

# Evaluation Metrics for Classification Models

21 November 2025 12:01

Evaluation metrics tell you how well your classification model performs. The choice of metric depends on the problem, dataset imbalance, and real-world constraints.

This tutorial covers:

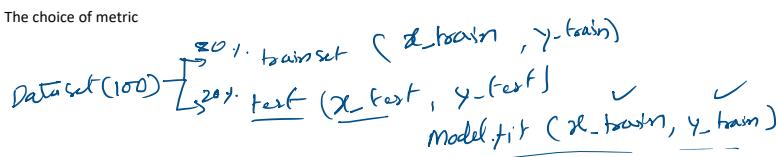
1. Core Concepts
2. Confusion Matrix
3. Basic Metrics (Accuracy, Precision, Recall, F1)
4. Advanced Metrics (AUC, ROC, PR curves, Log Loss)
5. Metrics for Imbalanced Data
6. Metrics for Multi-class & Multi-label Classification

## 1. ★ Core Concepts

Classification problems output **discrete labels** (e.g., spam/not spam, disease/healthy, object categories).

A model's performance is measured using:

- Predicted labels  $\rightarrow \hat{y}$
- Ground-truth labels  $\rightarrow y_{\text{true}}$
- Probabilities (for probabilistic metrics)  $\rightarrow \text{prob}$



$x_1$	$x_2$	$y$
age	salary	purchse
22	35,000	0
25	45,000	1
23	42,000	1
30	50,000	0
35	45,000	0
30	60,000	1

$y_{\text{pred}} = \text{Model} \cdot \text{predict}(X_{\text{test}})$

$y_{\text{test}} \rightarrow \text{Actual value/Ground truth}$

$y_{\text{pred}} \rightarrow \text{Predicted value/Labels}$

$y_{\text{test}} \left\{ \begin{array}{l} 1 \rightarrow \text{purchase} \\ 0 \rightarrow \text{not purchase} \end{array} \right.$

## 2. ⚗ Confusion Matrix — The Foundation

The confusion matrix is a  $2 \times 2$  table for binary classification.

		Predicted Value		Accuracy
		Yes	No	
Actual Value	Yes	TP	FN	Recall
	No	FP	TN	Specificity
Prevalence	4	16	30	6
Precision	4	16	10	6
Negative Predictive Value	6	144	2	6

From this table, all major metrics are derived.

## 3. 🔍 Basic Metrics

### 3.1 Accuracy

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

$$\frac{4+6}{12} = \frac{10}{12} = 0.66$$

$1 \rightarrow +ve$   
 $0 \rightarrow -ve$

✓ Good when classes are **balanced**

✗ Misleading for **imbalanced data**

Example: 99% accuracy when 99% samples belong to one class.

### 3.2 Precision

“How many predicted positives were actually positive?”

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

Useful when:

- False positives are costly (e.g., spam filter marking real mail as spam).

### 3.3 Recall (Sensitivity/TPR)

“How many actual positives did we correctly identify?”

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

Useful when:

- Missing positives is costly (e.g., cancer detection).

### 3.4 F1-Score

Harmonic mean of precision and recall.

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

It is useful when:

- You want a balance between precision and recall.
- The dataset is imbalanced.

For **multi-class** and **multi-label** problems, F1-score comes in different variants.

#### 1. Binary F1-Score

Used when:

- You have two classes (e.g., positive vs negative)
- You treat one class as “positive”

## 2. Macro F1-Score

$$\text{Macro } F1 = \frac{F1_1 + F1_2 + \dots + F1_k}{k}$$

- Computes F1 for each class independently
  - Averages them equally
  - Does NOT consider class imbalance
- ✓ Good when class performance should be treated equally  
✗ Not suitable for heavily imbalanced datasets

## 3. Micro F1-Score

$$\text{Micro } F1 = \frac{2 \times (TP_{total})}{2 \times TP_{total} + FP_{total} + FN_{total}}$$

- Aggregates global TP, FP, FN across all classes
  - Computes precision and recall at the global level
  - Equivalent to accuracy in multi-class classification
- ✓ Good when each sample has equal importance  
✗ Not good when you need per-class fairness

## 4. Weighted F1-Score

$$\text{Weighted } F1 = \sum_{i=1}^k w_i \times F1_i$$

where

$$w_i = \frac{\text{support of class } i}{\text{total samples}}$$

- Computes F1 for each class
  - Weighted by how many samples each class has (support)
- ✓ Best for imbalanced datasets  
✓ Still gives per-class insight

✗ Majority class can dominate the metric

## Summary Table of F1-Score Types

Type	Use Case	Considers Class Imbalance?	Notes
<b>Binary F1</b>	Single positive class	✗	For binary classification only
<b>Macro F1</b>	Equal class importance	✗	Good for fairness across classes
<b>Micro F1</b>	Overall performance	✓ (implicitly)	Treats all samples equally
<b>Weighted</b>	Imbalanced classes	✓	Most commonly used in real tasks
<b>F1</b>			

## Which F1 Should You Use?

Task Type	Recommended F1
Binary classification	Binary F1
Balanced multi-class	Macro
Imbalanced multi-class	Weighted
Want a single score emphasizing performance on all samples	Micro

## 4. Advanced Metrics

### 4.1 ROC Curve

- Plots True Positive Rate (Recall) vs False Positive Rate
- Shows model performance across thresholds.

