## 1. Binary Classification vs. Multiclass Classification

Classification is a supervised learning task where the goal is to predict a **discrete** class label. The distinction depends on the number of possible output classes.

| Feature | Binary Classification | Multiclass Classification |
|---|---|---|
| Output Classes | Exactly **two** possible outcomes. | **Three or more** possible outcomes. |
| Example Target | Is this email Spam? (Yes/No) | Which animal is in the picture? (Cat/Dog/Bird) |
| Purpose | Distinguish between two distinct classes (e.g., positive vs. negative). | Distinguish between multiple classes simultaneously. |
| Core Models | Logistic Regression, Simple SVMs. | Softmax Regression, K-Nearest Neighbors (KNN), Random Forest. |
| Evaluation | Metrics focus on one positive class (e.g., Precision, Recall, F1 for the 'Yes' class). | Metrics often use macro or micro averaging (e.g., Macro F1-score) across all classes. |

⊞ Export to Sheets ⧉

**Example:**

- **Binary**: Predicting whether a bank transaction is **Fraudulent** or **Legitimate**.
- **Multiclass**: Predicting the type of handwritten digit from 0 to 9 (10 classes). 🔗

## 2. Explain with Example Ensemble Learning in ML

**Ensemble Learning** is a powerful machine learning paradigm where multiple individual models (often called **weak learners** or **base models**) are trained and strategically combined to solve a particular computational problem. The core idea is that a group of diverse, imperfect models can collectively produce a prediction that is more accurate and robust than any single individual model.

### Key Principle

The improved performance stems from the assumption that the errors made by different base models are generally independent and random. By aggregating the predictions, these random errors tend to cancel each other out, leading to a reduction in model variance (Bagging) or bias (Boosting).

### Example: Random Forest (A Bagging Ensemble)

1. **Multiple Models**: A Random Forest trains hundreds of individual **Decision Trees** (the weak learners) on different subsets of the data and features.
2. **Diverse Predictions**: Each tree makes a prediction for a new data point (e.g., Tree 1 predicts 'Cat', Tree 2 predicts 'Dog', Tree 3 predicts 'Cat').
3. **Aggregation**: The Random Forest aggregates these predictions using **voting** (for classification) or **averaging** (for regression). If 70% of the trees vote 'Cat', the final ensemble prediction is 'Cat'.

This combined decision process produces a highly accurate and stable prediction.

## 3. Write Short Note on Metrics for Evaluating Classifier Performance

Classifier performance metrics are essential for quantifying how well a classification model distinguishes between classes.

1. **Accuracy**: The most straightforward metric, defined as the proportion of correct predictions (True Positives + True Negatives) out of all predictions.

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

- **Limitation**: Misleading in cases of class imbalance (e.g., a 99% accurate model in a dataset with 99% true negatives is useless).

2. **Precision** (Positive Predictive Value): Of all the instances the model predicted as positive, how many were actually positive. It is a measure of **exactness**.

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

3. **Recall** (Sensitivity/True Positive Rate): Of all the instances that were actually positive, how many did the model correctly identify. It is a measure of **completeness**.

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

4. $F_1$-**Score**: The harmonic mean of Precision and Recall. It balances the trade-off between the two, providing a single score that is particularly useful when you need to avoid both False Positives and False Negatives.

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

5. **ROC AUC (Receiver Operating Characteristic - Area Under the Curve)**: Measures the model's ability to distinguish between classes at various classification thresholds. A value closer to 1 indicates better performance.

---

## 4. Explain KNN with Example

**K-Nearest Neighbors (KNN)** is a simple, non-parametric, lazy learning algorithm used for both classification and regression.

### How KNN Works

1. **No Training Phase (Lazy Learning)**: The model does not learn a discriminative function from the training data. It simply stores the entire training dataset.

2. **Prediction**: When a new, unseen data point is provided, KNN performs the following:

   - **Distance Calculation**: Calculates the distance (usually Euclidean) between the new point and every point in the training data.

   - **Neighbor Selection**: Selects the $K$ training points that are closest to the new point (the "nearest neighbors").

   - **Classification/Regression**:

      - **Classification**: Assigns the new point the class label that is most common among its $K$ neighbors (a **majority vote**).

      - **Regression**: Assigns the new point the average value of the target variable of its $K$ neighbors.

### Example

Suppose we use $K = 5$ to classify a new email as 'Spam' or 'Not Spam'.

1. The algorithm finds the 5 closest labeled emails in the training data.

2. If 4 of the 5 neighbors are labeled 'Spam' and 1 is labeled 'Not Spam', the new email is classified as **'Spam'** by majority vote.

### Key Features

- The choice of $K$ is critical; a small $K$ can be sensitive to noise (high variance), while a large $K$ can miss local patterns (high bias).

- **Non-parametric**: It makes no assumptions about the underlying data distribution.

- **Computationally Expensive at Test Time**: Since all calculations happen during prediction, it can be slow with very large datasets.

