized measure** in the same units as the dependent variable.
- The dataset **has no extreme outliers**.

**Avoid RMSE when:**

- Dataset has **heavy outliers** (use MAE instead).
- You want an **error measure less biased by large deviations**.

# 4. R-squared (R² Score)

In regression analysis, we aim to build a model that **explains the relationship** between input features XXX and a continuous target variable yyy.

After building the model, we need to evaluate **how well it fits the data** — that is, how much of the variation in yyy can be explained by the model.

The **R-squared**, also known as the **Coefficient of Determination**, is a **statistical measure** that tells us **how well the regression predictions approximate the real data points**.

---

The **R-squared (R²)** measures the **proportion of the variance** in the dependent variable yyy that is **explained** by the independent variables XXX in the model.

It gives an idea of **the goodness of fit** of the model.

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

where:

- $SS_{res} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ → Residual Sum of Squares (SSR)
- $SS_{tot} = \sum_{i=1}^{n} (y_i - \bar{y})^2$ → Total Sum of Squares (SST)
- $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$ → mean of actual values

- $SS_{tot}$: total variation in $y$ — how much the values vary from their mean.
- $SS_{res}$: unexplained variation — the part not captured by the model.
- $SS_{tot} - SS_{res}$: explained variation — the part captured by the model.

Thus, $R^2$ quantifies the **fraction of total variation explained** by the model.

$$R^2 = \frac{\text{Explained Variation}}{\text{Total Variation}}$$

| Range | Interpretation |
|---|---|
| $R^2 = 1$ | Perfect fit — model explains 100% of variance |
| $R^2 = 0$ | Model explains none of the variance |
| $R^2 < 0$ | Model performs worse than a simple mean-based model |

- The closer $R^2$ is to 1, the better the model fits the data.
- A **negative R²** occurs if the model's predictions are **worse than just predicting the mean** for all data points.

**Example**

| Observation | Actual ($y_i$) | Predicted ($\hat{y}_i$) | Mean ($\bar{y}$) |
| --- | --- | --- | --- |
| 1 | 3 | 2.5 | 4 |
| 2 | 5 | 5.5 | 4 |
| 3 | 7 | 6 | 4 |

## Step 1: Compute SST (Total Sum of Squares)

$$SS_{tot} = (3-4)^2 + (5-4)^2 + (7-4)^2 = 1 + 1 + 9 = 11$$

## Step 2: Compute SSR (Residual Sum of Squares)

$$SS_{res} = (3-2.5)^2 + (5-5.5)^2 + (7-6)^2 = 0.25 + 0.25 + 1 = 1.5$$

## Step 3: Compute R²

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{1.5}{11} = 1 - 0.136 = 0.864$$
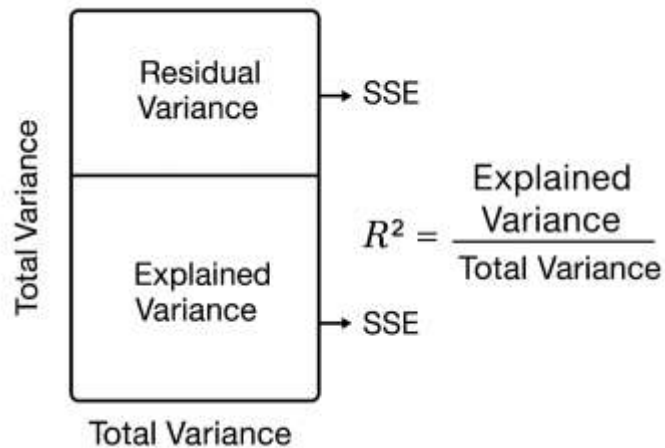
✅ Interpretation:

The model explains **86.4% of the variation** in the target variable.

## Geometric Interpretation

Imagine plotting actual values $y_i$ against predicted values $\hat{y}_i$:

- A **perfect fit** would lie exactly on the line $y = \hat{y}$.
- Deviations from this line represent **residuals** (errors).
- $R^2$ measures how tightly the points cluster around that perfect-fit line.

$$R^2 = \frac{\text{Explained Variance}}{\text{Total Variance}}$$

**Visual Understanding**

- If points lie **close to the regression line**, $R^2$ is high.
- If points are **scattered far away**, $R^2$ is low.
- A random prediction line or a horizontal mean line gives $R^2 \approx 0$.

# Properties of R²

| Property | Explanation |
|---|---|
| Range | (-∞, 1] |
| Best Value | 1 |
| Worst Case | Negative (model worse than mean) |
| Interpretation | Proportion of variance explained |
| Unit | Dimensionless |
| Dependency | Depends on number of predictors — increases with more variables |

# Advantages

☐ **Intuitive interpretation** — tells how well model explains variability.
☐ **Dimensionless** — independent of data units.
☐ **Widely used** — simple benchmark for comparing models.
☐ **Easy to compute** — derived from sums of squares.

