

# Basics of Machine Learning

1. Supervised ML Concepts
2. Unsupervised ML concepts
3. Overfitting, Bias, Validation, Metrics
4. Basic Data Handling skill

## Supervised Machine Learning (ML) Concepts:

Supervised ML involves training a model on **labeled data** (input-output pairs).

### Key Concepts:

-  Input (features) + Output (labels)
-  Model learns from training data to predict labels on unseen data
-  Examples: Classification (Spam/Not Spam), Regression (Price Prediction)

### Definition:

Supervised learning is a type of machine learning where an algorithm learns from labeled data. "Labeled" data means that each data point is associated with a corresponding output or "target" value.

The goal is to learn a mapping function that can predict the output for new, unseen input data.

### Key Concepts:

- Labeled Data: The cornerstone of supervised learning.
- Training Data: The dataset used to train the model.
- Features: The input variables used to predict the output.
- Labels/Targets: The output variables that the model is trying to predict.

### Algorithms:

Common supervised learning algorithms include:

- Linear Regression (for predicting continuous values)
- Logistic Regression (for classification)
- Decision Trees
- Random Forests
- Support Vector Machines (SVMs)
- Neural Networks

### Applications:

- Image classification
- Predicting house prices
- Medical diagnosis

- Spam detection

## **Unsupervised Machine Learning (ML):**

### **Definition:**

Unsupervised learning is a type of machine learning where an algorithm learns from unlabeled data.

The goal is to discover hidden patterns or structures in the data.

### **Key Concepts:**

Unlabeled Data: The input data without corresponding output labels.

Clustering: Grouping similar data points.

Dimensionality Reduction: Reducing the number of variables in a dataset.

### **Algorithms:**

Common unsupervised learning algorithms include:

- K-means clustering
- Hierarchical clustering
- Principal Component Analysis 1 (PCA)
- Autoencoders

### **Applications:**

- Customer segmentation
- Anomaly detection
- Image compression
- Recommendation systems

## **Validation**

Validation in **Machine Learning (ML)** refers to evaluating a model's performance to ensure that it generalizes well to unseen data.

### **1. Train-Test Split**

A simple method where the dataset is split into two subsets:

- **Training set:** Used to train the model.
- **Test set:** Used to evaluate the model's performance on unseen data.

## **# Stratified K-Fold Validation**

Stratified K-Fold validation is a cross-validation technique that preserves the class distribution in each fold, making it particularly useful for imbalanced datasets.

### **## How It Works**

1. **\*\*Splitting Process\*\*:** The dataset is divided into K folds (subsets) while maintaining the same percentage of samples for each class as in the original

dataset.

2. \*\*Iteration\*\*: The process runs K times, with each fold as the test set exactly once while the remaining K-1 folds form the training set.
3. \*\*Evaluation\*\*: The model is trained and evaluated K times, with the final performance typically being the average of all K evaluations.

## ## Key Advantages

- \*\*Preserves Class Distribution\*\*: Especially important for imbalanced datasets
- \*\*Reduces Variance\*\*: Provides more reliable performance estimates
- \*\*Comprehensive Evaluation\*\*: Every data point gets to be in the test set exactly once

## ## Implementation in Python

```
```python
from sklearn.model_selection import StratifiedKFold
from sklearn.datasets import load_iris
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

# Load dataset
X, y = load_iris(return_X_y=True)

# Initialize Stratified K-Fold
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

# Initialize model
model = LogisticRegression()

# Perform cross-validation
accuracies = []
for train_index, test_index in skf.split(X, y):
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]

    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    accuracies.append(accuracy_score(y_test, y_pred))

print(f"Average accuracy: {sum(accuracies)/len(accuracies):.4f}")
```

```

## **## When to Use**

- With classification problems (especially with imbalanced classes)
- When you need reliable performance estimates
- When dataset size is limited

## **## Comparison with Regular K-Fold**

- Regular K-Fold doesn't maintain class proportions
- Stratified K-Fold is generally preferred for classification
- Regular K-Fold may be sufficient for regression problems or balanced datasets

## **#### \*\*Intuitive Explanation of Stratified K-Fold Validation\*\***

Imagine you're a teacher dividing your class into **5 study groups** for a series of quizzes. Your class has:

- **60% "A" students** (high performers)
- **30% "B" students** (average performers)
- **10% "C" students** (struggling)

If you split them randomly into groups, some groups might accidentally get **too many "A" students**, while others get **too many "C" students**. This would make some quizzes unfairly easy or hard, and your evaluation of the groups' performance wouldn't be reliable.

