

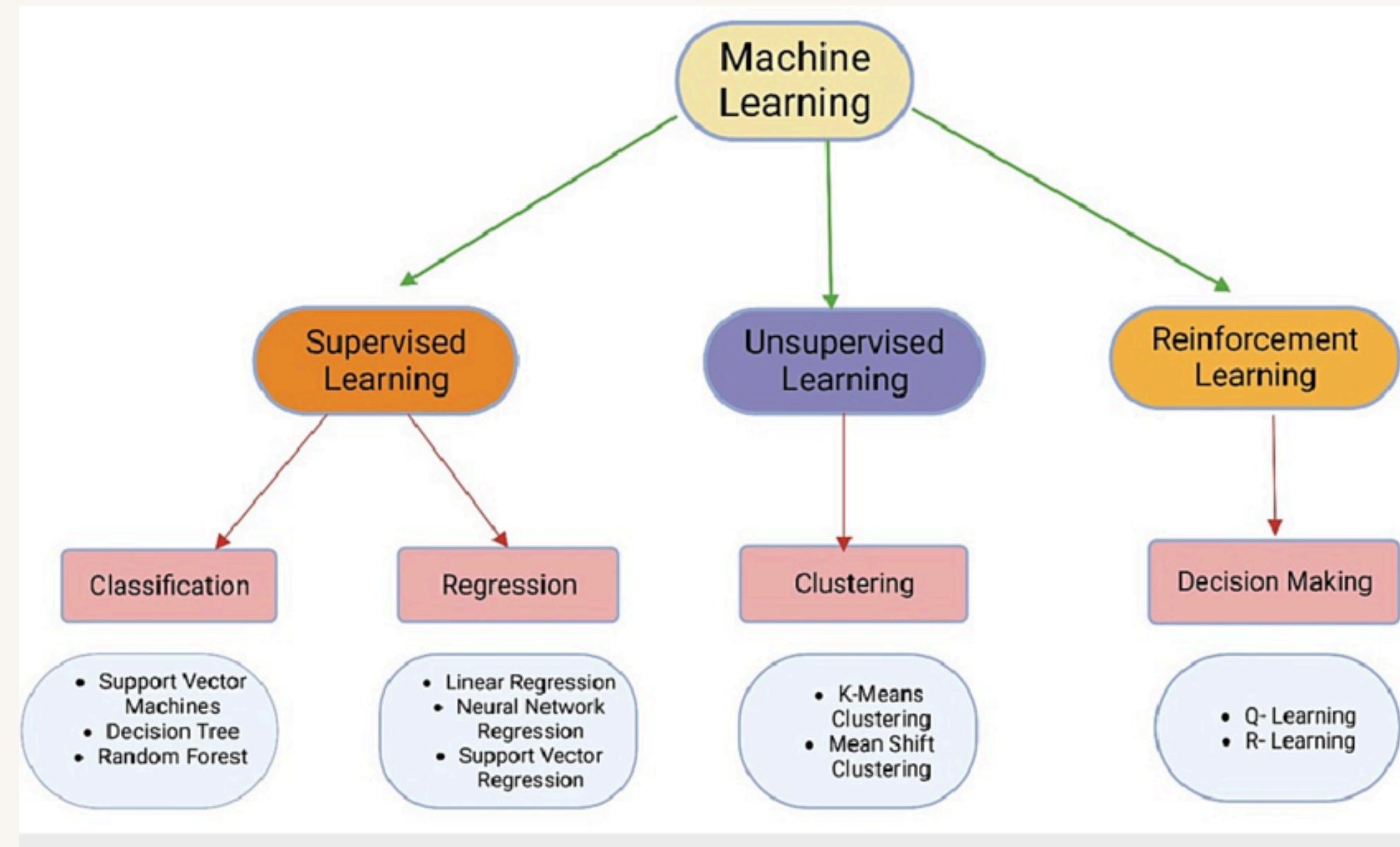
Artificial Intelligence & Machine Learning Certificate Course

Introduction to AI & ML Practical class

Sehan Arandara
Associate Software Engineer
SE Graduated @ SLIIT
AI ML Instructor

09.09.2025
Day 02

The Three Main Types of Machine Learning



A Machine Learning Pipeline

1. Data Collection
2. Data Pre-processing
3. Feature Engineering
4. Model Training
5. Model Evaluation
6. Model Deployment

Model 01: Classification

Age	Estimated_Salary	Previous_Purchases	Purchased
25	55000	2	0
35	75000	5	1
42	60000	1	0
28	80000	8	1
55	120000	10	1
22	40000	1	0
31	65000	3	0
48	95000	6	1
29	72000	4	0
38	85000	7	1
58	130000	12	1
21	25000	0	0
45	105000	9	1
33	78000	5	1
26	48000	2	0
39	92000	8	1
51	115000	11	1
24	45000	1	0
41	89000	6	1

overall goal: We want to predict if a customer will make a purchase. This is a "Yes" or "No" question.

Features : Age ,Estimated_Salary , Previous_Purchases

Target (The Output) : Purchased
1 means: Yes, the customer made a purchase.

0 means: No, the customer did not make a purchase.

Model 01: Prediction

Hours_Studied	Previous_Exam_Score	Attendance	Final_Exam_Score
5.5	65	85	72
6.2	70	91	78
4.8	58	78	65
7.5	82	95	88
3.1	45	60	52
8.9	90	98	95
2.5	38	55	45
6.8	75	92	81
9.5	95	99	98
4.1	50	70	58
5	60	82	69
7.2	80	94	85
3.8	48	65	55
8	85	96	90
2	35	50	40
6	68	88	75
5.8	66	86	73
4.5	55	75	62
7.8	84	95	89

overall goal: Predict the final exam score

Features : Hours_studied , Previous_exam_score , attendandace

Target (The Output) : final_Exam_score

Q n A

Artificial Intelligence & Machine Learning Certificate Course

Preprocessing & Model Evaluation in Machine Learning

Sehan Arandara
Associate Software Engineer
SE Graduated @ SLIIT
AI ML Instructor

16.09.2025
Day 4

ML Pipeline Recap

1. Data Collection
2. Data Preprocessing
3. Feature Engineering
4. Model Training
5. Model Evaluation
6. Model Deployment

Why Preprocessing?

- Real-world data is messy: missing values, categorical values, scaling issues.
- Ensures models can learn effectively.
- Improves the reliability of results.

Handling Missing Values

- Drop rows/columns with too many nulls.
- Fill with mean/median (numerical).
- Fill with mode (categorical).

Categorical Data

- Example: Gender, Smoker in Tabakar dataset.
 - ML models need numbers, not text.
 - Use One-Hot Encoding: Gender → [Male=1, Female=0].

Feature Scaling

- Features may have different ranges (Age vs Income).
 - Standardization (Z-score scaling).
 - Min-Max Normalization.
 - Ensures fair contribution in distance-based models.

Train-Test Split

- **Training Set:** Used for model learning.
- **Test Set:** Used for evaluation.
- **Common split:** 70%-80% train, 20%-30% test.

Model Evaluation Metrics

- **Regression:** MSE, RMSE, R^2 .
- **Classification:** Accuracy, Precision, Recall, F1-score.
- **Confusion Matrix:** TP, TN, FP, FN.

Overfitting vs Underfitting

- **Underfitting:** Model too simple, fails to capture patterns.
 - **Overfitting:** Model too complex, memorizes training data.
-
- **Solutions:** Regularization, More data, Early stopping, Cross-validation.

Saving a Model

- Trained models can be saved & reused.
 - Example: `joblib.dump(model, 'model.pkl')`.
 - Load later: `joblib.load('model.pkl')`.