## 5. Compare Bagging and Boosting used in Ensemble Learning

**Bagging** and **Boosting** are the two primary techniques used in ensemble learning to combine multiple weak learners.

| Feature | Bagging (Bootstrap Aggregating) | Boosting |
|---|---|---|
| Goal | Reduce **Variance** (address overfitting). | Reduce **Bias** (address underfitting). |
| Training Process | **Parallel**: Base models are trained independently and simultaneously on random subsets of the data (bootstrapped samples). | **Sequential/Adaptive**: Models are trained one after the other. Each new model focuses on the errors (residuals) of the preceding models. |
| Data Sampling | Uses **bootstrapping** (sampling with replacement) to create diverse training subsets. | Uses the *entire* dataset for each subsequent learner but **weights** samples based on how previous models performed on them. Misclassified points get higher weights. |
| Model Weights | All base models (e.g., trees) are typically weighted **equally** in the final prediction. | Models are weighted based on their performance; stronger models get **higher influence** in the final prediction. |
| Examples | **Random Forest**. | **AdaBoost**, Gradient Boosting Machines (GBM), XGBoost, LightGBM. |

⊞ Export to Sheets                                                ▢

## 6. Explain Kernel Methods which are suitable for SVM

**Kernel Methods** are central to **Support Vector Machines (SVMs)**, allowing them to perform classification tasks on data that is not linearly separable in the original input space.

### The Kernel Trick

The core idea is the **Kernel Trick**, which avoids the explicit, computationally expensive calculation of mapping the data into a high-dimensional feature space. Instead, it uses a **kernel function** (or similarity function) $K(\mathbf{x}_i, \mathbf{x}_j)$ to directly calculate the dot product between the mapped points in the higher dimension, without ever performing the mapping itself:

$$\text{Original Space Dot Product} : \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\text{High-Dimensional Space Dot Product} : \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$$

Suitable Kernel Functions  🔗

By choosing different kernel functions, SVMs can model relationships of varying complexity:

1. Linear Kernel:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$$

- **Use**: When data is linearly separable, performing the same function as standard Linear SVM.

2. **Polynomial Kernel**:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i \cdot \mathbf{x}_j + r)^d$$

- **Use**: Used to model relationships where the data can be separated by a curved or polynomial decision boundary. $d$ is the degree of the polynomial.

3. **Radial Basis Function (RBF) Kernel / Gaussian Kernel**:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma ||\mathbf{x}_i - \mathbf{x}_j||^2)$$

- **Use**: The most popular choice for non-linear SVMs. It maps samples non-linearly into a much higher-dimensional space. It is highly effective when there is no prior knowledge about the data.

---

## 7. What are different techniques used for outlier handling?

**Outliers** are data points that significantly deviate from other observations. They can skew model training and affect performance. Techniques for handling them fall into two main categories: detection and treatment.

**Detection Techniques**

1. **Statistical Methods**: Use Z-score (for normally distributed data) or IQR (Interquartile Range, for non-normally distributed data) to define boundaries and identify points outside those limits.

2. **Visualization**: Box plots, scatter plots, and histograms can visually reveal points far removed from the bulk of the data.

3. **Model-Based Methods**: Use unsupervised methods like **Local Outlier Factor (LOF)** or **Isolation Forest** to score data points based on their isolation or deviation from local clusters.

Treatment Techniques  🔗

1. **Removal (Dropping)**:

   - **When**: If the outlier is due to a **data entry error** or if the dataset is large and the number of outliers is small.

   - **Caution**: Removing too many points can lead to loss of valuable information and may distort the data distribution.

2. **Imputation (Capping/Winsorizing)**:

   - **Winsorizing**: Replace the outlier value with the value of a specified percentile (e.g., set all values above the 95th percentile to the 95th percentile value). This **caps** the extreme value instead of removing the entire observation.

   - **Imputation**: Replace the outlier with a measure of central tendency (mean, median) or by predicting a value using other features.

3. **Transformation**:

   - Use mathematical transformations like **logarithmic** or **square root** to compress the value range, making extreme values less severe.

4. **Model Selection**:

   - Use models that are inherently robust to outliers, such as **Tree-based methods**

(Decision Trees, Random Forests) or models based on rank/median (e.g., **Median Absolute Deviation** instead of standard deviation).

---

## 8. Explain Random Forest Algorithm with Example

**Random Forest** is an ensemble learning method primarily used for classification and regression. It operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or the mean prediction (regression) of the individual trees.

### How it Works (The Randomness)

The "randomness" is introduced in two ways to ensure the individual trees are diverse (decorrelated):

1. **Bagging (Row Sampling)**: Each tree is trained on a different **bootstrap sample** (random sampling with replacement) of the original training data. This is **Bootstrap Aggregating (Bagging)**.

2. **Feature Randomness (Column Sampling)**: When a decision tree is being grown, at each node, only a **random subset of features** is considered for finding the best split. This prevents a single, strong feature from dominating all trees.

### Algorithm Steps (Classification)

1. Select $k$ random samples from the dataset (bootstrapping).

2. Build a decision tree for each sample.

3. For each node in the tree, select only a random subset of features to determine the best split.

4. Repeat steps 1-3 until the desired number of trees is reached.

5. **Prediction**: For a new input, all individual trees predict a class, and the Random Forest selects the class with the **majority vote**.

