Machine Learning – Complete Notes

Table of Contents

✓ Introduction	3
↔ What Machine Learning Enables	3
ට Core Cycle of ML	3
Artificial Intelligence vs Machine Learning vs Deep Learning	4
P Important Definitions	4
Ⅲ What is Data Science?	4
😘 Types of Machine Learning	4
🔾 Examples of Each ML Type	5
려 Real-World Use of ML Types	5
Ⅲ Supervised Learning	6
Ⅲ Unsupervised Learning	6
Ⅲ Semi-Supervised Learning	7
Ⅲ Reinforcement Learning	7
Applications of Machine Learning	7
🕓 Machine Learning in Everyday Life	8
🕓 Training a Machine Learning Model	8
Examples to Understand ML	8
🕓 Scikit-learn (sklearn)	8
🔾 Installing Scikit-learn and Other Libraries	8
A) Regression (General Concept)	9
1. Linear Regression	9
2. Logistic Regression	10
3. Decision Tree	10
4. Naive Bayes	10
Model Evaluation Metrics – Classification & Regression	11
1. Regression Model Metrics	11
2. Classification Model Metrics	11
Ⅲ Complete Machine Learning Lifecycle	12

1. Define the Problem	
2. Data Collection	12
3. Data Preprocessing	12
4. Choose a Model	13
5. Split the Data	13
6. Evaluate the Model	13
7. Hyperparameter Tuning	13
8. Cross Validation	
9. Finalize the Model	
10. Deploy the Model	13
11. Monitor, Update, and Maintain the Model	13
Overfitting vs Underfitting	14
Data Preprocessing	15
Feature Transformation Techniques	16
⊕ Feature Encoding	17
প্রতামি Workflow of Training & Testing:	18
Support Vector Machines (SVM)	18
♦ What is a Hyperplane?	18
≪ Key Components of SVM:	19
ত্তি Types of Kernels in SVM:	19
K-Nearest Neighbours (KNN)	20
KNN for Classification:	20
NNN for Regression:	21
▶ Distance Metrics in KNN:	21
Distance Metrics In ML (KNN)	22
Secondary Euclidean Distance	22
Manhattan Distance	23
IIII Minkowski Distance	24
Hamming Distance	25
Decision Trees	
✓ Tree Splits	
Entropy vs Gini Impurity	

Decision Tree Building Process	28
	28
Regression Trees (CART)	28
Sklearn Implementation	29
Ensemble Methods in Machine Learning	29
Types of Ensemble Methods	30
🕸 Random Forest	31
✓ Applications of Ensemble Methods	34

✓ Introduction

Who should learn ML?

Researchers, students, scientists, entrepreneurs, businesspeople — almost everyone today should learn machine learning.

What is Machine Learning (ML)?

A subfield of Artificial Intelligence (AI) that allows computers to learn from data and make decisions without being explicitly programmed.

What Machine Learning Enables

- Learning from data
- Making intelligent decisions
- Creating predictive models
- Utilizing past experiences (like humans)

Core Cycle of ML

- 1. Data Collection
- 2. Learning Patterns from Data
- 3. Making Predictions
- 4. Decision Making Based on Predictions

"Data is the new oil" - Companies/governments with data hold power and value.

Artificial Intelligence vs Machine Learning vs Deep Learning

Artificial Intelligence (AI): Engineering of building intelligent systems. ML and DL are subtypes.

Machine Learning (ML): Training machines using data to make predictions. **Deep Learning (DL):** Advanced ML using deep neural networks for solving complex problems.

Generative AI: Models that generate text, images, video (e.g., ChatGPT, DALL·E, Sora). **Artificial General Intelligence (AGI):** Hypothetical future where AI mimics all aspects of human intelligence.

Important Definitions

Machine Learning (ML): Ability of machines to learn from data and improve over time without being explicitly programmed.

Deep Learning: Advanced ML involving neural networks that mimic how the human brain learns.

AI Ecosystem:

Data Analytics -> AI -> ML -> Deep Learning -> Generative AI

What is Data Science?

A broad field combining:

- Statistics
- Mathematics
- Machine Learning
- Deep Learning
- AI

Roles: Data Scientist, ML Engineer, AI Researcher, Data Analyst, etc.

Types of Machine Learning

- 1. **Supervised Learning** Trained on labeled data (with answers).
- 2. **Unsupervised Learning** No labels; machine finds patterns.
- 3. **Semi-Supervised Learning** Mix of labeled and unlabeled data.
- 4. **Reinforcement Learning** Learns via reward/punishment (trial & error).

