

## 1: Introduction to Model Evaluation

- 90% of ML models never make it to production
- Failure due to misapplication, not ML tech
- Causes: non-ML-ready data, wrong problem selection, poor ops collaboration
- Importance of model evaluation in addressing these issues

## 2: Why Model Evaluation is Critical

- Verifies effectiveness and reliability
- Prevents real-time application failures
- Identifies and mitigates data leakage

**3 .Data leakage** occurs when **information from outside the training dataset** — especially from the target variable (label) — is **accidentally included in the training process**.

This leads to the model learning patterns it **shouldn't have access to** at prediction time, resulting in:

- **Overly optimistic performance during training/testing**
- **Poor generalization to real-world, unseen data**
  -
- Enables continuous performance monitoring

## 4: Types of Predictive Models

- Regression: Continuous output
- Classification: Nominal or binary output

## 5: Confusion Matrix Overview

- 2x2 matrix for binary classification
- Measures precision, recall, accuracy, specificity, AUC-ROC

## 6: Confusion Matrix Terms

- TP: Correctly predicted positive
- TN: Correctly predicted negative
- FP (Type I Error): False alarm
- FN (Type II Error): Missed detection
- Metrics: Accuracy, Precision, NPV, Recall, Specificity

### Accuracy

#### Definition:

The ratio of correctly predicted observations to the total observations.

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

#### **When to Use:**

Good when **class distribution is balanced**.

#### **Example:**

- Model predicted 90 out of 100 patients correctly → Accuracy = 90%

## **Precision**

#### **Definition:**

Out of all positive predictions, how many were **actually positive**?

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

#### **When to Use:**

Important when **false positives** are costly (e.g., predicting cancer when there isn't one).

#### **Example:**

- Model predicted 10 people as having disease, but only 6 actually had it → Precision =  $6/10 = 60\%$

## **Recall (Sensitivity or True Positive Rate)**

#### **Definition:**

Out of all actual positives, how many did the model **correctly identify**?

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

#### **When to Use:**

Important when **missing positive cases is dangerous** (e.g., fraud detection, disease diagnosis).

#### **Example:**

- There were 20 sick patients, and model found 15 → Recall = 15/20 = 75%

## Specificity (True Negative Rate)

### Definition:

Out of all actual negatives, how many were **correctly classified**?

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

### When to Use:

Important when **false alarms are costly**, like falsely flagging a person as a criminal.

### Example:

- 80 healthy people; model wrongly flagged 10 as sick → Specificity = 70/80 = 87.5%

## 7: F1-Score

- Harmonic mean of Precision and Recall
- Useful for imbalanced datasets

### F1-Score?

The **F1-Score** is the **harmonic mean** of **Precision** and **Recall**. It provides a **single score** that balances both **false positives** and **false negatives**.

---

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


---

### Why Harmonic Mean (and not average)?

- The **harmonic mean** is more conservative than the arithmetic mean.
- It punishes extreme values (i.e., if either Precision or Recall is low, F1 is low).

- This makes F1-Score **especially useful when you need a balance** between Precision and Recall.
- 

## 📌 When to Use F1-Score?

✓ Best suited for **imbalanced datasets** — when:

- One class significantly outweighs the other.
  - Accuracy may be misleading (e.g., predicting all as negative might give high accuracy but is useless).
- 

## 💡 Example Scenario – Medical Diagnosis:

Let's say out of 100 patients:

- 10 actually have a disease (positive class)
- Model identifies 6 as positive
  - Out of these, 5 are actually sick (True Positives)
  - 1 is healthy (False Positive)

⌚ Precision:

$$= \frac{5}{5 + 1} = 0.83$$

⌚ Recall:

$$= \frac{5}{10} = 0.5$$

⌚ F1-Score:

$$= 2 \times \frac{0.83 \times 0.5}{0.83 + 0.5} \approx 0.625$$

🔍 Even though **precision is high, recall is low**, and the **F1-Score reflects that**.

---

## Summary:

Metric	Focus	Use When...
Precision	False Positives are costly	You want fewer <b>false alarms</b>
Recall	False Negatives are costly	You want to catch <b>all positives</b>
<b>F1-Score</b>	Balance between P & R	You need a balance; classes are <b>imbalanced</b>

## 8: Gain & Lift Charts

- Rank ordering of probabilities
- Steps:
  1. Calculate probabilities
  2. Rank and decile
  3. Compute response rate
  4. Plot Lift vs. Population
- Shows model effectiveness in distinguishing responders

### Gain & Lift Charts?

Gain and Lift Charts are **visual tools** used to evaluate how well your **classification model ranks the likelihood** of a certain outcome — for example, whether a customer will respond to a campaign or default on a loan.