## **#### \*\*What Stratified K-Fold Does:\*\***

Instead of random splitting, **Stratified K-Fold ensures that each group has the same ratio of A, B, and C students as the whole class** (60% A, 30% B, 10% C).

### **This way:**

- Every quiz (test fold) is **representative** of the whole class.
- No group gets an unfair advantage or disadvantage.
- The average performance across all quizzes gives a **true estimate** of how well the teacher's methods work.

---

## **### \*\*Real-World Machine Learning Example\*\***

Suppose you have:

- A dataset of **1000 patients**, where:
  - **900 are healthy** (class 0)

- \*\*100 have a disease\*\* (class 1)

If you use \*\*normal K-Fold\*\*, some folds might get \*\*too few (or zero) disease cases\*\*, making model evaluation unreliable.

### **But \*\*Stratified K-Fold\*\* ensures:**

- Each fold has \*\*90% healthy & 10% disease cases\*\*, just like the original data.
- The model is fairly tested on all types of data.

---

### **### \*\*Key Takeaway\*\***

Stratified K-Fold is like a \*\*fair referee\*\* ensuring every test round has the right mix of classes, so your model's performance isn't skewed by luck.

#### **\*\*Use it when:\*\***

- ✓ You have \*\*imbalanced classes\*\* (e.g., fraud detection, rare diseases).
- ✓ You want \*\*reliable performance estimates\*\* without randomness affecting results.

## **### \*\*Practical Example of Stratified K-Fold Validation\*\***

Let's walk through a real-world scenario step by step using Python and a sample dataset.

---

### **### \*\*Scenario: Predicting Loan Defaults (Imbalanced Data)\*\***

We'll use a dataset where:

- \*\*Most people repay loans\*\* (class 0 = 90%)
- \*\*A few default\*\* (class 1 = 10%)

### **#### \*\*Step 1: Load and Inspect the Data\*\***

```
```python
import numpy as np
from sklearn.datasets import make_classification

# Create an imbalanced dataset (90% non-default, 10% default)
X, y = make_classification(
    n_samples=1000,
    n_classes=2,
    weights=[0.9, 0.1], # 90% class 0, 10% class 1
    random_state=42
)

# Check class distribution
```

```
print("Class distribution:", np.bincount(y))
```
**Output:**  

```
Class distribution: [900 100] # 900 non-defaults, 100 defaults
```
---
```

#### #### \*\*Step 2: Apply Stratified K-Fold\*\*

We'll split the data into \*\*5 folds\*\* while preserving the 90:10 ratio in each fold.

```
```python
from sklearn.model_selection import StratifiedKFold

skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

for fold, (train_idx, test_idx) in enumerate(skf.split(X, y)):
    y_train, y_test = y[train_idx], y[test_idx]
    print(f"Fold {fold + 1}:")
    print(f" Train class distribution: {np.bincount(y_train)}")
    print(f" Test class distribution: {np.bincount(y_test)}")
```
**Output:**  

```

```

Fold 1:

```
Train class distribution: [720 80] # 90% class 0, 10% class 1
Test class distribution: [180 20] # Same ratio!
```

Fold 2:

```
Train class distribution: [720 80]
Test class distribution: [180 20]
... (all folds maintain the 90:10 ratio)
```
---
```

#### \*\*Key Observation:\*\*

- Each fold's test set has \*\*exactly 10% defaults\*\*, just like the original data.
- Without stratification, some folds might randomly get \*\*0 defaults\*\* or \*\*too many\*\*, leading to biased evaluation.