	Types of ML	~ ~; **
1 • Wd • Tea • Pre	ERVISED orks under supervision acher teaches ediction tcome	• No supervision • No teacher • Self learning • No labelling of data • Find patterns by itself
3 • Mix	IISUPERVISED ture of 1 and 2 ne data is labelled most is not	REINFORCEMENT Hit and Trial Learning Learn from mistakes Reward and punishment rule Prediction based on reward and punishment Depends on feedback

S Examples of Each ML Type

Supervised Learning:

- Input & expected output provided
- Algorithms: Linear Regression, Logistic Regression, Decision Trees

Unsupervised Learning:

- No labels, finds hidden patterns
- Task: Clustering (e.g., grouping fruits)

Semi-Supervised Learning:

- Few labeled + many unlabeled examples

Reinforcement Learning:

- Feedback-based learning (e.g., auto-parking, games)

Real-World Use of ML Types

- Google Search: Unsupervised learning
- ChatGPT: Reinforcement learning (thumbs up/down)
- Sora & AI agents: Generative AI based on deep learning

Supervised Learning

Supervised learning is where you provide both input data and labeled output to the machine. The machine learns from this labeled data and predicts outcomes on unseen data. Two main types of supervised learning problems:

- **Classification**: Predicts a categorical outcome (e.g., 'Yes' or 'No', types of fruit, survived/died in Titanic).
- **Regression**: Predicts a numerical value (e.g., house prices, BMI, age).

Examples:

- Binary classification: Yes/No, Male/Female
- Multiclass classification: Class of travel (1st, 2nd, 3rd)
- Regression: Predicting price of a house based on area size.

Algorithms used:

- Logistic Regression (for classification)
- Linear Regression (for regression)
- Decision Trees
- Random Forest
- KNN
- Naive Bayes
- Support Vector Machines

III Unsupervised Learning

Unsupervised learning involves input data without labeled responses. The machine identifies hidden patterns or groupings in the data.

- No teacher, no labels.
- The system tries to find structure in the data (e.g., clustering, association).

Example:

- Grouping people by shirt color in a classroom.
- Google search categorization (e.g., searching 'Babar Azam' brings cricket-related content).

Common Algorithms:

- K-Means Clustering
- Hierarchical Clustering
- DBSCAN

Semi-Supervised Learning

This is a hybrid of supervised and unsupervised learning.

- Some data is labeled, most is unlabeled.
- Helpful when labeling all data is expensive.
- The model uses both labeled and unlabeled data to train better.

Example:

- Tagging a few images manually and letting the model label the rest using pattern recognition.

Reinforcement Learning

Reinforcement learning is based on feedback (reward and punishment).

- Machine learns from actions via trial and error.
- No explicit labels; the agent learns a policy to maximize cumulative reward.

Examples:

- Autonomous car parking: Learns proper parking by receiving rewards for correct moves.
- Dog training: Fetching a ball for a treat.
- ChatGPT using thumbs up/down as feedback.

Types:

- Model-Free (e.g., Q-Learning): No environment model, learns by direct interaction.
- Model-Based: Builds a model of the environment and uses it for decision-making.

Applications of Machine Learning

- Personalized music recommendations (e.g., Arijit Singh, Naat, Tilawat, Qawali)
- News/article suggestions based on reading history
- Facebook face tagging (automated by ML)
- E-commerce recommendations (e.g., Amazon, Alibaba)
- Disease detection in humans, animals, and plants
- Cancer detection, genetic disorder identification
- Personalized nutrition through apps analyzing your blood reports
- Microbiome prediction and analysis
- Robotics, Software Engineering, Price predictions (flight tickets)
- Self-driving cars, parking systems
- Ride-hailing apps (e.g., Uber, Careem, InDrive) using geo-tagging

Machine Learning in Everyday Life

- Most apps we use today are powered by machine learning.
- AI mimics human behavior, and ML enables this by training on data.
- Without data, no machine learning. Without ML, no AI.
- Data is the new oil raw material for AI systems.

A Training a Machine Learning Model

- Machine Learning = Data-driven training
- Requires a proper algorithm to help the machine understand patterns.
- Analogy: You can't understand German if you're not trained in it.
- The model learns from the data and makes predictions.
- Final model can be saved and used to predict unknown values.

Examples to Understand ML

Example 1: Bridal Dress Price Prediction

- Inputs (X): Fabric cost, stitching cost, machine price, shop rent, discount, bonus
- Output (y): Final price of bridal dress
- Model learns these patterns to predict future prices.

Example 2: Samosa Pricing in Ramadan

- Inputs: Ingredient prices (potato, oil, mint), packaging, gas
- Output: Final price of samosa
- By tracking daily profits and sales, the model learns ideal pricing.

Scikit-learn (sklearn)

- A Python library used for traditional machine learning algorithms
- Includes algorithms for Supervised, Unsupervised, and some Reinforcement Learning
- Easy to use model training in 5-7 lines of code
- Example tasks: Classification, Regression, Clustering

A Installing Scikit-learn and Other Libraries

1. Create Conda Environment:

Conda create -n python.my python=3.10

2. Activate Environment:

Conda activate python.my

3. Install scikit-learn:

pip install scikit-learn

4. Install supporting libraries:

pip install pandas matplotlib seaborn ipykernel openpyxl

A) Regression (General Concept)

Regression is a technique to model and analyze relationships between a dependent variable and one or more independent variables.

Types:

- Linear Regression (Simple, Multiple)
- Polynomial Regression
- Ridge, Lasso Regression
- Decision Tree Regressor, Random Forest Regressor

Use Case:

- Predicting house prices, stock prices, etc.

1. Linear Regression

Linear Regression is a supervised learning algorithm used for predicting a continuous value. It finds the best-fit straight line (regression line) through the data points to predict future values.