### Example

Imagine a Random Forest of 100 trees predicting whether a customer will click an online ad.

- The new customer's features are fed to all 100 trees.

- Tree 1 predicts **Click**, Tree 2 predicts **No Click**, ..., Tree 100 predicts **Click**.

- If 65 trees predict **Click** and 35 trees predict **No Click**, the final Random Forest prediction is **Click**.

---

## 9. Calculate Macro Average Precision, Macro Average Recall and Macro Average F-score

The confusion matrix for the 4-class classification problem (A, B, C, D) is:

| True ↓ / Pred → | A | B | C | D | Total Actual |
|---|---|---|---|---|---|
| A | 100 | 80 | 10 | 10 | 200 |
| B | 0 | 9 | 0 | 1 | 10 |
| C | 0 | 1 | 8 | 1 | 10 |
| D | 0 | 1 | 0 | 9 | 10 |
| Total Predicted | 100 | 91 | 18 | 21 | 230 |

⊞ Export to Sheets

The metrics are calculated for each class ($N = 4$ classes) and then averaged (Macro-average).

| Class | $TP$ | Precision $P = \frac{TP}{\text{Total Pred}}$ | Recall $R = \frac{TP}{\text{Total Actual}}$ | $F_1 = \frac{2PR}{P+R}$ |
|---|---|---|---|---|
| A | 100 | $\frac{100}{100} = 1.0$ | $\frac{100}{200} = 0.5$ | 0.6667 |
| B | 9 | $\frac{9}{91} \approx 0.0989$ | $\frac{9}{10} = 0.9$ | 0.1782 |
| C | 8 | $\frac{8}{18} \approx 0.4444$ | $\frac{8}{10} = 0.8$ | 0.5714 |
| D | 9 | $\frac{9}{21} \approx 0.4286$ | $\frac{9}{10} = 0.9$ | 0.5806 |

⊞ Export to Sheets                          ⟡

The calculated Macro Averages are:

| Metric | Calculation | Result (Rounded) |
|---|---|---|
| **Macro Average Precision** | $\frac{1.0+0.0989+0.4444+0.4286}{4}$ | **0.4930** |
| **Macro Average Recall** | $\frac{0.5+0.9+0.8+0.9}{4}$ | **0.7750** |
| **Macro Average F-score** | $\frac{0.6667+0.1782+0.5714+0.5806}{4}$ | **0.4992** |

⊞ Export to Sheets                          ⟡

---

## 10. Different Variants of Multi-Class Classification

Most binary classification algorithms (like SVMs or Logistic Regression) cannot natively handle more than two classes. Therefore, they must be adapted using strategies that decompose the multiclass problem into multiple binary problems.

i) One-vs-Rest (OvR) or One-vs-All (OvA)  🔗

- **Principle**: If there are $N$ classes, $N$ separate binary classifiers are trained.
- **Training**: For each class $i$, one classifier is trained to distinguish that class (the 'Positive' class) from *all other classes combined* (the 'Negative' class).
- **Prediction**: A new instance is passed to all $N$ classifiers. The instance is assigned the class corresponding to the classifier that outputs the highest confidence score.
- **Example**: For classifying digits (0-9), one classifier predicts "Is it a 0?", a second predicts "Is it a 1?", and so on, up to "Is it a 9?".

ii) One-vs-One (OvO)  🔗

- **Principle**: For $N$ classes, $\frac{N(N-1)}{2}$ binary classifiers are trained.
- **Training**: Every possible pair of classes requires a dedicated classifier. This classifier is trained only on the data points belonging to those two classes.
- **Prediction**: A new instance is passed to all $\frac{N(N-1)}{2}$ classifiers. A **majority vote** is taken, and the class that wins the most binary contests is assigned as the final label.
- **Example**: For classifying Cat, Dog, Bird (3 classes), we train three classifiers: Cat vs. Dog, Cat vs. Bird, and Dog vs. Bird. If 'Cat vs. Dog' predicts Cat, and 'Cat vs. Bird' predicts Cat, Cat receives two votes and is the final prediction.

---

## 11. Write Short Note on Adaboost

**AdaBoost** (Adaptive Boosting) is one of the earliest and most popular **Boosting** ensemble algorithms. It is highly effective and typically uses simple decision trees (often just a **decision stump**, a tree with only one split) as its weak learners.

**Key Features**

1. **Sequential Training**: AdaBoost trains a sequence of weak learners where each successive learner is trained to correct the mistakes of the previous one.

2. **Adaptive Weighting**:

   - **Sample Weighting**: Data points that were **misclassified** by the previous model are given **higher weights** so that the next weak learner focuses more on them.

   - **Learner Weighting**: Each weak learner is assigned a **weight** in the final ensemble based on its accuracy (performance). More accurate learners are given higher influence.

3. **Final Prediction**: The final prediction is a weighted sum (or weighted majority vote) of all the weak learners, where the weights are determined by the learners' individual accuracy.

**Advantage**

AdaBoost is simple, fast, and remarkably effective at reducing **bias** and converting a set of low-accuracy weak learners into a single, high-accuracy strong classifier.