#### #### \*\*Step 3: Train a Model with Stratified K-Fold\*\*

Let's use Logistic Regression and evaluate accuracy across folds.

```
```python
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

```
model = LogisticRegression()
accuracies = []

for train_idx, test_idx in skf.split(X, y):
    X_train, X_test = X[train_idx], X[test_idx]
    y_train, y_test = y[train_idx], y[test_idx]

    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    accuracies.append(acc)
    print(f"Accuracy: {acc:.2f}")

print(f"\nAverage accuracy: {np.mean(accuracies):.2f}")
```
**Output:**

```
Accuracy: 0.92
Accuracy: 0.91
Accuracy: 0.93
Accuracy: 0.90
Accuracy: 0.92

Average accuracy: 0.92
```
---
```

### ### \*\*Comparison: Stratified vs. Random K-Fold\*\* #### \*\*Problem with Random K-Fold (No Stratification)\*\*

```
```python
from sklearn.model_selection import KFold

kf = KFold(n_splits=5, shuffle=True, random_state=42)

for fold, (train_idx, test_idx) in enumerate(kf.split(X)):
    y_test = y[test_idx]
    print(f"Fold {fold + 1} test classes: {np.bincount(y_test)}")
```

```

### \*\*Possible Output (Unreliable Splits):\*\*

```
```
Fold 1 test classes: [183 17] # 8.5% defaults (close, but not exact)
Fold 2 test classes: [177 23] # 11.5% defaults (oversampled)
Fold 3 test classes: [195 5] # 2.5% defaults (too few!)
```

```

### **\*\*Issue:\*\***

- Fold 3 has \*\*only 5 defaults\*\* (too few for reliable evaluation).
- Model performance would vary wildly across folds.

---

### **### \*\*When Should You Use Stratified K-Fold?\*\***

1. \*\*Classification Problems\*\* (especially imbalanced datasets).
2. \*\*Small Datasets\*\* (where random splits might miss minority classes).
3. \*\*Reliable Evaluation\*\* (when you need consistent class ratios in each fold).

### **\*\*Avoid it for:\*\***

- Regression problems (no classes to stratify).
- Perfectly balanced datasets (normal K-Fold is fine).

---

### **### \*\*Final Summary\*\***

- \*\*Stratified K-Fold = Fair Splitting\*\* (maintains class ratios).
- \*\*Normal K-Fold = Random Splitting\*\* (risks skewed evaluation).
- \*\*Use Case:\*\* Fraud detection, medical diagnosis, or any imbalanced classification.

Stratified K-Fold is indeed a powerful tool, especially when dealing with real-world "messy" data where classes are imbalanced. Let me add a few more intuitive insights and a bonus practical scenario to really cement the concept.

---

### **### \*\*Bonus: The "Hospital Lab Test" Analogy\*\***

#### **Imagine you're a researcher testing a new disease-detection AI:**

- Your hospital has \*\*1,000 patient blood samples\*\*:
  - \*\*950 healthy\*\* (class 0)
  - \*\*50 with a rare disease\*\* (class 1)

If you split the data randomly into 5 folds (like normal K-Fold), you might accidentally put:

- \*\*All 50 disease cases into just 1-2 folds\*\* (other folds get 0 disease cases!)
- Your AI would seem "perfect" in folds with no disease cases (predicting "healthy" every time), but this is dangerously misleading.

#### **\*\*Stratified K-Fold acts like a careful lab technician:\*\***

- It ensures every fold gets exactly \*\*10 disease cases\*\* (5% of each fold, matching the original 5% disease rate).
- Now when your AI scores 95% accuracy, you know it's truly because it works

—not because it lucked out on easy folds.

---

### ### \*\*Advanced Practical Tip: Stratification with Multiple Classes\*\*

What if you have \*\*3+ classes\*\* (e.g., "healthy", "mild disease", "severe disease")? Stratified K-Fold handles this too!