Equation: y = mx + c

Where:

- -y = predicted value
- m = slope (how much y changes with x)
- -x = independent variable
- \mathbf{c} = intercept (value of y when x = 0)

☑ Visualizing Regression:

- Scatter plot is used to plot points.
- A regression line (best fit line) is drawn to minimize errors (distance from actual

Key Concepts:

- Dependent variable (target) is continuous
- Simple Linear Regression: one independent variable
- Multiple Linear Regression: more than one independent variable

Evaluation Metrics:

- **1. Mean Squared Error (MSE):** Average of squared errors. Lower is better.
- 2. Root Mean Squared Error (RMSE): Square root of MSE. Same unit as output.
- **3.** R² Score (Coefficient of Determination): Indicates how well data fits the model. Closer to 1 is better.

↑ When to Use Linear Regression?

- The relationship between x and y is linear.
- The errors (residuals) are normally distributed.
- No multicollinearity among input variables (for multiple regression).
- Homoscedasticity: constant variance of residuals.

2. Logistic Regression

Logistic Regression is used for binary classification problems (e.g., yes/no, spam/not spam). It models the probability that a given input belongs to a certain class using a sigmoid function.

Equation: $y = 1 / (1 + e^{-(-z)})$, where z = wx + b

Key Concepts:

- Output is between 0 and 1
- Used for classification, not regression
- Can be extended to multi-class using One-vs-Rest (OvR)

Evaluation Metrics:

- Accuracy, Precision, Recall, F1 Score
- Confusion Matrix

3. Decision Tree

Decision Tree is a supervised learning algorithm used for both classification and regression tasks.

It splits the dataset into branches based on feature values to reach a decision outcome.

Key Concepts:

- Root Node, Internal Node, Leaf Node
- Gini Impurity, Entropy, Information Gain

Advantages:

- Easy to understand and interpret
- Works with both categorical and numerical data

Disadvantages:

- Can overfit on noisy data
- Prone to creating complex trees (can be fixed with pruning)

4. Naive Bayes

Naive Bayes is a probabilistic classifier based on Bayes' Theorem with the 'naive' assumption of independence between features.

Bayes Theorem:

P(A|B) = (P(B|A) * P(A)) / P(B)

Key Concepts:

- Assumes that features are independent
- Fast and works well with high-dimensional data
- Performs well with text data (e.g., spam detection, sentiment analysis)

Types:

- Gaussian Naive Bayes
- Multinomial Naive Bayes
- Bernoulli Naive Bayes

Model Evaluation Metrics – Classification & Regression

1. Regression Model Metrics

Used when predicting numerical/continuous values (e.g., house price, temperature):

• Mean Absolute Error (MAE)

Average of absolute errors.

Formula: MAE = $(1/n) * \Sigma | y - \hat{y} |$ (*Closer to 0 is best fit*)

• Mean Squared Error (MSE)

Average of squared errors. Penalizes large errors more.

Formula: MSE = $(1/n) * \Sigma(y - \hat{y})^2$ (*Closer to 0 is best fit*)

• Root Mean Squared Error (RMSE)

Square root of MSE. Same units as the target.

Formula: RMSE = \sqrt{MSE} (*Closer to 0 is best fit*)

• R-squared (R² Score)

Proportion of variance in the dependent variable predictable from independent variables.

variables.

Formula: $R^2 = 1 - (SS_{res} / SS_{tot})$

(Closer to 1 means better fit)

2. Classification Model Metrics

Used when predicting categories (e.g., Yes/No, Spam/Not Spam):

Accuracy

Proportion of total correct predictions.

Formula: (TP + TN) / (TP + TN + FP + FN)

Precision

How many predicted positives are actually positive.

Formula: TP / (TP + FP)

(High precision = low false positives)

• Recall (Sensitivity or True Positive Rate)

How many actual positives are correctly predicted.

Formula: TP / (TP + FN)

(High recall = low false negatives)

• F1 Score

Harmonic mean of precision and recall.

Formula: 2 * (Precision * Recall) / (Precision + Recall)

(Best when you want a balance between precision and recall)

Confusion Matrix

A table layout that shows TP, TN, FP, FN. Helps understand the types of classification errors.

• ROC Curve & AUC (Area Under Curve)

Shows trade-off between True Positive Rate and False Positive Rate. AUC closer to 1 means a better model.

Complete Machine Learning Lifecycle

1. Define the Problem

Clearly define what problem you're trying to solve. Is your target variable categorical (classification) or continuous (regression)?

2. Data Collection

Collect relevant data — either primary (original) or secondary (from existing sources). Better data quality leads to better models.

3. Data Preprocessing

Clean raw data: handle missing values, outliers, incorrect formats, and perform Exploratory Data Analysis (EDA).

4. Choose a Model

Select a model based on the problem type — e.g., classification, regression, or clustering.

5. Split the Data

Split the data into training and testing sets, usually in an 80/20 ratio.

6. Evaluate the Model

Evaluate the model's performance using metrics: Accuracy, Precision, Recall, F1-score (for classification) or MAE, MSE, R² (for regression).

7. Hyperparameter Tuning

Adjust internal model parameters to improve performance. Use techniques like **GridSearchCV** or **RandomizedSearchCV**.