They help you answer:

"Is my model good at identifying the most likely responders?"

---

### Key Concepts:

- **Gain** tells you how much more likely you're to capture a positive outcome compared to random selection.
- **Lift** is the ratio of the gain achieved using the model vs. random targeting.

---

### Steps to Build a Gain/Lift Chart:

#### Step 1: Predict Probabilities

Use your classification model (e.g., Logistic Regression, Random Forest) to **predict the probability** of the positive class for each data point.

Example:

Say you have 1,000 customers, and your model predicts their probability of responding to a campaign.

---

## Step 2: Sort and Decile the Data

Sort the customers by **predicted probability in descending order**, and divide them into **10 equal-sized groups (deciles)**.

- **Decile 1:** Top 10% of customers with highest predicted response probability
  - ...
  - **Decile 10:** Bottom 10% of customers
- 

## Step 3: Calculate Response Rate

For each decile, calculate:

- Number of **actual responders**
- % of total responders captured

Let's say:

- Total responders = 200 out of 1000 customers
  - Decile 1 contains 70 responders
  - So, Gain at Decile 1 =  $(70 / 200) \times 100 = 35\%$
- 

## Step 4: Compute Lift

Lift=Response rate in decile/Average response rate  
$$\text{Lift} = \frac{\text{Response rate in decile}}{\text{Average response rate}}$$

Example:

- Avg response rate =  $200 / 1000 = 20\%$
- Response rate in Decile 1 =  $70 / 100 = 70\%$

$$\text{Lift} = \frac{\text{Response rate in decile}}{\text{Average response rate}}$$

💡 Meaning: Customers in Decile 1 are **3.5x more likely** to respond than if chosen randomly.

---

## Step 5: Plot the Charts

- **Gain Chart:** % Responders (Y-axis) vs. % Population (X-axis)
  - **Lift Chart:** Lift value (Y-axis) vs. Decile Number (X-axis)
- 

📈 Interpreting the Charts:

### ✓ A Good Model:

- Gain Curve rises steeply initially
- Lift is highest in early deciles (e.g., Decile 1-3)
- Shows that the model is effectively **ranking responders ahead of non-responders**

### ✗ A Poor Model:

- Gain Curve is close to diagonal (random)
- Lift is close to 1 in all deciles

## 9: Kolmogorov-Smirnov Chart (K-S Chart)

- Measures separation between positive and negative distributions
- Max difference indicates effectiveness

In most classification models the K-S will fall between 0 and 100, and that the higher the value the better it is at separating the positive from negative cases.

The chart typically includes:

- **Cumulative % of Positives** (Y-axis)
- **Cumulative % of Negatives**
- **X-axis:** Sorted population (usually in percentiles or deciles)
- The **maximum vertical distance** between these two curves is the **K-S statistic**

## Interpreting the K-S Statistic

- The **higher** the K-S value, the **better** the model is at distinguishing classes.
- Typical range:
  - 0: **No separation** (bad model)
  - 1 (or 100%): **Perfect separation** (ideal but rare)
  - In practice:
    - K-S > 0.3 is considered acceptable
    - K-S > 0.4 is strong

## 10: AUC-ROC

- ROC: Sensitivity vs. (1-Specificity)

### ROC?

**ROC** stands for **Receiver Operating Characteristic** curve.

It's a **graphical plot** that shows the **performance of a binary classification model** across all classification thresholds.



### What does ROC plot?

It plots:

- Y-axis: True Positive Rate (TPR) = Sensitivity = Recall

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

- X-axis: False Positive Rate (FPR) = 1 – Specificity

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

Each point on the ROC curve represents a different threshold used to convert predicted probabilities

### AUC?

**AUC** stands for **Area Under the Curve (ROC Curve)**.

- It's a **single number summary** of the ROC curve.
- AUC value ranges from **0 to 1**:
  - **1.0** = Perfect model
  - **0.5** = No better than random guessing

- < 0.5 = Worse than random (bad model)
- AUC: Area under the ROC curve
- Good for class imbalance; threshold independent

**AUC** stands for **Area Under the Curve (ROC Curve)**.

- It's a **single number summary** of the ROC curve.
- AUC value ranges from **0 to 1**:
  - **1.0** = Perfect model
  - **0.5** = No better than random guessing
  - < 0.5 = Worse than random (bad model)

Metric	Description
ROC	Curve of TPR vs. FPR for various thresholds
AUC	Numerical area under ROC curve (0.5–1.0)
High AUC	Better model discrimination ability
Advantage	Works well with class imbalance, threshold-free

## 11: Log Loss

- Negative log of correct predicted probabilities
- Evaluates probabilistic model performance
- Penalizes false predictions more heavily

### Log Loss?

**Log Loss** measures the performance of a **classification model** where the predicted output is a **probability value between 0 and 1**.