```
```python
# Dataset with 3 imbalanced classes (70%, 20%, 10%)
X, y = make_classification(
    n_samples=1000, n_classes=3,
    weights=[0.7, 0.2, 0.1], random_state=42
)

skf = StratifiedKFold(n_splits=5)
for train_idx, test_idx in skf.split(X, y):
    print("Test fold classes:", np.bincount(y[test_idx]))
```

Output (each fold maintains 70:20:10 ratio):
```
Test fold classes: [140  40  20] # 70% / 20% / 10%
Test fold classes: [140  40  20]
...
```

```

---

### ### \*\*When Stratification Can't Help\*\*

While stratification fixes class imbalance, it \*\*won't\*\* help with:

1. **Small datasets with tiny minorities** (e.g., 1,000 samples but only 5 fraud cases → each fold gets just 1 fraud case, which is too few to learn from).
  - \*Solution:\* Oversample minority class or use synthetic data (SMOTE).

2. **Time-series data** (where order matters, and random splitting leaks future data into past training).

- \*Solution:\* Use `TimeSeriesSplit` instead.

---

### ### \*\*Key Takeaways\*\*

1. **Stratified K-Fold = "Fair Exam Proctor"**

- Ensures every test fold reflects real-world class ratios.

2. **Use it by default for classification** (unless you have a specific reason not to).

3. \*\*Watch out for:\*\*
  - Extremely rare classes (may need synthetic data).
  - Non-classification problems (e.g., regression).

## ### \*\*LOCO (Leave-One-Covariate-Out) Validation Explained Intuitively\*\*

\*\*LOCO\*\* is a model-agnostic interpretability method that tests how much a \*specific feature\* impacts your model's predictions—by systematically removing it and measuring the change in performance.

---

### ### \*\*🍎 Simple Analogy: The "Apple Pie Recipe" Test\*\*

**Imagine you're judging an apple pie contest where each pie uses:**

- Apples (A)
- Cinnamon (B)
- Sugar (C)

To determine \*\*how critical cinnamon (B) is\*\* to the pie's taste:

1. \*\*Bake a normal pie\*\* (with A+B+C) → Score: 9/10
2. \*\*Remove cinnamon (B)\*\* → New pie (A+C) → Score drops to 6/10
3. \*\*Conclusion\*\*: Cinnamon contributes \*\*+3 points\*\* to the score.

This is exactly how LOCO works for machine learning features!

---

### ### \*\*How LOCO Works (Step-by-Step)\*\*

1. \*\*Train your model\*\* on all features (e.g., `'[Age, Income, Debt]'`).
2. \*\*Remove one feature at a time\*\* (e.g., drop `Income`).
3. \*\*Retrain the model\*\* and measure performance change (e.g., accuracy drop).
4. \*\*Large performance drop?\*\* → That feature was important!

```
```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Original model (all features)
model_all = RandomForestClassifier().fit(X_train, y_train)
baseline_acc = accuracy_score(y_test, model_all.predict(X_test))

# LOCO for 'Income' feature
```

```

X_train_remove = X_train.drop(columns=['Income'])
X_test_remove = X_test.drop(columns=['Income'])

model_remove = RandomForestClassifier().fit(X_train_remove, y_train)
new_acc = accuracy_score(y_test, model_remove.predict(X_test_remove))

print(f"Baseline Accuracy: {baseline_acc:.2f}")
print(f"Accuracy without 'Income': {new_acc:.2f}")
print(f"Importance of 'Income': {baseline_acc - new_acc:.2f}")
```

```

### **\*\*Output\*\*:**