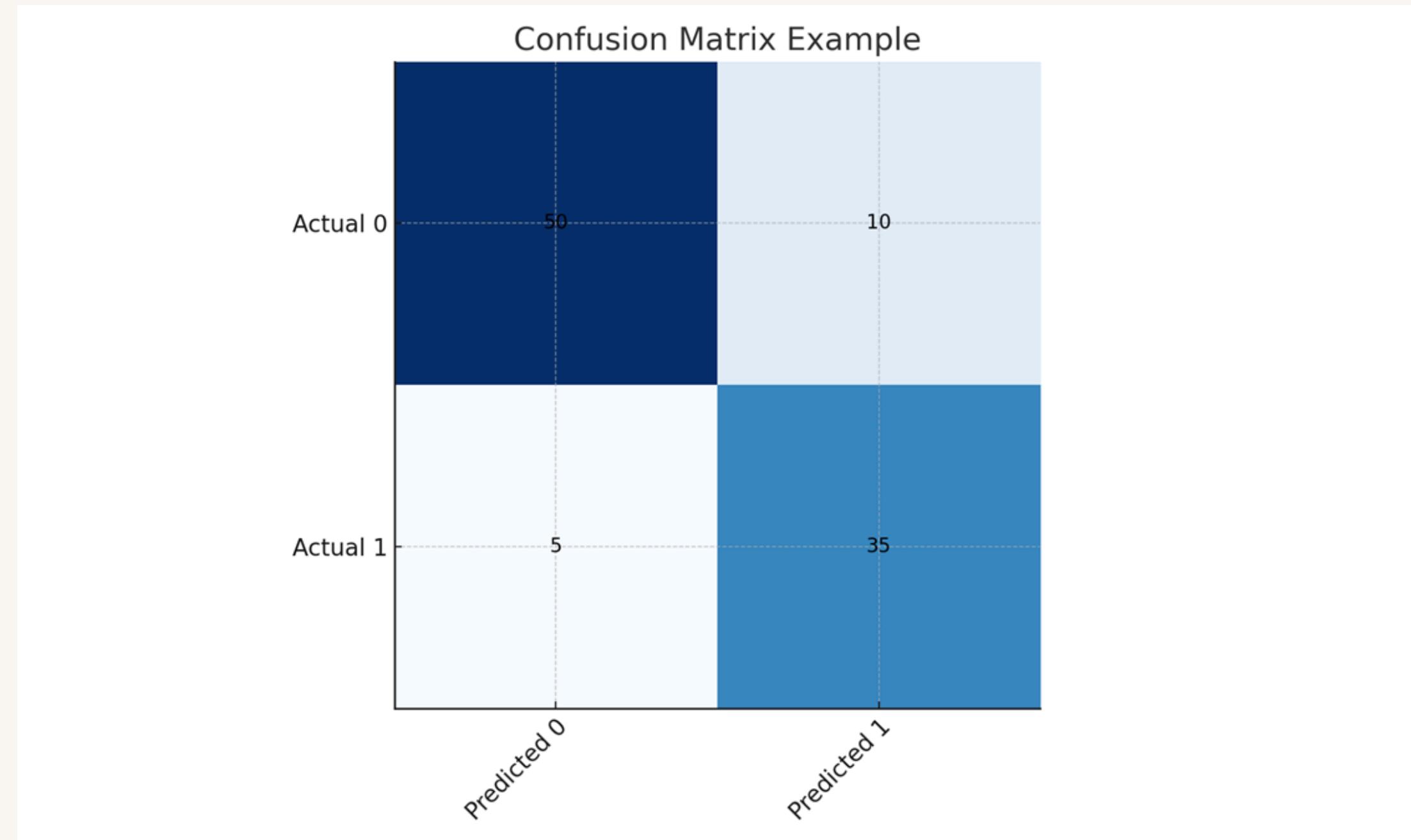
Pipeline with Tabakar Dataset

1. Load data.
2. Handle missing values.
3. Encode categorical data.
4. Feature scaling.
5. Train-test split.
6. Train model (Logistic Regression, etc.).
7. Evaluate with confusion matrix.
8. Save model.

Next Steps

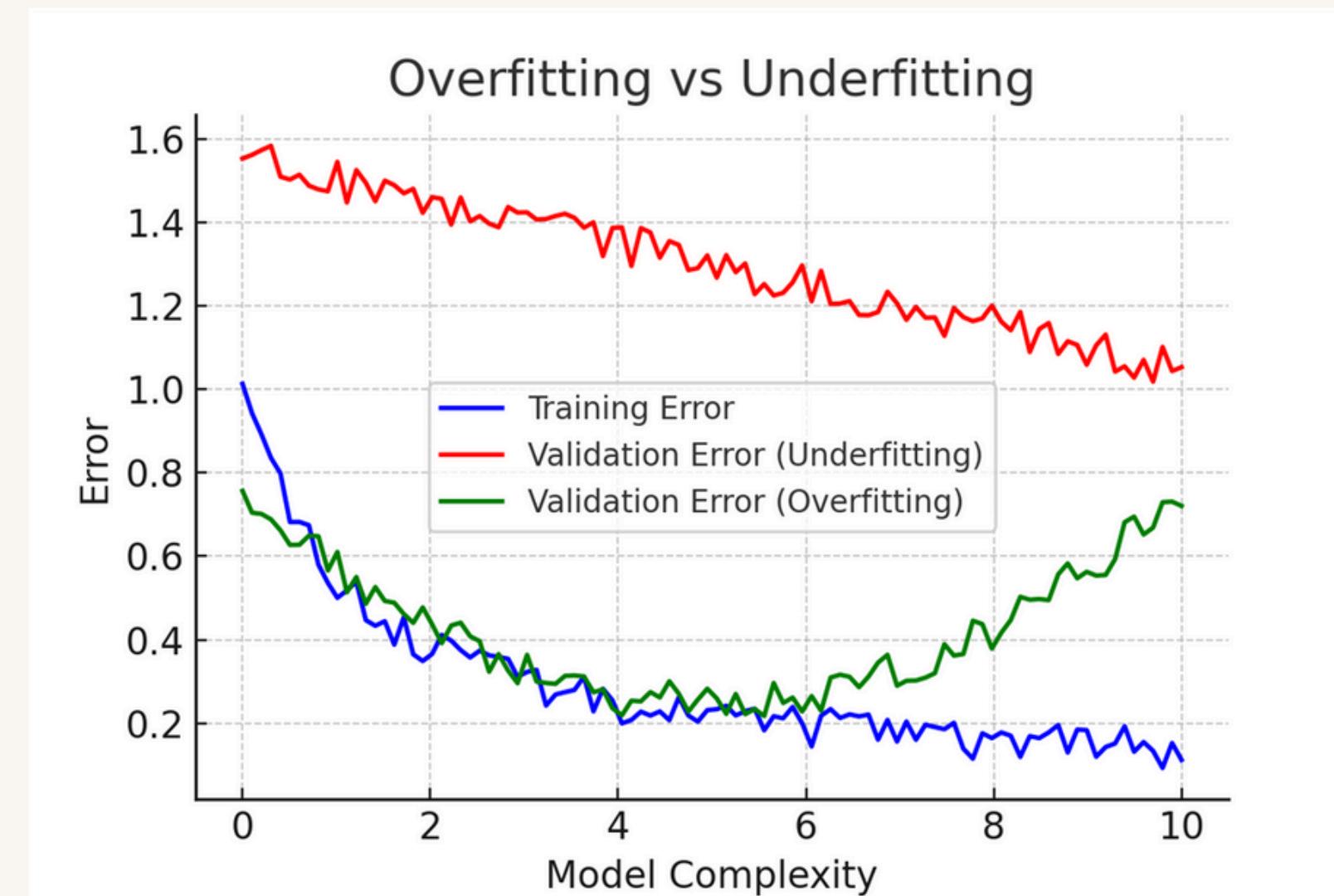
- Hands-on coding session with Tabakar dataset.
- Practice preprocessing + evaluation.
- Compare overfitting/underfitting examples.