#### AUC-ROC

Area under ROC curve

$AUC \approx P(\text{positive class has higher score than negative class})$

✓ Best for:

- Balanced data
- Binary classification

### 4.2 Precision–Recall Curve

Plots Precision vs Recall at different thresholds.

✓ Better than ROC for highly imbalanced data.

#### Average Precision (AP)

Area under PR curve.

### 4.3 Log Loss (Cross-Entropy Loss)

Measures how well probabilities match actual labels.

$$\text{LogLoss} = -\frac{1}{N} \sum (y \log(p) + (1-y) \log(1-p))$$

✓ Probabilistic metric

✓ Lower is better

Used in:

- Kaggle competitions
- Well-calibrated probability models

## 5. Metrics for Imbalanced Classification

For skewed datasets (fraud detection, medical diagnosis), use:

Sensitivity (Recall)

$$TP / (TP + FN)$$

Specificity

$$TN / (TN + FP)$$

Balanced Accuracy

$$\frac{\text{Recall}_{\text{positive}} + \text{Recall}_{\text{negative}}}{2}$$

### F2 or F0.5 Score

Weights recall higher or precision higher.

### Matthews Correlation Coefficient (MCC)

Best single metric for imbalanced problems.

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

✓ Range: -1 to +1

✓ Robust even when classes are extreme (e.g., 99:1)

## 6. Multi-Class Metrics

For multi-class classification, metrics are computed using:

- Macro (unweighted mean of per-class scores)
- Micro (global counts)
- Weighted (weighted by class frequency)

Common metrics:

- Macro F1
- Weighted accuracy
- Top-k accuracy (e.g., Top-5 accuracy used in ImageNet)
- Cohen's Kappa (measures inter-class agreement)

## 7. Multi-Label Classification Metrics

In multi-label problems, predictions assign **multiple labels** to one sample.

Metrics include:

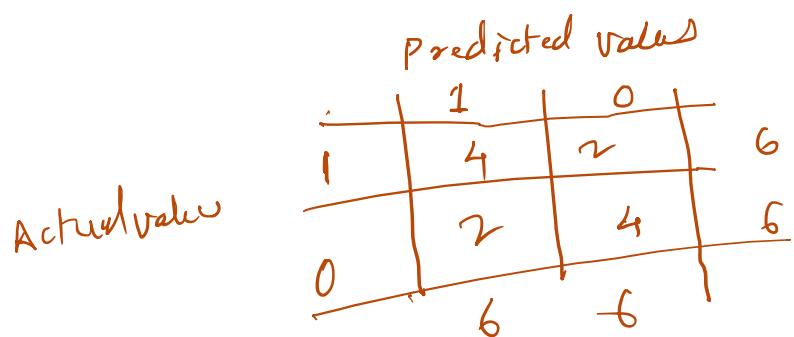
- Hamming Loss
- Jaccard Index
- Subset Accuracy (strictest—exact match required)

## Summary Table

Metric	Good For	Bad For	Notes
Accuracy	Balanced data	Imbalanced data	Overall correctness
Precision	FP costly	FN costly	"How precise are positive predictions?"
Recall	FN costly	FP acceptable	"How many positives did we catch?"
F1	Precision = Recall	Probabilistic metrics	Harmonic mean
AUC-ROC	Balanced binary	Extreme imbalance	Rank-based
AUC-PR	Imbalanced data	Balanced	Focus on minority class
Log Loss	Probability calibration	Hard-label only	Penalizes wrong confidence
MCC	Imbalanced data	Very large datasets	Best single metric

Patient	True Label (Actual y)	Predicted Probability ( $\hat{y}$ )	$\hat{y} > 0.5$	Pred label	
1	1	0.95	True	1	
2	0	0.85	True	1	
3	1	0.90	True	1	
4	0	0.10	False	0	
5	1	0.80	True	1	
6	0	0.05	False	0	
7	1	0.60	True	1	
8	1	0.30	False	0	
9	0	0.70	True	1	
10	0	0.40	False	0	
11	1	0.20	False	0	
12	0	0.25	False	0	

$\rightarrow TP \rightarrow 4$   
 $\rightarrow TN \rightarrow 4$   
 $\rightarrow FP \rightarrow 2$   
 $\rightarrow FN \rightarrow 2$