It **penalizes incorrect predictions** — especially if the model is confident (i.e., assigns a probability close to 1) but wrong.

$$\text{Log Loss} = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \log(p_i) + (1 - y_i) \cdot \log(1 - p_i)]$$

Where:

- $N$ : total number of observations
- $y_i$ : actual label (0 or 1)
- $p_i$ : predicted probability of class 1 (positive class)

## 12: Gini Coefficient

- Derived from AUC:  $\text{Gini} = 2 * \text{AUC} - 1$
- Measures discriminatory power of model
- The **Gini Coefficient** is a metric used to evaluate **binary classification models**, especially in fields like credit scoring, marketing, and insurance.
- It is derived from the **AUC (Area Under the ROC Curve)** and measures how well the model can **differentiate between classes** (typically, positive vs. negative outcomes).

$$\text{Gini} = 2 \times \text{AUC} - 1$$

- If  $\text{AUC} = 1.0$  (perfect model), then  $\text{Gini} = 1$  (perfect separation).
- If  $\text{AUC} = 0.5$  (random guessing), then  $\text{Gini} = 0$ .
- If  $\text{AUC} < 0.5$ , then  $\text{Gini} < 0$  (worse than random model).

## 13: RMSE (Root Mean Squared Error)

- Popular in regression tasks
- Sensitive to outliers
- Penalizes large errors more heavily
- Formula:  $\sqrt{\text{sum}(\text{errors}^2) / N}$
- **Root Mean Squared Error (RMSE)** is one of the most commonly used evaluation metrics in **regression tasks**. It measures the **average magnitude of the prediction errors**, giving a relatively high weight to **large errors** due to the squaring part.
- It helps us understand **how far off our model's predictions are from the actual values**, on average.

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

Where:

- $N$  = number of observations
- $\hat{y}_i$  = predicted value
- $y_i$  = actual value

## 14: R-squared / Adjusted R-squared

### R-Squared Value Range

$$0 \leq R^2 \leq 1$$

- $R^2 = 1 \rightarrow$  Perfect fit: The model explains **100%** of the variability
- $R^2 = 0 \rightarrow$  No fit: The model explains **none** of the variability

Higher the  $R^2$ , better the model (but not always! Overfitting can happen.)

## Why Cross-Validation?

- When you train a model and test it on the **same data**, you may get **overly optimistic results**.
- Cross-validation allows you to **test the model on unseen data**, ensuring the model generalizes well.

### Benefits:

- **No data leakage** from training into evaluation
- **In-time validation** simulates how the model will perform in real scenarios
- Ensures model isn't just memorizing training data

## Drawbacks of Simple Validation

### What is Simple Validation?

- Splitting dataset once into **train/test sets** (e.g., 70/30 split)

## Issues:

1. **Reduces training data:**
    - o Less data to train = possibly weaker model
  2. **Leads to biased models:**
    - o The test set might not represent the full data distribution (random noise)
    - o You might get different results with different splits
  3. **Unstable metrics:**
    - o Model accuracy may vary too much with each split
- 

## K-Fold Cross-Validation

### What is it?

- Split data into **k equal parts** (folds)
- Train the model on **k-1 folds**, test on the **remaining 1 fold**
- **Repeat k times**, changing the test fold each time

### Result:

- Final model performance is the **average across all k runs**

### Advantages:

- **Every data point is used for both training and testing**
- **Reduces variance** in performance estimates
- **Minimizes selection bias** compared to single split

Example ( $k = 5$ ):

- Run 1: Train on folds 2–5, Test on fold 1
  - Run 2: Train on folds 1,3–5, Test on fold 2
  - ...and so on
- 

## How K-Fold Helps

### Key Insights:

- **Robustness:** Reliable performance across all folds implies strong model
  - **Overfitting detection:** If accuracy is high on train folds but low on test folds, overfitting is likely
  - **Consistency:** Low variation in metrics across folds = better generalization
-

## Best Practices in Model Evaluation

Step	Practice
 Data Splitting	Use <b>train</b> , <b>validation</b> , and <b>test</b> sets properly. Example: 60/20/20
 Metric Selection	Choose based on problem type: e.g., Accuracy, F1-Score (classification); RMSE, R <sup>2</sup> (regression)
 Monitor Overfitting	Use validation loss vs. training loss curves
 Prevent Data Leakage	Never let test data influence training (e.g., feature scaling must be done <b>after split</b> )
 Continuous Monitoring	Keep track of <b>model drift</b> in production and re-train if performance drops