8. Cross Validation

Use k-fold cross-validation to test model performance on different subsets and avoid bias or overfitting.

9. Finalize the Model

Choose the best performing model and finalize it for deployment.

10. Deploy the Model

Deploy the model into production through a web app, mobile app, or software.

11. Monitor, Update, and Maintain the Model

Continuously monitor model performance, update it with new data, and retrain when necessary.

Overfitting vs Underfitting

Aspect	Overfitting	Underfitting
Definition	Model learns both patterns and noise in training data too well	Model is too simple to capture patterns in the data
Training Data Accuracy	Very high	Very low or moderate
Test Data Accuracy	Low (poor generalization to unseen data)	Low (poor on both training and test data)
Model Complexity	Too complex (many parameters, captures unnecessary patterns)	Too simple (misses important patterns)
Analogy	Tailored shirt that only fits one person	One-size-fits-all shirt that fits no one properly
Example in ML	Decision tree with no pruning, deep neural network without regularization	Linear regression applied to non-linear data
Key Reason	Model memorizes training data rather than learning generalized relationships	Model fails to learn the actual structure of the data
Fixes	Use cross-validation, regularization (L1, L2), pruning, simplify model	Use more complex models, add features, reduce bias

Data Preprocessing

ш -	dia Freprocessi	118
Step No.	Main Step	Sub-Processes & Full Explanation
1	Data Cleaning	 - Handle Missing Values: Replace with mean/median/mode, or drop rows/columns with too many missing values. - Remove Noisy Data: Apply smoothing or transformations to reduce random variability. - Handle Outliers: Detect using IQR/Z-score and remove or cap/extreme treatment. - Resolve Inconsistencies: Correct wrong data entries, typos, and ensure unit standardization (e.g., km vs miles).
2	Data Integration	 Merging Datasets: Combine multiple data sources or formats into one dataset. Duplicate Detection/Removal: Eliminate repeated or redundant rows that inflate dataset size. Consolidation: Unify records from different sources (e.g., survey + census) after resolving naming/format conflicts.
3	Data Transformation	 Scaling: Adjust feature values to a fixed range (e.g., 0-1) using MinMaxScaler or StandardAero. Normalization: Transform values to follow a Gaussian distribution (mean=0, std=1). Aggregation: Merge multiple features into one (e.g., Family Size = Siblings + Spouse + Parents/Children). Generalization: Convert specific numeric data into broader categories (e.g., Age 28 → "Adult").
4	Data Reduction	 Dimensionality Reduction: Reduce number of input features using techniques like PCA or feature importance (e.g., remove irrelevant columns). Numerosity Reduction: Replace verbose text with numeric codes (e.g., Male = 1, Female = 0). Data Compression: Store data in a more compact format for faster loading and computation.
5	Data Discretization	 Convert Continuous to Categorical: Break continuous features into bins/groups (e.g., Age 0-12 = Child, 13-19 = Teen). Encoding: Convert categorical data to numeric form (e.g., one-hot or label encoding). Clustering for Grouping: Use unsupervised algorithms to group similar records (e.g., via K-means).

Feature Transformation Techniques

Step / Technique	Description	Methods / Tools Used	When to Use
Standardization	Rescales data to have mean = 0, std = 1 (Z-score scaling)	StandardScaler	Data with normal distribution or for algorithms needing zero mean
Min-Max Scaling	Scales data between 0 and 1	MinMaxScaler	When features have known min/max or neural networks
Max-Abs Scaling	Scales features between -1 and 1 using max absolute value	MaxAbsScaler	Sparse datasets, algorithms needing zero-centered data
Robust Scaling	Scales data using median and IQR, resistant to outliers	RobustScaler	When dataset has many outliers
Normalization	Scales each data sample (row) to unit norm (length = 1)	Normalizer	Cosine similarity, KNN
Log Transformation	Applies logarithm to reduce large range/skewnes	np.log(x) or LogTransformer	For skewed, non-negative data

Step / Technique	Description	Methods / Tools Used	When to Use
	s		
Square Root Transform	Applies square root to reduce variance	np.sqrt(x)	Count-based features or high range features
Box-Cox Transformation	Makes data more Gaussian; works on positive values	PowerTransforme r(method='box- cox')	If data is not normal and only has positive values
Yeo-Johnson	Similar to Box- Cox but handles zero and negative values	PowerTransforme r(method='yeo- johnson')	When data has negatives
Quantile Transformer	Transforms feature to uniform or normal distribution	QuantileTransfor mer(output_distri bution='normal'/' uniform')	Non-parametric distributions

Feature Encoding

Encoding Type	Best For	Avoid When
Label Encoding	No order among categories, tree- based models	Linear models – might misinterpret magnitude
One-Hot Encoding	No order, few categories	Too many categories - causes dimensionality
Ordinal Encoding	Natural ordering in values	If no real order exists
Binary Encoding	High cardinality features	Very small or very large categories

Encoding Type Best For Avoid When

Frequency Quick numerical conversion When frequency doesn't reflect

Encoding feature value

♥ Workflow of Training & Testing:

- **Training Data:** Used to train the model (X_train, y_train).
- **Testing Data:** Used to evaluate the model (X_test, y_test).