Accuracy      Instance      (1800)  
 $(950 \rightarrow \text{Positive (1)})$   
 $50 \rightarrow \text{Negative (0)}$

Predicted  $\xrightarrow{1000}$  Positive (1)

$TP \rightarrow 950$

$FP \rightarrow 50$

$FN \rightarrow 0$

$TN \rightarrow 0$

$$\text{Acc} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{950 + 0}{950 + 0 + 50 + 0} = \frac{950}{1000} = 95\%$$

$$\frac{TP + TN}{TP + TN + FP + FN} = \frac{950}{1000} = 95\%$$

ROC

Patient	True Label (Actual y)	Predicted Probability ( $\hat{y}$ )	$\hat{y} > 0.1$	$\hat{y} > 0.2$	$\hat{y} > 0.3$	$\hat{y} > 0.9$
1	1	0.95	1	1	1	1
2	0	0.85	1	1	1	1
3	1	0.90	1	1	1	1
4	0	0.10	0	0	0	0
5	1	0.80	1	1	1	1
6	0	0.05	0	0	0	0
7	1	0.60	1	1	1	1
8	1	0.30	1	1	1	0
9	0	0.70	1	1	1	1
10	0	0.40	1	0	1	1
11	1	0.20	1	0	0	0
12	0	0.25	1	1	0	0

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

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

$$\hat{y} > 0.1 \quad TPR = \frac{6}{6+0} = 1$$

$$\hat{y} > 0.2 \quad FPR = \frac{4}{2+6} = 0.66$$

$$\hat{y} > 0.3 \quad FP \rightarrow 4$$

$$\hat{y} > 0.3 \quad TP = 5 \quad TPR = \frac{5}{5+1} = \frac{5}{6} = 0.83$$

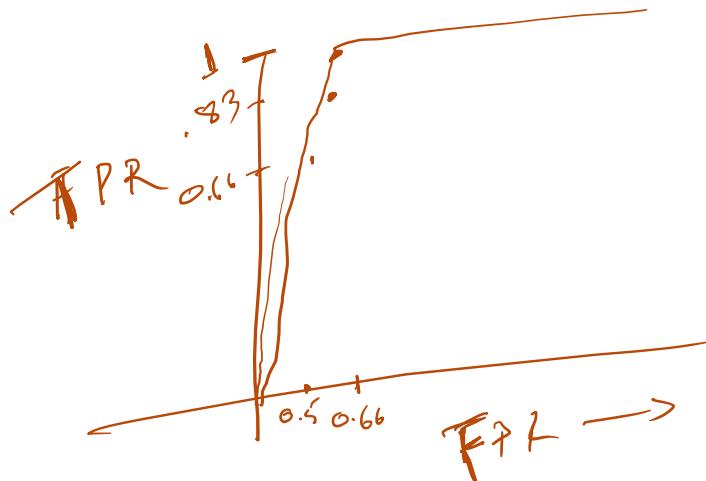
$$\hat{y} > 0.3 \quad TN = 2 \quad FPR = \frac{4}{4+2} = \frac{4}{6} = 0.66$$

$$FP = 4$$

$$FN = 1$$

$$\hat{y} > 0.9 \quad TPR = \frac{TP}{TP + FN} = \frac{4}{4+2} = \frac{4}{6} = 0.66$$

$$\begin{array}{l} \text{TP} \rightarrow 4 \\ \text{TN} \rightarrow 3 \\ \text{FP} \rightarrow 3 \\ \text{FN} \rightarrow 2 \end{array} \quad TPR = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad FPR = \frac{\text{FP}}{\text{FP} + \text{TN}} = \frac{3}{3+3} = \frac{1}{2} = 0.5$$



#### Task 4: ROC Curve Table

Generate a table of **TPR** and **FPR** at different thresholds (e.g., 0.9, 0.8, ..., 0.1), and create a rough sketch of the ROC curve (or compute AUC numerically if you prefer).

#### Task 5: Compute AUC-ROC

Using the TPR and FPR values from Task 4:

- Plot the **ROC curve** (TPR vs FPR).
- Use the **trapezoidal rule** to compute AUC:

$$(0, 0.1) \quad \text{AUC} = \int_0^1 TPR(FPR) dFPR \approx \sum \left( \frac{TPR_i + TPR_{i+1}}{2} \cdot (FPR_{i+1} - FPR_i) \right)$$

$$(0.1, 0.2)$$

$$0.1 \rightarrow TPR = \frac{1}{0.66}$$

$$0.2 \rightarrow TPR = \frac{0.83}{0.66}$$

$$\frac{1 + 0.83}{2} (0.66 - 0.6)$$

$$(0.2, 0.3)$$

$0.9 \rightarrow 1$