Confusion Matrix Example



Overfitting vs Underfitting

Overfitting vs Underfitting



Q n A

Artificial Intelligence & Machine Learning Certificate Course

Supervised Machine Learning & Pipeline Deep Dive

Sehan Arandara
Associate Software Engineer
SE Graduated @ SLIIT
AI ML Instructor

18.09.2025
Day 5

Recap: What is Supervised Learning?

- Learns from labeled data (input features + target output).
- Two main tasks:
 - Regression → predict continuous values (e.g., house prices).
 - Classification → predict categories (e.g., spam/not spam).
- Algorithms you've already introduced: Linear Regression, Logistic Regression.

Step 1 – Data Collection

- Define the problem clearly (e.g., “Will the customer churn?”).
- Gather labeled data from surveys, APIs, sensors, databases, etc.
- Ensure enough quantity + quality for training.

Step 2 – Data Preprocessing

- Handle Missing Values → drop rows/columns, or impute.
- Categorical Data → convert text into numbers (One-Hot Encoding).
- Feature Scaling → standardization or normalization.
- Train-Test Split → separate unseen data for fair evaluation.

Step 3 – Model Training & Cost Function

- Cost Function measures how far predictions are from actuals.
- Regression: Mean Squared Error (MSE).
- Classification: Log Loss / Cross-Entropy.
- The algorithm adjusts parameters to minimize the cost function.

Step 4 – Evaluation

- Classification Metrics: Accuracy, Precision, Recall, F1, Confusion Matrix.
- Regression Metrics: MSE, RMSE, R^2 .
- Cross-validation for more reliable evaluation

Step 5 – Overfitting & Underfitting

- Underfitting → Model too simple (high bias, low variance).
- Overfitting → Model too complex, memorizes training data (low bias, high variance).
- Solutions:
 - Gather more data.
 - Regularization.
 - Use simpler models.
 - Pruning (Decision Trees).
 - Ensemble methods (Random Forest).

Step 6 – Algorithms

- Decision Tree
 - Splits data into rules based on features.
 - Easy to visualize & interpret.
 - Risk: overfitting if tree is too deep.
- Random Forest
 - An ensemble of Decision Trees.
 - Each tree is trained on a random subset of data & features.
 - Reduces overfitting, improves accuracy.
 - More stable than a single tree.

Step 7 – Deployment

- Save trained model (joblib / pickle).
- Build inference function to predict on single inputs.
- Deploy via Flask/FastAPI, or embed into apps.

Regression vs Classification



Regression

- Goal: Predict a continuous numerical value.
- Examples:
 - Predict house prices.
 - Forecast tomorrow's temperature.
 - Estimate exam score.
- Cost Function: Mean Squared Error (MSE), RMSE, MAE.
- Evaluation Metrics: R^2 , RMSE, MSE.
- Algorithms: Linear Regression, Decision Tree Regressor, Random Forest Regressor.

Classification

- Goal: Predict a discrete category (class).
- Examples:
 - Spam vs Not Spam.
 - Customer will purchase (Yes/No).
 - Predict survival on Titanic dataset.
- Cost Function: Log Loss (Cross-Entropy).
- Evaluation Metrics: Accuracy, Precision, Recall, F1-score, Confusion Matrix.
- Algorithms: Logistic Regression, Decision Tree Classifier, Random Forest Classifier.

Regression Metrics Interpretation

- MSE / RMSE / MAE → Lower is better.
 - Compare with the scale of the target values.
 - Example: $\text{MSE} = 0.9 \rightarrow$ Bad if target is 0–1, Good if target is 0–1000.
- R^2 (Coefficient of Determination) → Higher is better.
 - Range: $-\infty$ to 1.
 - 0.0 → No better than average.
 - 0.5 → Model explains 50% variance (moderate).
 - 0.9 → Excellent performance.

Classification Metrics Interpretation

- Accuracy → Good if classes are balanced.
- Precision → Focus on reducing false positives (e.g., fraud detection).
- Recall → Focus on reducing false negatives (e.g., cancer detection).
- F1-score → Balances Precision & Recall (range 0–1).
- Confusion Matrix → Visual inspection of TP, TN, FP, FN.

Overfitting vs Underfitting

Underfitting

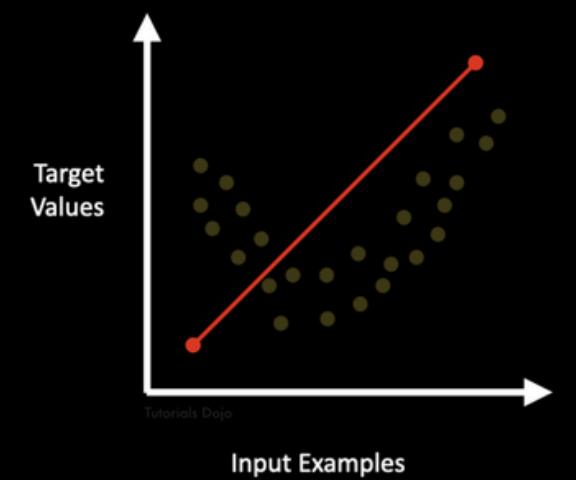
- The model is too simple to capture the data's patterns.
- Happens when:
 - Using a weak algorithm (e.g., linear regression for non-linear data).
 - Not enough features.
 - Too much regularization (forcing model to be too simple).
- Symptoms:
 - Low training accuracy.
 - Low test accuracy.

Overfitting

- The model is too complex, memorizes training data instead of generalizing.
- Happens when:
 - Model has too many parameters (e.g., deep tree).
 - Too few training examples.
 - Noisy data (outliers, irrelevant features).
- Symptoms:
 - Very high training accuracy.
 - Low test accuracy.

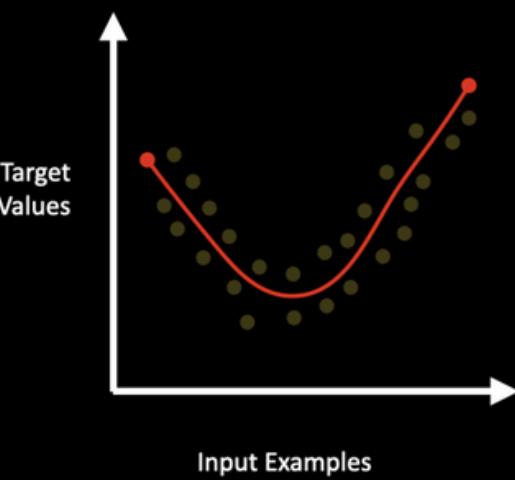
Underfitting vs Overfitting

UNDERFITTING



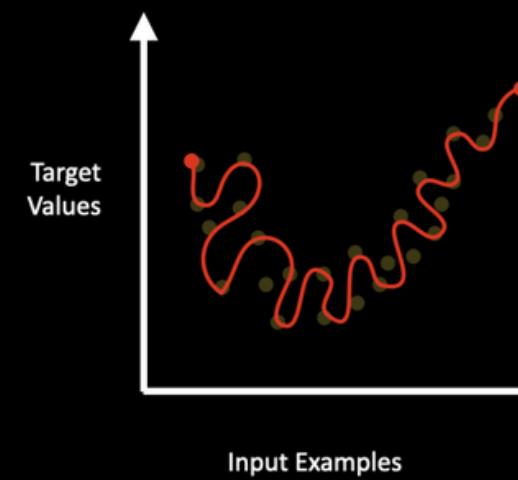
The model performs poorly on the training data.