Workflow:

- **1. Fit the model:** model.fit(X_train, y_train)
- 2. Predict: y_pred = model.predict(X_test)
- 3. Evaluate:
 - **For Regression:** compare y_test with y_pred using regression metrics.
 - For Classification: compare labels using classification metrics.

Support Vector Machines (SVM)

Support Vector Machines (SVM) are supervised learning algorithms used primarily for classification tasks but also applicable in regression and outlier detection.

SVM works by finding the optimal hyperplane that best separates the data points of different classes.

What is a Hyperplane?

A hyperplane is a decision boundary that separates different classes.

- In 2D, the hyperplane is simply a straight line.
- In 3D, it is a flat surface (plane).
- In higher dimensions, it is a multidimensional surface.

Real-life examples:

- A divider between salt and chili powder.

- A separator card used in grocery store drawers to divide lentils.
- A paper separator in books or the Quran.

Key Components of SVM:

- **1. Hyperplane:** The boundary line or surface that separates different classes.
- **2. Support Vectors:** The data points nearest to the hyperplane that influence its position and orientation.
- 3. Margin: The distance between the support vectors of each class and the hyperplane. SVM tries to maximize this margin.

Applications of SVM:

- Image Classification
- Text Categorization (e.g., Spam Detection)
- Protein Classification in Bioinformatics
- Handwriting Recognition

SVM Modes:

- SVC (Support Vector Classification)
- SVR (Support Vector Regression)
- Outlier Detection (limited cases)

Types of Kernels in SVM:

1. Linear Kernel:

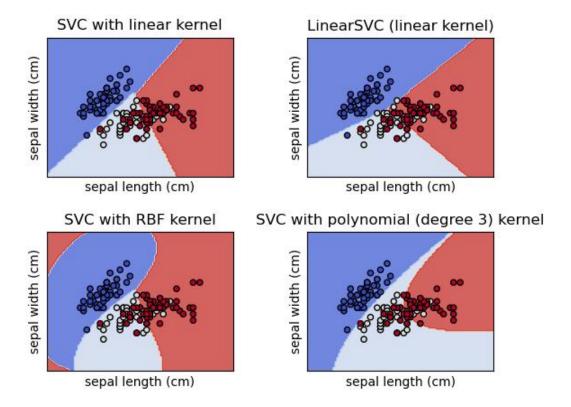
- Straight-line separation.
- Used when data is linearly separable.

2. RBF Kernel (Radial Basis Function):

- Circular or radial separation.
- Best for non-linear boundaries.

3. Polynomial Kernel:

- Curved boundaries depending on the degree (e.g., quadratic, cubic).
- Can fit more complex data patterns.



K-Nearest Neighbours (KNN)

K-Nearest Neighbours (KNN) is a supervised learning algorithm used for both classification and regression.

It works on the principle that: "A data point is known by the company it keeps."

KNN uses the 'K' closest training samples (neighbours) to make predictions.

- KNN = K Nearest Neighbours
- The value of 'K' determines how many neighbours will influence the result.

KNN for Classification:

Imagine a 2D plane with red, green, and blue points representing three different classes.

Now, a new unknown (white) point is introduced, and we want to classify it. Steps:

1. Choose a value for K (e.g., K=1, 2, or 5).

- 2. Find the K nearest points using distance metrics.
- 3. Count the number of points from each class within these neighbours.
- 4. Assign the unknown point to the class with the majority (mode).

Example:

- K=1: Closest point is green \rightarrow classify as green.
- **K=2:** Two nearest points \rightarrow majority green \rightarrow classify as green.
- **K=5:** If 3 green and 2 red \rightarrow classify as green.

KNN for Regression:

In regression, instead of using the mode, we use the mean or median of the values from K nearest neighbours.

Example:

- A new (yellow) point appears near three known points with values 2, 3, and 1.
- Mean = $(2 + 3 + 1)/3 = 2 \rightarrow assign 2$ to the new point.
- Median = $2 \rightarrow$ also can be used.

KNN can be used for numeric value prediction based on surrounding known values.

Choosing the Right K:

- K should not be too small (might lead to overfitting).
- K should not be too large (might include irrelevant neighbours).
- Must be an odd number in binary classification to avoid ties.

Distance Metrics in KNN:

- 1. Euclidean Distance Straight line distance.
- 2. Manhattan Distance Grid-like (city block) distance.
- 3. Minkowski Distance Generalized distance metric.
- 4. Hamming Distance Used for categorical data.

These are used to find how "close" one point is to another.

✓ Advantages of KNN:

- 1. Very simple and intuitive.
- 2. No training phase required.
- 3. Versatile Can be used for both classification and regression.
- 4. Multiple distance functions available.

★ Disadvantages of KNN:

- 1. Computationally expensive Must compute distance to all points.
- 2. Sensitive to imbalanced data Majority class can dominate.
- 3. Sensitive to irrelevant features Affects distance calculation.
- 4. Requires careful selection of K.

Applications of KNN:

- Recommendation systems (e.g., YouTube, e-commerce)
- Fraud detection in banking
- Healthcare classification systems
- Financial services and loan prediction
- Spam filtering

Distance Metrics In ML (KNN)



Security Euclidean Distance

Euclidean Distance is the straight-line distance between two points in a 2D or multidimensional space. It is the most common distance metric and is also known as the Pythagorean Distance because it is derived using the Pythagorean Theorem.