# Limitations

⚠ **Cannot detect overfitting** — R² increases even if new variables don't add real predictive value.

⚠ **Not suitable for non-linear models** — only measures linear fit strength.

⚠ **Doesn't indicate bias** — model may have a high R² but poor predictive accuracy.

⚠ **Sensitive to outliers** — outliers can inflate R².

### R² and Correlation Coefficient

For **simple linear regression** (one variable), R² is the **square of the Pearson correlation coefficient (r)** between the predicted and actual values:

$$R^2 = r^2$$

This means:

- $R^2$ captures how strongly the model's predictions correlate with true outcomes.
- $r$ can be positive or negative; $R^2$ is always non-negative.

# Relationship with Error Metrics

| Metric | Meaning | Ideal Value |
|--------|---------|-------------|
| **MAE** | Average magnitude of errors | 0 |
| **MSE** | Average squared errors | 0 |
| **RMSE** | Standard deviation of prediction errors | 0 |
| **R²** | Proportion of explained variance | 1 |

⬥ While MAE, MSE, and RMSE measure **error magnitude**,

⬥ R² measures **explanatory power**.

# When to Use R²

☐ Use R² when:

- Evaluating **linear regression** models.
- Comparing models with **the same dependent variable**.
- Assessing **goodness of fit** (not prediction accuracy).

☐ Avoid R² when:

- The model is **nonlinear**.
- You're evaluating **out-of-sample performance** (use RMSE instead).
- You have **imbalanced data** or **heteroscedasticity**.

## High R² (≈ 0.9–1.0):

Points are tightly clustered around the regression line.

## Low R² (≈ 0.0–0.3):

Points are widely scattered — model explains little variance.

```lua
High R²                    Low R²
 |  * * * * *                *         *
 |   * * *                 *     *        *
 |    *                    *     *
 +- - - - - - - -         +- - - - - - - - -
```

---

## 5. Adjusted R²

To overcome overfitting problems, we use **Adjusted R²**, which penalizes the addition of unnecessary variables.

$$R_{adj}^2 = 1 - \left( \frac{(1 - R^2)(n - 1)}{n - k - 1} \right)$$

where:

- $n$ = number of observations
- $k$ = number of predictors

**Key Point:**

- If a new variable **improves** the model, Adjusted R² increases.
- If it **does not help**, Adjusted R² decreases.

# Typical R² Ranges (Rules of Thumb)

**R² Range Interpretation**
0.9 – 1.0   Excellent fit

**R² Range Interpretation**

0.7 – 0.9  Good fit

0.5 – 0.7  Moderate fit

0.3 – 0.5  Weak fit

< 0.3  Poor fit

| Aspect | R-squared (R² Score) |
|---|---|
| Definition | Proportion of variance in target explained by model |
| Formula | $R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$ |
| Range | $(-\infty, 1]$ |
| Best Value | 1 |
| Units | Dimensionless |
| Purpose | Measures model's goodness of fit |
| Improves With | Better explanatory power |
| Related To | MSE, correlation coefficient |
| Variants | Adjusted R² for multiple regression |

# Key Points

- **R² measures goodness of fit**, not prediction accuracy.
- **High R²** means the model explains a large portion of the data's variance.
- **R² = 0** means model is no better than predicting the mean.
- **Adjusted R²** corrects for adding irrelevant variables.
- Always interpret R² **alongside RMSE or MAE** for a complete evaluation.

---

# Metrics for Clustering Problems: Silhouette Score, Davies-Bouldin Index

# 1. Introduction

Unlike supervised learning, clustering is **unsupervised**, meaning there are **no ground-truth labels**.
Thus, we can't use standard metrics like accuracy, precision, or R².

To evaluate how well our clusters are formed, we use **internal evaluation metrics**, which rely only on the data and the clustering results.

Two of the most widely used internal metrics are:

- **Silhouette Score**
- **Davies–Bouldin Index**

Both help assess the **quality, compactness, and separation** of clusters.

---

# 2. Goals of Clustering Evaluation

A good clustering should have:

1. **High intra-cluster similarity** → points within the same cluster are close together.
2. **Low inter-cluster similarity** → clusters are well separated from each other.

Hence, a good clustering algorithm aims for:

- **Cohesion (compactness):** how close points in the same cluster are.
- **Separation:** how far clusters are from each other.

---

# 3. Silhouette Score

---

## 3.1 Definition

The **Silhouette Score (s)** measures how well each data point fits within its assigned cluster compared to other clusters.

It combines both **cohesion** and **separation** into a single value.