BALANCED



The model performs optimally as expected

OVERFITTING



The model performs well on the training data but does not perform well on the evaluation data.

Regression (Prediction) Algorithms

- Linear Regression → Simple & interpretable, best for linear data.
- Polynomial Regression → Captures curved trends.
- Decision Tree Regressor → Handles non-linear relationships.
- Random Forest Regressor → Reduces overfitting, very accurate.
- Support Vector Regressor (SVR) → Works well in high dimensions.
- Gradient Boosting (XGBoost, LightGBM) → Powerful on structured data.

Classification Algorithms

- Logistic Regression → Fast baseline for binary classification.
- K-Nearest Neighbors (KNN) → Classifies based on neighbors.
- Decision Tree Classifier → Easy to visualize, prone to overfitting.
- Random Forest Classifier → Robust and accurate.
- Support Vector Machine (SVM) → Strong boundaries between classes.
- Naïve Bayes → Works well with text data.
- Neural Networks → Best for complex tasks (images, speech).



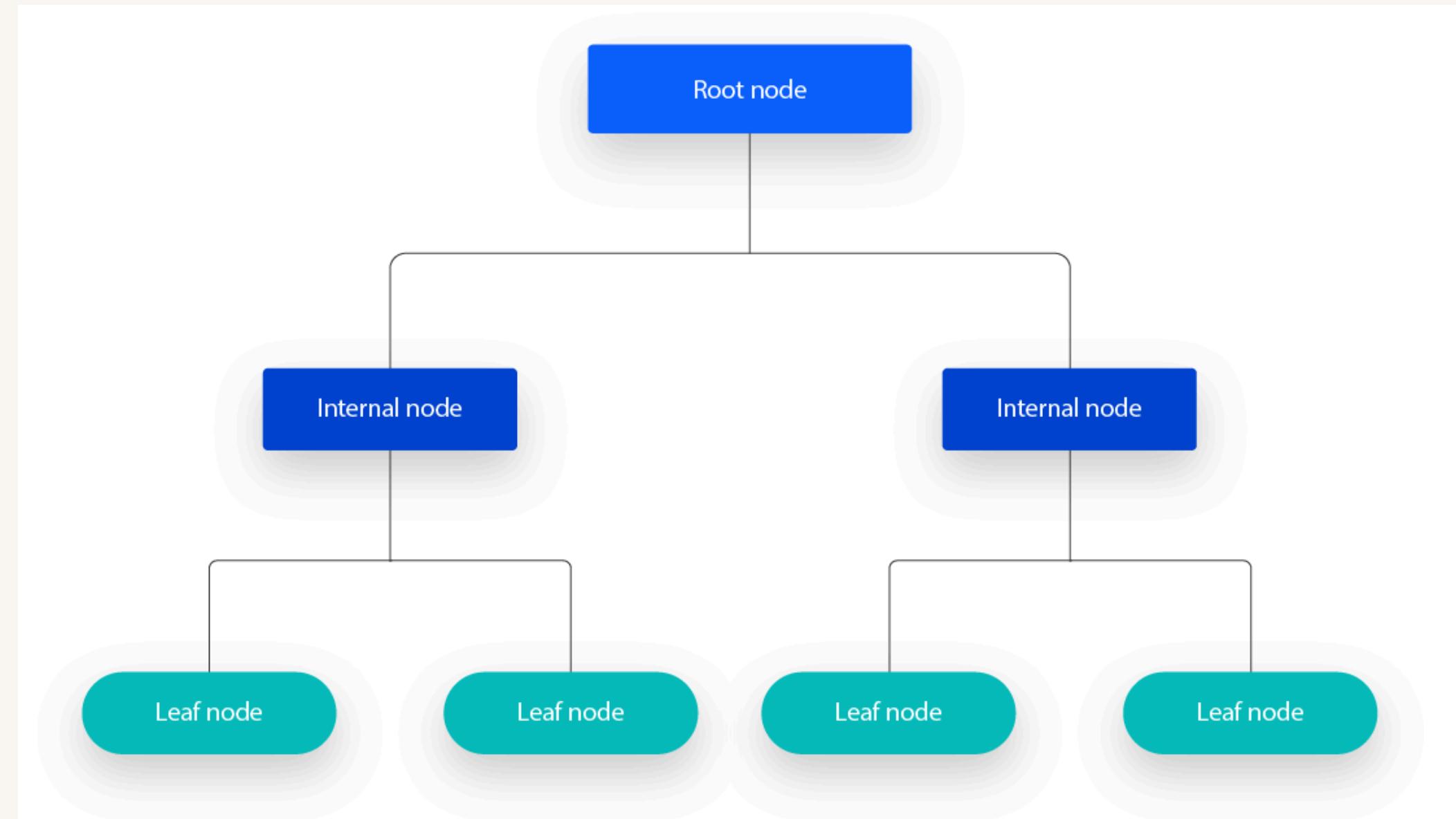
Decision Trees

- A tree-like model of decisions.
- Each node → a question (e.g., Age > 30?).
- Each branch → an answer (Yes/No).
- Each leaf → a prediction (class or value).
- Easy to understand & visualize.
- Can easily overfit (too many splits).

How Decision Trees Work

- Start with all data at the root.
- Choose the feature & split that best separates the target.
 - Classification → Gini Index, Entropy.
 - Regression → Mean Squared Error.
- Repeat splitting recursively until stopping condition.
- Leaves give final prediction.