Consider two points: P(x1, y1) and Q(x2, y2).

Using the Pythagorean Theorem:

Distance PQ =
$$\sqrt{[(x^2 - x^1)^2 + (y^2 - y^1)^2]}$$

This formula can be extended to n-dimensional space:

$$D(P, Q) = \sqrt{[\sum (i = 1 \text{ to } n) (xi - yi)^2]}$$

Example:

If
$$P = (1, 2)$$
 and $Q = (3, 4)$, then

$$D(P, Q) = \sqrt{[(3-1)^2 + (4-2)^2]} = \sqrt{[4+4]} = \sqrt{8} \approx 2.83$$

Manhattan Distance

Manhattan Distance is the sum of the absolute differences between the coordinates of two points. It measures the distance travelled along axes at right angles (like navigating a grid of streets in Manhattan).

Formula:

$$D(P, Q) = |x2 - x1| + |y2 - y1|$$

General form for n dimensions:

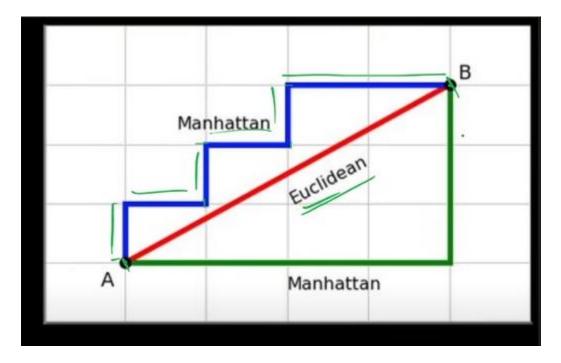
$$D(P, Q) = \sum (i = 1 \text{ to } n) |xi - yi|$$

Example:

If
$$P = (1, 2)$$
 and $Q = (3, 4)$, then $D(P, Q) = |3 - 1| + |4 - 2| = 2 + 2 = 4$

Key Note:

- Manhattan Distance is useful when movement is restricted to horizontal and vertical paths.



Minkowski Distance

Minkowski distance is a generalized metric used to measure the distance between two points in a normed vector space. It uses a parameter 'p' which defines the type of distance metric to be used. It includes **Euclidean** and **Manhattan** distances as special cases.

The formula is:

$$D(x, y) = (\Sigma |xi - yi|^p)^(1/p)$$

Where:

- x and y are two data points (vectors)
- p is a parameter:
 - If $p = 1 \rightarrow Manhattan Distance$
 - If $p = 2 \rightarrow Euclidean Distance$
 - For other values → General Minkowski Distance

Example:

Let
$$x = (1, 2)$$
, $y = (3, 4)$ and $p = 2$

$$D(x, y) = ((|1 - 3|^2 + |2 - 4|^2))^{(1/2)}$$

= (4 + 4)^(1/2)
= $\sqrt{8} \approx 2.83$

Changing p value adapts this formula to different distance metrics.

Hamming Distance

Hamming Distance measures the number of differing characters between two strings of <u>equal</u> length.

It is used in scenarios where the data is categorical or in binary/string form.

Example:

Let str1 = "AAMMAR", str2 = "AAMMAD"

Step-by-step comparison:

- A vs A \rightarrow same
- A vs A \rightarrow same
- M vs $M \rightarrow same$
- M vs $M \rightarrow same$
- A vs A \rightarrow same
- R vs D \rightarrow different

Only one character differs \rightarrow Hamming Distance = 1

Another Example:

str1 = "EUCLIDEAN", str2 = "MANHATTAN"

Length = 9

- Only 2 characters match (A and N) \rightarrow 7 differences \rightarrow Hamming Distance = 7

Used mostly in NLP, text similarity, error detection, etc.

Decision Trees

A **Decision Tree** is a supervised machine learning algorithm used for **classification** and **regression** tasks.

It splits the data into subsets based on the value of input features, forming a tree-like structure.

- **Root Node** → Represents the entire dataset.
- **Decision/Internal Nodes** → Represent tests on features.
- **Leaf/Terminal Nodes** → Represent output class labels or prediction values.

• **Branches** → Represent the outcome of a test.

✓ Advantages of Decision Trees

- Easy to understand and interpret.
- Handles both **numerical** and **categorical** data.
- Requires **less data preprocessing** (no scaling or normalization).
- Can model **non-linear relationships**.
- Works well even with **missing values**.

X Disadvantages

- Prone to Overfitting if not pruned properly.
- Can be **unstable** (small change in data can lead to different tree).
- **Biased** towards features with more levels.
- Greedy approach may not result in the **optimal** tree globally.

☑ Tree Splits

1. Entropy (Information Entropy)

Entropy is a measure of **disorder or impurity** in the dataset.

Formula:

$$\label{eq:continuity} \begin{split} &\operatorname{Entropy}(S) = -\sum_{i=1}^{i=1} \operatorname{npilog}(p_i) \operatorname{Entropy}(S) = -\sum_{i=1}^{i=1} \operatorname{npilog}(p_i) \operatorname{Entropy}(S) = -\sum_{i=1}^{i=1} \operatorname{npilog}(p_i) \operatorname{log2}(p_i) \end{split}$$

Where:

- pip_ipi = Proportion of class i in the dataset S.
- Entropy is **0** when the data is pure (only one class).
- Entropy is **maximum** when the classes are equally mixed.