For each data point $i$:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

Where:

- $a(i)$: average distance between $i$ and all other points in the **same cluster** (intra-cluster distance).
- $b(i)$: smallest average distance from $i$ to all points in **other clusters** (nearest-cluster distance).

# 3.2 Range and Interpretation

| Silhouette Value (s(i)) | Interpretation |
|---|---|
| +1 | Perfectly matched to its own cluster, far from others |
| 0 | On or very close to the boundary between clusters |
| −1 | Misclassified; assigned to the wrong cluster |

The **overall Silhouette Score** for a dataset is the **mean** of all individual $s(i)$ values.

$$S = \frac{1}{n} \sum_{i=1}^{n} s(i)$$

- $a(i)$: measures **tightness** of the cluster (lower is better).
- $b(i)$: measures **separation** between clusters (higher is better).
- Thus, a large positive $s(i)$ means good clustering, while a negative one indicates overlap or misassignment.

**Silhouette Score Range Clustering Quality**
| | |
|---|---|
| 0.71 – 1.00 | Excellent |
| 0.51 – 0.70 | Good |
| 0.26 – 0.50 | Fair |
| $\leq 0.25$ | Poor |

Suppose we cluster data into two groups:

- For a data point $i$:

$a(i) = 2.0$ (average distance to points in same cluster)

$b(i) = 5.0$ (average distance to nearest cluster)

$$s(i) = \frac{5 - 2}{\max(2, 5)} = \frac{3}{5} = 0.6$$

✅ This point is **well-clustered**, as 0.6 indicates a clear boundary between clusters.

Silhouette plots show the **distribution of silhouette scores per cluster**.

- Wide bars indicate well-separated clusters.
- Negative or small bars indicate overlapping clusters.

---

### 3.7 Advantages

☐ Intuitive and easy to interpret.
☐ Works with any distance metric (Euclidean, Manhattan, cosine, etc.).
☐ Can be used to select the **optimal number of clusters (k)** — the k that maximizes the silhouette score.

---

### 3.8 Limitations

⚠☐ Computationally expensive for large datasets (requires distance between many points).
⚠☐ May not work well with **non-convex** clusters.
⚠☐ Sensitive to noise and outliers.

# 4.Davies–Bouldin Index (DBI)

The **Davies–Bouldin Index** measures the **average similarity** between each cluster and its most similar one.

It considers both **intra-cluster distance (compactness)** and **inter-cluster distance (separation)**.

$$DBI = \frac{1}{k} \sum_{i=1}^{k} \max_{j \neq i} \left( \frac{S_i + S_j}{M_{ij}} \right)$$

Where:

- $k$ = number of clusters
- $S_i$ = average distance of all points in cluster $i$ to its centroid (cluster scatter)
- $M_{ij}$ = distance between centroids of clusters $i$ and $j$

## 4.2 Intuitive Meaning

- $S_i$: how **spread out** cluster $i$ is
- $M_{ij}$: how **far apart** clusters $i$ and $j$ are
- The ratio $\frac{S_i + S_j}{M_{ij}}$ measures how **similar** the clusters are.
- The **maximum** over $j$ finds the most **similar (closest)** cluster to $i$.
- The **average** of all these values gives DBI.

### 4.3 Range and Interpretation

| Davies–Bouldin Index (DBI) | Interpretation |
|---|---|
| 0 | Perfect clustering (ideal case) |
| < 1 | Good separation and compact clusters |
| > 1 | Poor clustering; overlapping or scattered clusters |

⬍ **Lower DBI is better.**

# 4.4 Example (Conceptual)

Suppose we have 3 clusters:

| Cluster | $S_i$ | Closest Cluster Distance $M_{ij}$ | $R_i = \max_j \frac{S_i + S_j}{M_{ij}}$ |
|---------|-------|-----------------------------------|------------------------------------------|
| $C_1$ | 0.5 | 3.0 | 0.4 |
| $C_2$ | 0.6 | 2.8 | 0.5 |
| $C_3$ | 0.4 | 3.2 | 0.3 |

$$DBI = \frac{0.4 + 0.5 + 0.3}{3} = 0.4$$

☑ Interpretation: **Low DBI = good clustering quality.**

### 4.5 Advantages

☐ Quantifies **both compactness and separation**.
☐ **Lower is better**, easy to compare between models.
☐ Works with any clustering algorithm (K-means, DBSCAN, etc.).