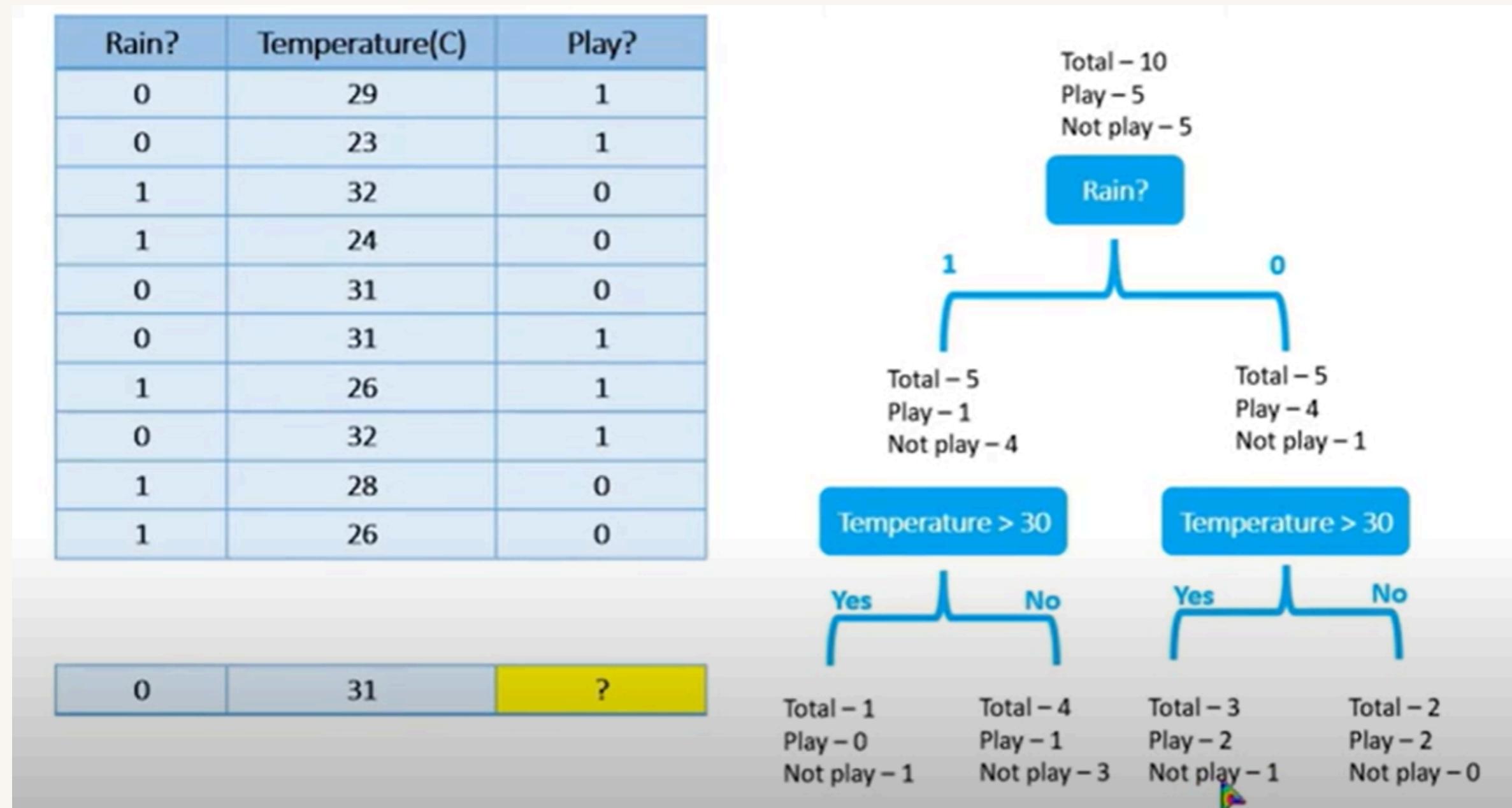
How Decision Trees Work



How Decision Trees Work



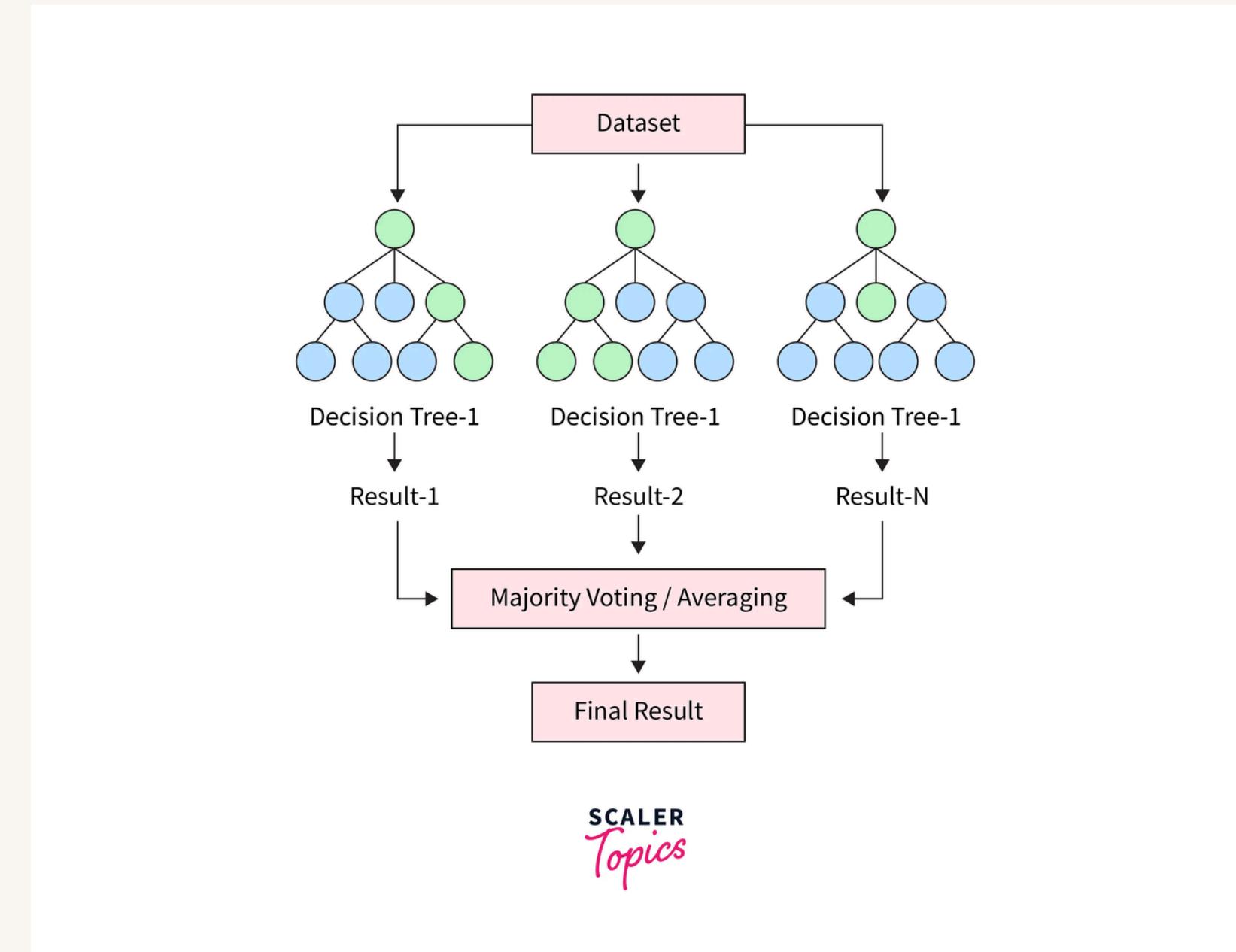
Example :





Random Forest

- An ensemble of many Decision Trees.
- Each tree trained on a random subset of data + features.
- Final prediction = majority vote (classification) or average (regression).
- Reduces overfitting of single trees.
- More accurate & robust.
- Harder to interpret than a single tree.