```

```
Baseline Accuracy: 0.89
Accuracy without 'Income': 0.82
Importance of 'Income': 0.07 # Big drop → Income is critical!
```

```

---

### **### \*\*Key Use Cases for LOCO\*\***

1. **Feature Importance Ranking**
  - Find which features **\*actually\*** impact predictions (vs. correlation plots).
  
2. **Debugging Bias**
  - E.g., If removing `Gender` drops accuracy, your model may be unfairly using it.
  
3. **Model Simplification**
  - Identify redundant features (no performance drop when removed).

---

### **### \*\*LOCO vs. SHAP/LIME\*\***

| Method   | Pros                      | Cons                                   |
|----------|---------------------------|--|
| **LOCO** | Simple, easy to implement | Slow (retrains model for each feature) |
| **SHAP** | Theoretical guarantees    | Computationally expensive              |
| **LIME** | Works for any model       | Local explanations only                |

### **\*\*Rule of Thumb\*\*:**

- Use **\*\*LOCO\*\*** for a quick global importance check.
- Use **\*\*SHAP/LIME\*\*** for detailed local explanations.

---

### ### \*\*When to Avoid LOCO\*\*

- \*\*High-dimensional data\*\* (1000+ features → retraining 1000 models is impractical).
- \*\*Feature interactions\*\* (LOCO may underestimate importance if features depend on each other).

\*\*Alternative\*\*: For high-dimensional data, use permutation importance (similar idea but faster).

---

### ### \*\*Try It Yourself!\*\*

Pick a Kaggle dataset (e.g., Titanic survival prediction) and:

1. Train a baseline model.
2. Apply LOCO to columns like `Age`, `Fare`, or `Sex`.
3. See which feature removal hurts accuracy the most!

## AUC-ROC (Area Under the Receiver Operating Characteristic Curve)

- **Definition:** AUC measures the ability of the model to distinguish between classes. The ROC curve is a plot of **True Positive Rate (Recall)** vs **False Positive Rate**.
- **Use:**
  - **AUC = 1:** Perfect model.
  - **AUC = 0.5:** No better than random guessing.
  - **AUC > 0.5:** Better than random, but not perfect.

## Quick Summary on Evaluation:

- **Accuracy:** Overall correctness of the model.
- **Precision:** Correctness of positive predictions.
- **Recall:** Ability to find all positive cases.
- **F1 Score:** Balance between precision and recall.
- **AUC-ROC:** Measures the model's ability to separate classes.

These metrics provide different insights into model performance, especially when the dataset is imbalanced or when certain types of errors are more costly.

## Bias vs Variance:

- **Bias:** Error due to wrong assumptions.
- **Variance:** Error due to model sensitivity to small changes in data.
-

## Correct Understanding:

| Type                | Bias      | Variance      | Description                           |
|---------------------|-----------|---------------|---------------------------------------|
| <b>Underfitting</b> | High bias | Low variance  | Model is too simple, misses patterns  |
| <b>Overfitting</b>  | Low bias  | High variance | Model is too complex, memorizes noise |

 So:

-  **High variance** → Model changes a lot with different data → **Overfitting**
-  **Low variance** → Model is stable → May still **underfit** if it has high bias

## Common “Wrong Assumptions” in High Bias Models:

| Assumption Type            | Reality                         | Impact             |
|----------------------------|---------------------------------|--------------------|
| Linear relationship        | Non-linear patterns             | Misses patterns    |
| Few features matter        | Many features influence outcome | Poor predictions   |
| Data is clean and balanced | Data is noisy or imbalanced     | Misleading results |

## Unsupervised Algorithms

| Task   | Algorithm Examples            |
|--|-------------------------------|
| <b>Clustering:</b> Group similar data points together.   | K-Means, DBSCAN, Hierarchical |
| <b>Dimensionality Reduction:</b><br>Reduce the number of features while keeping important info.                              | PCA, t-SNE, Autoencoders      |
| <b>Association Mining/Association Rule Learning:</b> Find interesting relationships (rules) between items in large datasets. | Apriori                       |

# METRICS FOR EVALUATING MACHINE LEARNING MODELS

## Regressions

- MSE
- MAE
- R Square
- Adjusted R Square

## Classifications

- ROC-AUC
- Log -Loss
- Confusion Metrics