### 4.6 Limitations

⚠ Sensitive to **cluster shape and outliers**.
⚠ Assumes clusters are **convex and isotropic** (K-means assumption).
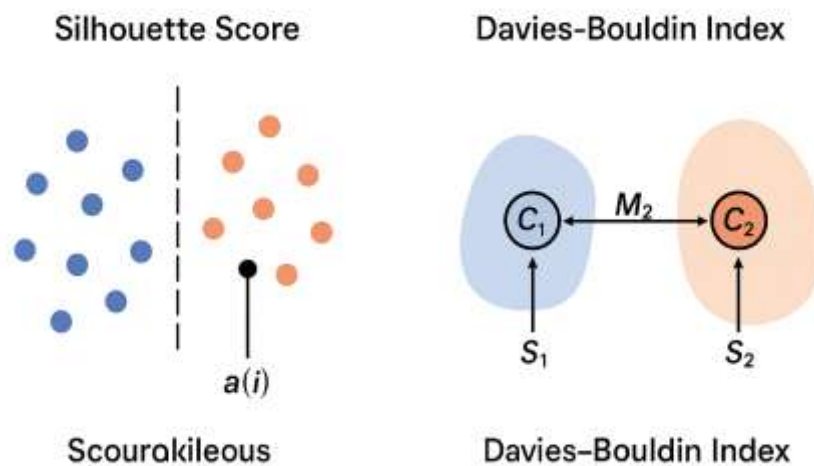⚠ Not intuitive compared to Silhouette Score.

## 5. Comparison: Silhouette Score vs Davies–Bouldin Index

| Aspect | Silhouette Score | Davies–Bouldin Index |
|--------|------------------|----------------------|
| **Definition** | Measures how similar a point is to its cluster vs others | Measures average similarity between clusters |
| **Range** | −1 to +1 | 0 to ∞ |
| **Ideal Value** | Closer to +1 | Closer to 0 |
| **High Value Means** | Better clustering | Worse clustering |
| **Measures** | Both cohesion & separation | Compactness vs separation |
| **Computation** | Based on all pairwise distances | Based on centroid distances |
| **Ease of Interpretation** | Very intuitive | Less intuitive |

| Aspect | Silhouette Score | Davies–Bouldin Index |
|---|---|---|
| Use Case | Evaluating overall clustering quality, finding optimal k | Comparing model compactness and separation efficiency |

## Metrics for Clustering Problems



Silhouette Score — Scourakileous

Davies–Bouldin Index

# 6. Choosing Between Them

| If you want... | Use |
|---|---|
| A **visual, interpretable** metric | Silhouette Score |
| A **compact, mathematical** metric for algorithmic comparison | Davies–Bouldin Index |
| To **tune the number of clusters (k)** | Silhouette Score |
| To **evaluate compactness and separation** numerically | DBI |

## Summary Table

| Metric | Ideal Value | Range | Measures | High Value Means | Low Value Means |
|---|---|---|---|---|---|
| **Silhouette Score** | +1 | −1 to +1 | Cohesion + Separation | Good clustering | Poor clustering |
| **Davies–Bouldin Index** | 0 | 0 to ∞ | Compactness + Separation | Poor clustering | Good clustering |

# Key Points

- **Silhouette Score** and **Davies–Bouldin Index** are **internal validation metrics** — they do not require labels.

- **Silhouette Score**:
  - o Higher = better.
  - o Intuitive and widely used.
- **Davies–Bouldin Index**:
  - o Lower = better.
  - o Quantifies similarity among clusters.
- Both metrics are useful for **comparing clustering results** and **tuning k** in algorithms like K-Means.

--------------------------------------------------------

# Ranking / Information Retrieval Metrics

## A. Precision@k

- Precision calculated for top-k predicted items.

## B. Recall@k

- Recall calculated for top-k results.

## C. Mean Average Precision (MAP)

- Averages precision over multiple queries.

## D. NDCG (Normalized Discounted Cumulative Gain)

- Measures ranking quality with position-based relevance:
  $$\text{DCG}_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i+1)}$$
  $$\text{NDCG}_p = \frac{DCG_p}{IDCG_p}$$

# Choosing the Right Metric

| Task Type | Common Metrics |
|---|---|
| Binary Classification | Accuracy, Precision, Recall, F1, ROC-AUC |
| Multi-class Classification | Macro/Micro-averaged F1, Confusion Matrix |
| Imbalanced Classification | F1 Score, PR Curve, ROC-AUC |
| Regression | MAE, RMSE, R², Adjusted R² |

| Task Type | Common Metrics |
|---|---|
| Ranking | NDCG, MAP, Precision@k |

# Cross-Validation for Robust Evaluation

- **K-Fold Cross-Validation**
- **Stratified K-Fold** (for classification)
- Reduces variance and bias in model evaluation.

# Common Pitfalls

- Relying only on accuracy (especially with class imbalance).
- Ignoring domain-specific costs (e.g., false negatives in cancer detection).
- Not validating metrics on **unseen** test data.

# Summary

- No single metric is best for all problems.
- Always consider **context, data imbalance, and business impact**.
- Combine multiple metrics for a complete evaluation.