Q n A

Artificial Intelligence & Machine Learning Certificate Course

Unsupervised machine learning

Sehan Arandara
Associate Software Engineer
SE Graduated @ SLIIT
AI ML Instructor

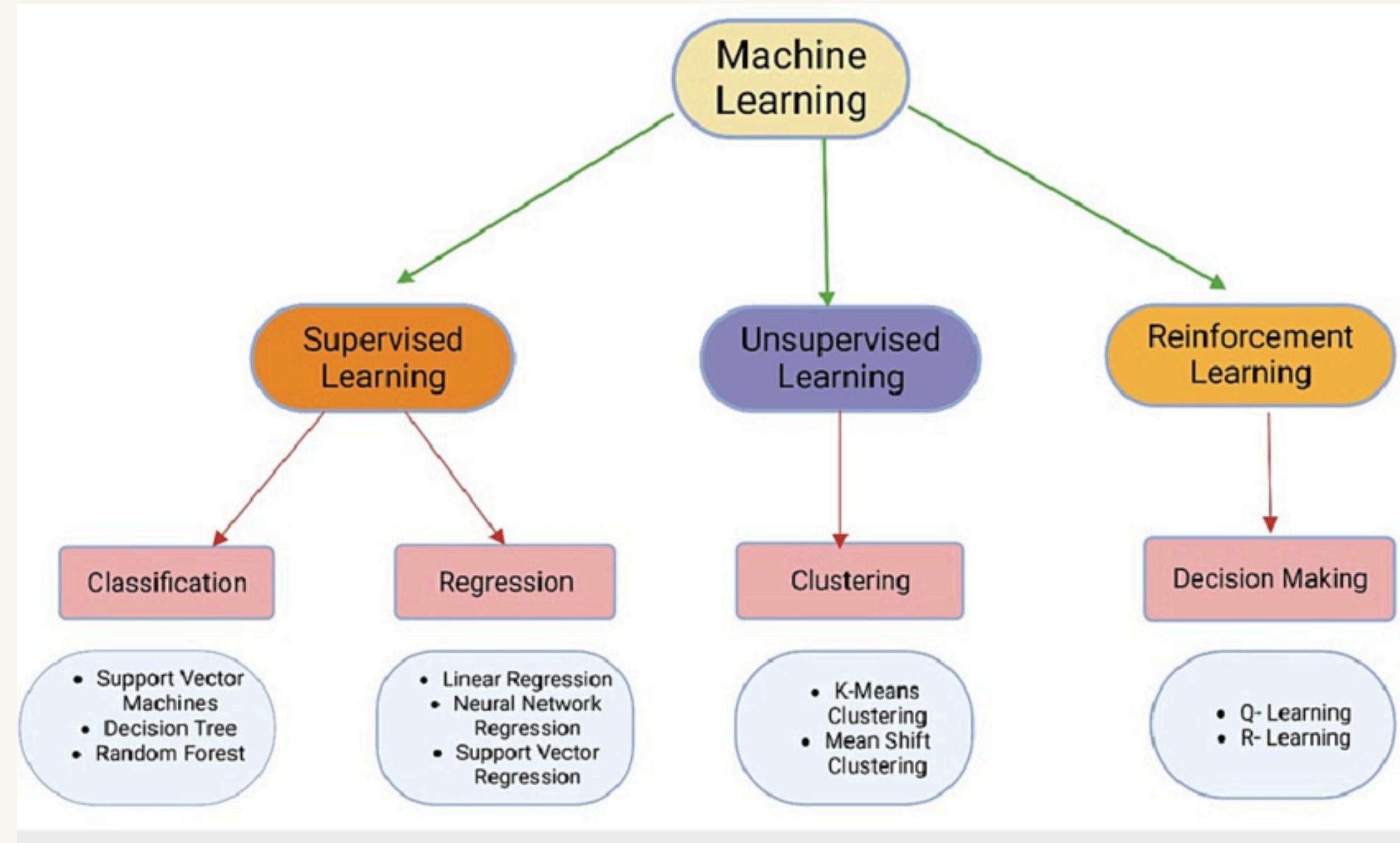
26.09.2025

Day 7

Recap: Labeled vs. Unlabeled Data

- Labeled Data
 - data set that defined the features and the target
- Unlabeled Data
 - This is raw data with no predefined labels or answers. You only have the input data. The algorithm must analyze this data to find its own structure, groupings, or patterns.

Recap : Types of ML



What is Unsupervised Machine Learning?

- An AI model learns from unlabeled data.
- Goal: To find hidden patterns, structures, and relationships on its own.
- It's "unsupervised" because there is no teacher or answer key to guide it.
- The algorithm acts like a data detective 🕵️, exploring the data for clues

Unsupervised Task: Clustering

- What it is: Automatically grouping similar data points together.
- Goal: Data points in the same cluster are similar to each other.
- Real-World Use Case: Customer Segmentation
 - Group customers by their purchasing behavior (e.g., "Bargain Hunters," "Loyal Fans").
 - Allows for effective, targeted marketing.

Unsupervised Task: Association

- What it is: Discovering "if-then" rules in your data.
- Famous Example: Market Basket Analysis
- Finds rules like: "If a customer buys diapers, they are also likely to buy beer."
- Used for product recommendations ("Frequently bought together") and store layout optimization.

Unsupervised Task: Dimensionality Reduction

- What it is: Simplifying complex data by reducing the number of features (variables).
- Goal: Keep the most important information while removing noise.
- Real-World Use Case: Data Visualization
 - Makes it possible to plot high-dimensional data (e.g., 50+ features) in 2D or 3D to see its structure.

Clustering in Machine Learning

- Clustering is an unsupervised machine learning technique that groups similar data points together into clusters based on their characteristics, without using any labeled data. The objective is to ensure that data points within the same cluster are more similar to each other than to those in different clusters, enabling the discovery of natural groupings and hidden patterns in complex datasets.

Clustering in Machine Learning

- Goal: Discover the natural grouping or structure in unlabeled data without predefined categories.
- How: Data points are assigned to clusters based on similarity or distance measures.
- Similarity Measures: Can include Euclidean distance, cosine similarity or other metrics depending on data type and clustering method.
- Output: Each group is assigned a cluster ID, representing shared characteristics within the cluster.

Clustering in Machine Learning



Clustering Use cases

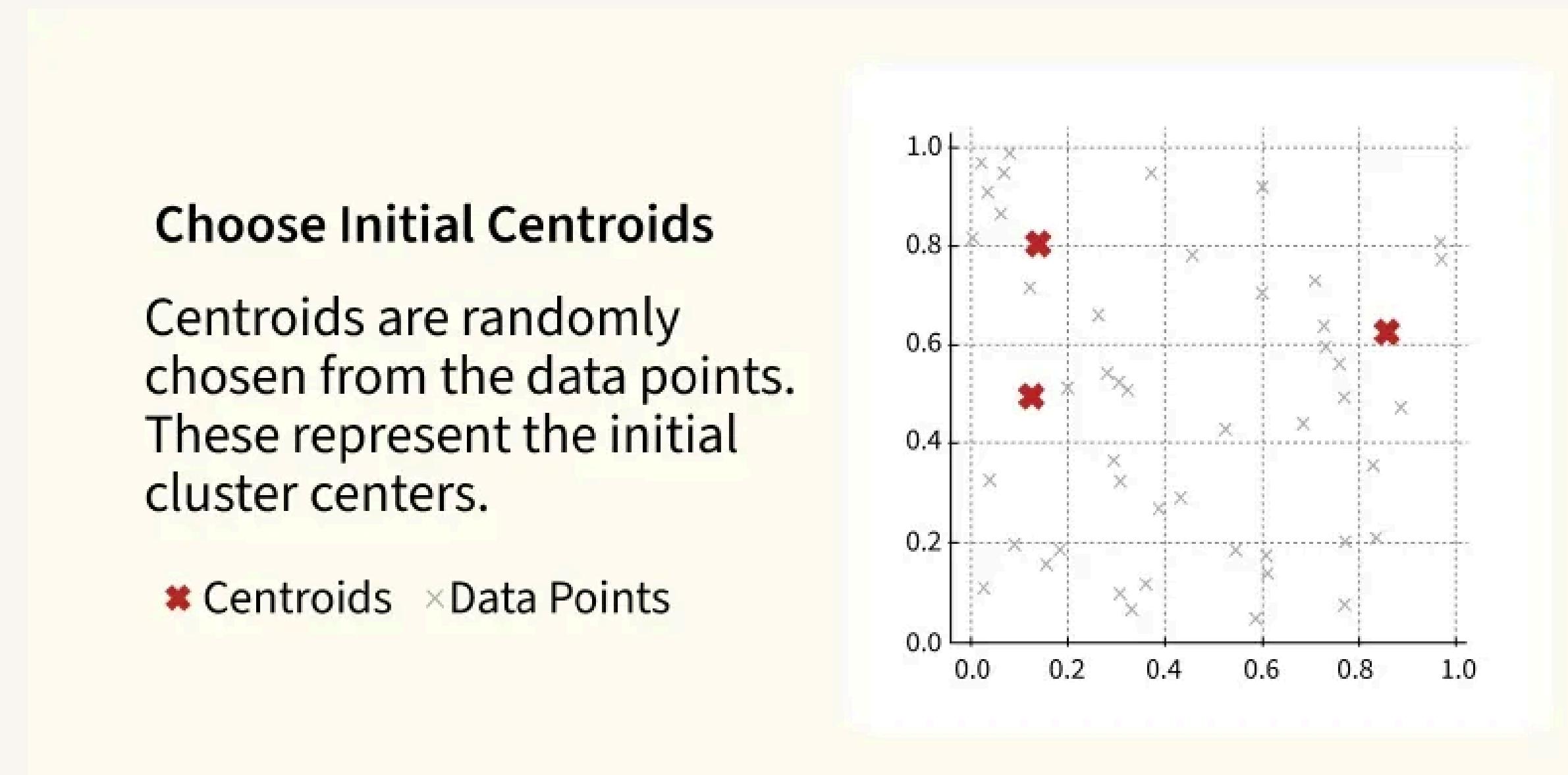
- Customer Segmentation: Grouping customers based on behavior or demographics for targeted marketing and personalized services.
- Anomaly Detection: Identifying outliers or fraudulent activities in finance, network security and sensor data.
- Image Segmentation: Dividing images into meaningful parts for object detection, medical diagnostics or computer vision tasks.
- Recommendation Systems: Clustering user preferences to recommend movies, products or content tailored to different groups.
- Market Basket Analysis: Discovering products frequently bought together to optimize store layouts and promotions.