2. Information Gain (IG)

Information Gain measures the **reduction in entropy** after a dataset is split on a feature.

Formula:

Where:

- AAA = Attribute being split.
- SvS_vSv = Subset for value v of attribute A.
- IG helps choose the best feature for a node.

3. Gini Impurity

An alternative to entropy used by CART (Classification and Regression Trees). It measures the **probability of a wrong classification**.

Formula:

$$Gini(S)=1-\sum i=1npi2Gini(S)=1-\setminus sum_{i=1}^{n}p_i^2Gini(S)=1-i=1\sum n-pi2$$

Where:

- pip_ipi = Probability of class i.
- Gini = $0 \rightarrow Pure node$
- Gini = Max when classes are evenly distributed

Entropy vs Gini Impurity

Metric	Entropy	Gini Impurity
Formula	$-\sum pilog_2(pi)-\sum pi log_2(pi)$	$1-\sum pi21 - \sum pi21 - \sum pi2$
Use	ID3 Algorithm	CART Algorithm
Interpretation	Information content / impurity	Likelihood of incorrect

Metric	Entropy	Gini Impurity
		label
Performance	Slightly more computational cost	Faster and easier to compute
Similarity	Both measure purity of a split and often result in similar trees	

Decision Tree Building Process

- 1. Start with the full dataset.
- 2. **Calculate entropy or gini** for the current node.
- 3. Evaluate every feature:
 - o For each feature, calculate the IG or gini gain.
 - o Choose the feature with **highest gain** to split.
- 4. **Split the dataset** based on the chosen feature.
- 5. **Repeat recursively** for each child node until:
 - All data is pure (entropy/gini = 0).
 - o Max depth or min samples criteria is reached.

Pruning Decision Trees

Pruning helps reduce **overfitting** by removing unnecessary branches.

Types:

- **Pre-pruning (Early Stopping)**: Stop tree growth based on some threshold.
- **Post-pruning**: Build the full tree first, then remove branches that don't improve accuracy.

Regression Trees (CART)

Instead of class labels, regression trees predict continuous values.

• Use **Mean Squared Error (MSE)** as the impurity measure:

```
MSE=1n\Sigma i=1n(yi-y^-)2MSE = \frac{1}{n} \sum_{i=1}^n (y_i - y_i^-)^2 MSE=n1 \quad i=1\sum_{i=1}^n (y_i - y_i^-)^2
```

% Sklearn Implementation

from sklearn.tree import DecisionTreeClassifier

model = DecisionTreeClassifier(criterion='entropy', max_depth=3)

model.fit(X_train, y_train)

Ensemble Methods in Machine Learning

Ensemble methods are machine learning techniques that combine multiple base models to create one powerful predictive model.

- Just like asking multiple experts to solve a problem and taking a consensus.
- These techniques aim to improve accuracy, robustness, and generalization.

Why Use Ensemble Methods?

- To reduce overfitting or underfitting.
- To increase accuracy and stability.
- To balance weaknesses of individual models.
- To handle data variance and noise.

Real-life Analogy

Imagine you take opinions from:

- A Logistic Regression model,
- A Decision Tree,
- And a Support Vector Machine (SVM).

If LR gives 51% accuracy, DT gives 72%, and SVM gives 68%, you can combine the strengths of each.

Whichever predictions are repeated or agreed upon are considered more reliable. This is the core idea of ensemble learning.

Basic Process of Ensemble Learning

- 1. Input Data
- 2. Train multiple models on data (with some variation)
- 3. Make predictions from all models
- 4. Aggregate predictions (by voting or averaging)

Types of Ensemble Methods

There are three main ensemble techniques:

1. A Bagging (Bootstrap Aggregating)

Definition:

Bagging involves creating multiple subsets of the data using random sampling with replacement, training a model on each subset, and aggregating their predictions.

West Points:

- Reduces variance
- Helps avoid overfitting
- All models are trained in parallel
- Random Forest is the most popular example

Analogy:

Imagine dividing 150kg of rice into bags of 10kg to transport it easily. Or think of Eid meat being distributed evenly in bags to different houses—each household cooks it differently, but the source is the same.

Random Forest

- It is a **bagging algorithm** using **decision trees**.
- Operates in parallel and aggregates results.
- Adds an additional layer of **randomness**:

At each split, selects a **random subset of features** instead of all features.

Key Components:

Component Description

Bootstrap Sample Random sampling of data with replacement

Parallel Trees Trains multiple trees at the same time

Random Features Uses random feature selection per split

Aggregation Uses **mode (majority vote)** for classification, **mean** for regression

Strengths:

- Robust and less prone to overfitting
- Can handle **non-linear** data
- Supports **feature importance**
- Efficient on large datasets

Applications:

- Credit scoring
- Stock market prediction
- Medical diagnosis
- E-commerce recommendations

Best Practices:

Parameter Tip

Number of Trees More trees → better accuracy, but more computation

Parameter Tip

Tree Depth Shallow trees reduce overfitting; deeper trees may overfit

Feature Quality Ensure selected features are relevant and informative

2. S Boosting

Definition:

Boosting is an iterative technique that turns weak learners into strong ones. Each new model is trained to fix the errors made by the previous model.