## Unsupervised Models

- Rand Index
- Mutual Information
- Dunn's Index
- Silhouette Coefficient

## Others

- CV Error
- Heuristic methods to find K
- BLEU Score(NLP)



[www.datasciencewizards.ai](http://www.datasciencewizards.ai)



[contact@datasciencewizards.ai](mailto:contact@datasciencewizards.ai)

## ◆ 1. Classification Metrics

Used when predicting **categories/labels** (e.g., spam vs. not spam).

### 📌 Accuracy

- **Definition:** Correct predictions / Total predictions
- **Use When:** Classes are balanced.

### 📌 Precision

- **Definition:** True Positives / (True Positives + False Positives)
- **Use When:** False positives are costly (e.g., email spam filter).

### 📌 Recall (Sensitivity / True Positive Rate)

- **Definition:** True Positives / (True Positives + False Negatives)
- **Use When:** False negatives are costly (e.g., disease detection).

### 📌 F1-Score

- **Definition:** Harmonic mean of precision and recall
- **Use When:** You want a balance between precision and recall.

### 📌 Confusion Matrix

- **Definition:** A table showing TP, TN, FP, FN
- **Use:** Gives a detailed view of model performance.

### ROC Curve & AUC

- **ROC Curve:** Plots TPR vs. FPR
- **AUC:** Area under the ROC curve; closer to 1 is better.

## ◆ 2. Regression Metrics

Used when predicting **continuous values** (e.g., house price prediction).

### Mean Absolute Error (MAE)

- **Definition:** Average of absolute differences between predicted and actual values
- **Pros:** Easy to understand; less sensitive to outliers.

### Mean Squared Error (MSE)

- **Definition:** Average of squared differences between predicted and actual values

- **Pros:** Penalizes large errors more than MAE.

### Root Mean Squared Error (RMSE)

- **Definition:** Square root of MSE
- **Pros:** Same unit as target variable.

### R<sup>2</sup> Score (Coefficient of Determination)

- **Definition:** Proportion of variance explained by the model
- **Range:**  $-\infty$  to 1 (1 = perfect prediction)

## ◆ Evaluation Metrics:

**Classification:**

- **Accuracy** =  $(TP + TN) / \text{Total}$
- **Precision** =  $TP / (TP + FP)$  → How many predicted positives were correct
- **Recall** =  $TP / (TP + FN)$  → How many actual positives were identified
- **F1-score** =  $2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$

**Regression:**

- **MAE:** Mean Absolute Error
- **MSE:** Mean Squared Error
- **RMSE:** Root Mean Squared Error
- **R<sup>2</sup> score:** Percentage of variance explained by the model

## ◆ Evaluation Metrics

### For Classification Models

| Metric | Formula | Description |
|--------|---------|-------------|
|--------|---------|-------------|

|                  |   |   |
|------------------|---|---|
| <b>Accuracy</b>  | $(TP + TN) / (TP + TN + FP + FN)$   | Measures overall correctness. Use when classes are balanced.  |
| <b>Precision</b> | $TP / (TP + FP)$  | Out of all predicted positives, how many are correct. Use when <b>false positives</b> are costly (e.g., spam detection).    |
| <b>Recall</b>    | $TP / (TP + FN)$  | Out of all actual positives, how many were identified. Use when <b>false negatives</b> are costly (e.g., cancer detection). |
| <b>F1-score</b>  | $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$ | Harmonic mean of precision and recall. Best when you need balance between precision and recall.                             |

 **Example (Confusion Matrix):**

|                 | <b>Predicted Positive</b> | <b>Predicted Negative</b> |
|-----------------|---------------------------|---------------------------|
| Actual Positive | TP                        | FN                        |
| Actual Negative | FP                        | TN                        |

 **For Regression Models**

| Metric                                 | Description   |
|--|---|
| <b>MAE (Mean Absolute Error)</b>       | Average of absolute errors: `mean(`   |
| <b>MSE (Mean Squared Error)</b>        | Average of squared errors:<br>`mean((actual - predicted)^2)`<br>◆ Penalizes large errors more (good if large errors are bad).                   |
| <b>RMSE (Root Mean Squared Error)</b>  | Square root of MSE: `sqrt(MSE)`<br>◆ In the same unit as the original data.   |
| <b>R<sup>2</sup> Score (R-squared)</b> | <b>1 - (SS_res / SS_total)</b><br>◆ Explains how well the model explains variance in the data.<br>◆ Ranges from 0 to 1 (closer to 1 is better). |