K means Clustering

K means Clustering

- K-Means Clustering is an unsupervised machine learning algorithm that helps group data points into clusters based on their inherent similarity .
- Example : an online store can use K-Means to segment customers into groups like "Budget Shoppers," "Frequent Buyers," and "Big Spenders" based on their purchase history.
- "k" represents the number of groups or clusters we want to classify our items into.

K means Clustering

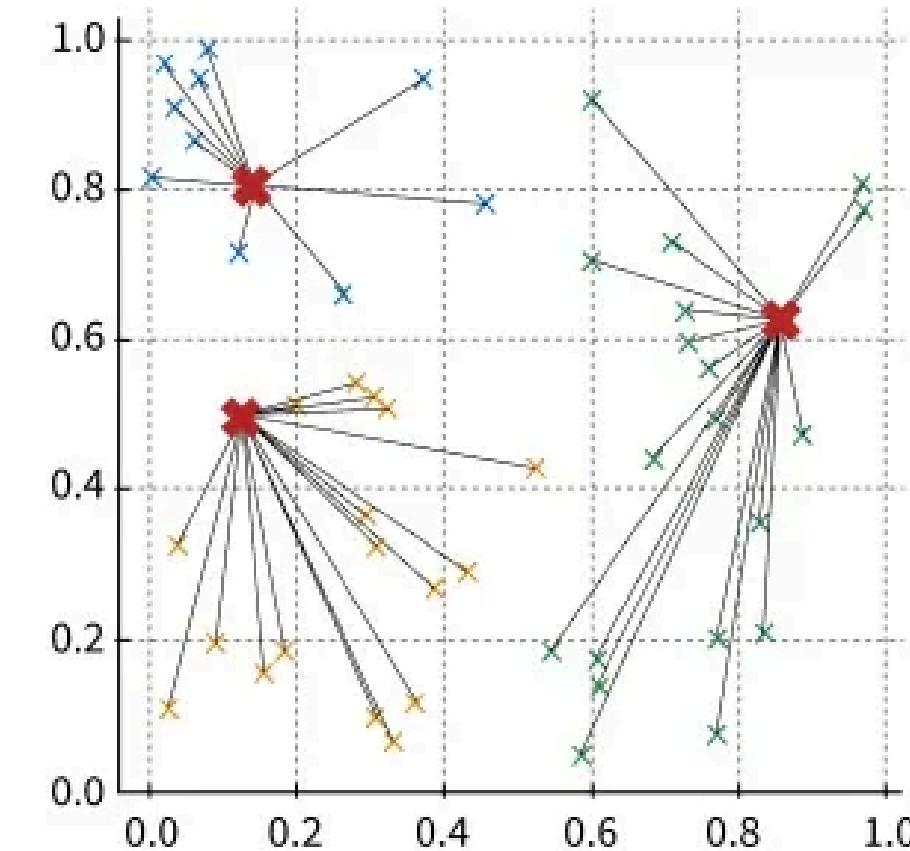


K means Clustering

Assign Points to Nearest Centroid

Each point is assigned to the nearest centroid, forming clusters

- ✖ Centroids ✕ Cluster 1
- ✖ Cluster 2 ✗ Cluster 3

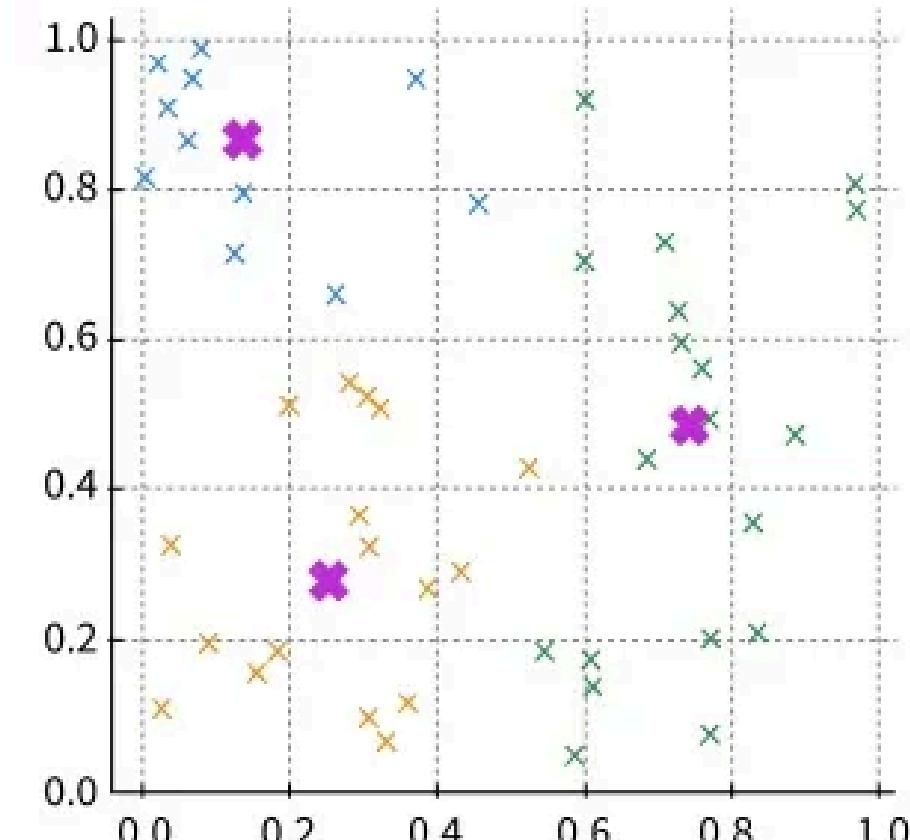


K means Clustering

Update Centroids

Centroids are recalculated as the mean of the points in each cluster

- ✿ New Centroids
- ✖ Cluster 1
- ✖ Cluster 2
- ✖ Cluster 3



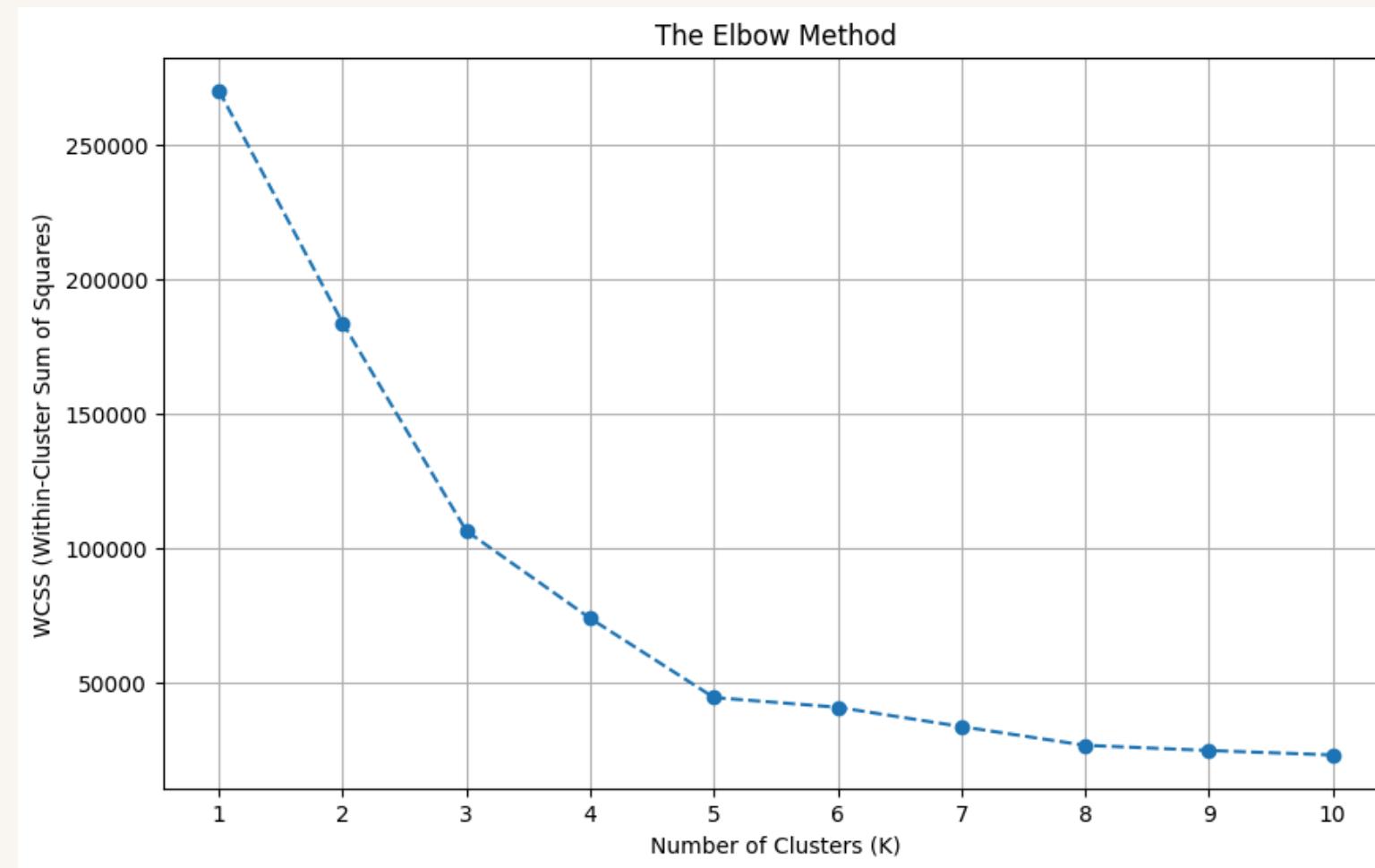
Finding the Optimal 'K' (The Elbow Method)



- How do we know if we should use K=3, K=4, or K=5 clusters? The Elbow Method helps us decide .