Wey Points:

- Focuses more on difficult or misclassified observations
- Models are trained sequentially
- · Weights are adjusted after each round
- Reduces both bias and variance

Types of Boosting Algorithms

1. AdaBoost (Adaptive Boosting)

- o Adjusts weights of misclassified samples.
- o Used for binary classification (e.g., face detection).

2. Gradient Boosting

- Fits new models to residual errors of previous models.
- Minimizes loss using gradient descent.
- o Variants: **XGBoost**, **LightGBM**.

3. CatBoost

- o Optimized for categorical data.
- Handles missing values automatically.

4. Stochastic Gradient Boosting

o Uses random subsets of data/features to reduce overfitting.

5. LPBoost (Linear Programming Boosting)

o Minimizes exponential loss using linear programming.

Analogy:

Imagine a football coach giving special attention to weak players during practice to turn them into strong players.

Or think of applying multiple coats of paint until the surface becomes smooth and clean.

Algorithm	What is it?	Year	Creator(s)	Example Use-Case	Pros	Cons
AdaBoost (Adaptive Boosting)	A machine learning ensemble method that adjusts the weights of the classifier and the training data to correct the previous mistakes.	1995	Yoav Freund and Robert Schapire	Used for binary classification tasks like face detection.	- Simple to implement Less prone to overfitting Good generalization.	- Sensitive to noisy data and outliers Performance depends on the weak learner.
Gradient Boosting	An approach to machine learning based on boosting which converts weak learners to strong ones, focusing on minimizing the loss function.	1999	Jerome H. Friedman	Applied in ranking algorithms, like search engines.	Highly effective in predictive accuracy. Can optimize on different loss functions. Works well with heterogeneous features.	- Prone to overfitting if not properly tuned More parameters to tune Stower to train due to sequential tree building.
Stochastic Gradient Boosting	A variation of gradient boosting which incorporates randomness in the training process to reduce variance and overfitting.	Early 2000s	Jerome H. Friedman	Used in scenarios where model variance is a concern, such as in financial modeling.	- Reduces the risk of overfitting. - Improves model generalization.	- May increase training time Can be less intuitive to tune.
LPBoost (Linear Programming Boosting)	A boosting algorithm that uses linear programming for combining weak learners.	2001	Warmuth and Liao	Useful in classification tasks with a focus on margin maximization.	Focuses on maximizing the margin. Can perform well on sparse data.	More sensitive to noise. Computational complexity carbe high.
TotalBoost	An algorithm that, like AdaBoost, focuses on increasing the margin but also adjusts for misclassified examples more aggressively.	2007	Warmuth, et al.	Applied in areas requiring robust performance despite noisy data.	Aggressively adjusts for misclassifications. Can be robust in noisy environments.	- Can be computationally intensive Risk of overemphasis on outliers.

Algorithm	What is it?	Year	Creator(s)	Example Use-Case	Pros	Cons
XGBoost (eXtreme Gradient Boosting)	An optimized distributed gradient boosting library designed to be efficient, flexible, and portable.	2016	Tianqi Chen and Carlos Guestrin	Commonly used in winning Kaggle competitions for structured or tabular data.	- Fast to train Handles missing data Supports regularization to reduce overfitting.	- More complex and requires more parameter tuning. - Can overfit if the number of trees is too large.
LightGBM	A gradient boosting framework that uses tree-based learning algorithms, designed for distributed and efficient training, especially on large datasets.	2017	Microsoft	Effective for large datasets and used in fields like genomics for feature selection.	- Faster training on large datasets Lower memory usage Supports categorical features directly.	- May overfit on small datasets. - Less community and resource support compared to XGBoost.
CatBoost (Categorical Boosting)	An open-source machine learning algorithm that can handle categorical data directly and is based on gradient boosting.		Yandex	Excelling in handling categorical data directly, like in banking for credit risk modeling.	- Good with categorical data without extensive preprocessing Less prone to overfitting User-friendly.	- Slower predictions compared to other boosting algorithms. - Relatively new, so less tested in a variety of scenarics

3. Stacking (Stacked Generalization)

Definition:

Stacking involves training multiple base models and then combining their predictions using a meta-model.

West West W

- Uses different algorithms (e.g., SVM + DT + LR)
- Meta-model learns how to combine their predictions
- Works in 2 levels:
 - Level 0: Base learners
 - Level 1: Meta-model

Analogy:

Like combining answers from three students and then a professor (meta-model) decides which combination is best.

Tree-Based Ensembles

- Many ensemble methods work with decision trees.
- When you combine multiple decision trees, you get a Random Forest (like a forest of trees).
- In boosting, multiple shallow trees are trained sequentially to form a strong model.

Applications of Ensemble Methods

- Finance: Credit scoring, fraud detection
- Healthcare: Disease prediction, diagnosis
- E-Commerce: Recommender systems
- Agriculture & Environment: Risk prediction, yield forecasting
- Social media / Platforms: Content recommendation and ranking