- We'll run the K-Means algorithm for a range of different 'K' values (say, 1 to 10)
- we'll calculate something called the WCSS (Within-Cluster Sum of Squares). WCSS is basically the sum of the squared distances between each point and its cluster's center.
- A smaller WCSS means the clusters are more compact and therefore better.

Finding the Optimal 'K' (The Elbow Method) 💪

- When we plot WCSS against the number of clusters, the graph will look like an arm. The point where the graph "bends" or forms an elbow is our best 'K'. This is the point of diminishing returns, where adding another cluster doesn't significantly improve the compactness.

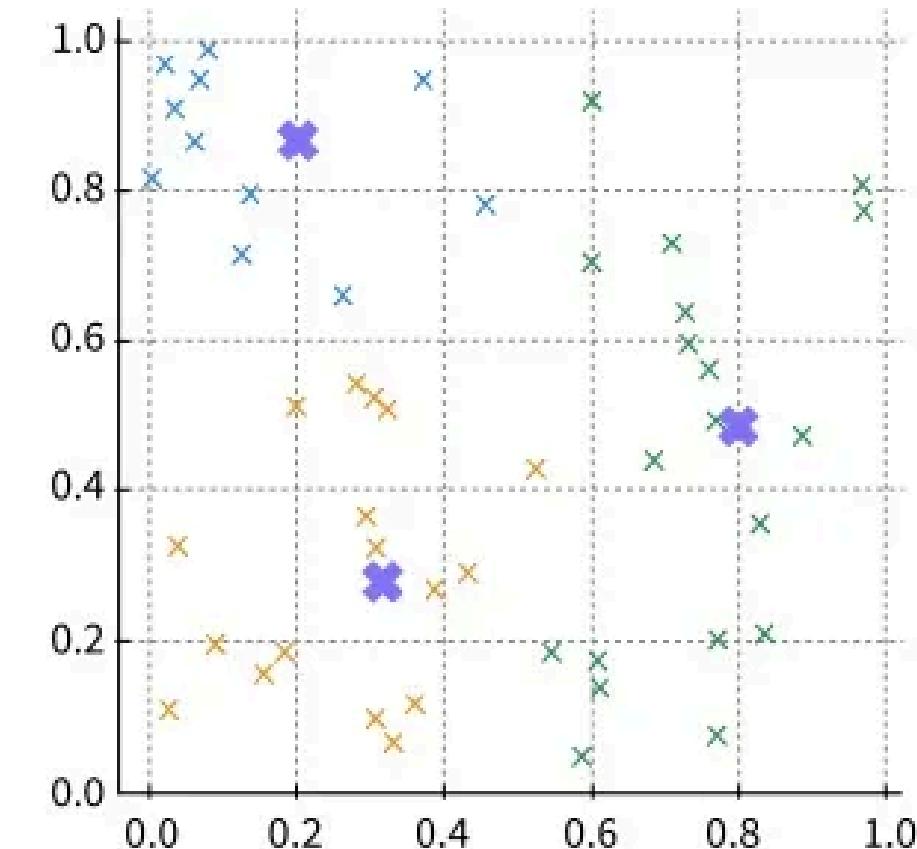


K means Clustering

Repeat Until Convergence

This process repeats until the centroids stabilize and do not move further.

- ✿ Final Centroids
- ✖ Cluster 1
- ✖ Cluster 2
- ✖ Cluster 3



Why Use K-Means Clustering?

- K-Means is popular in a wide variety of applications due to its simplicity, efficiency and effectiveness. Here's why it is widely used:
 - Data Segmentation
 - Image Compression
 - Anomaly Detection

Q n A

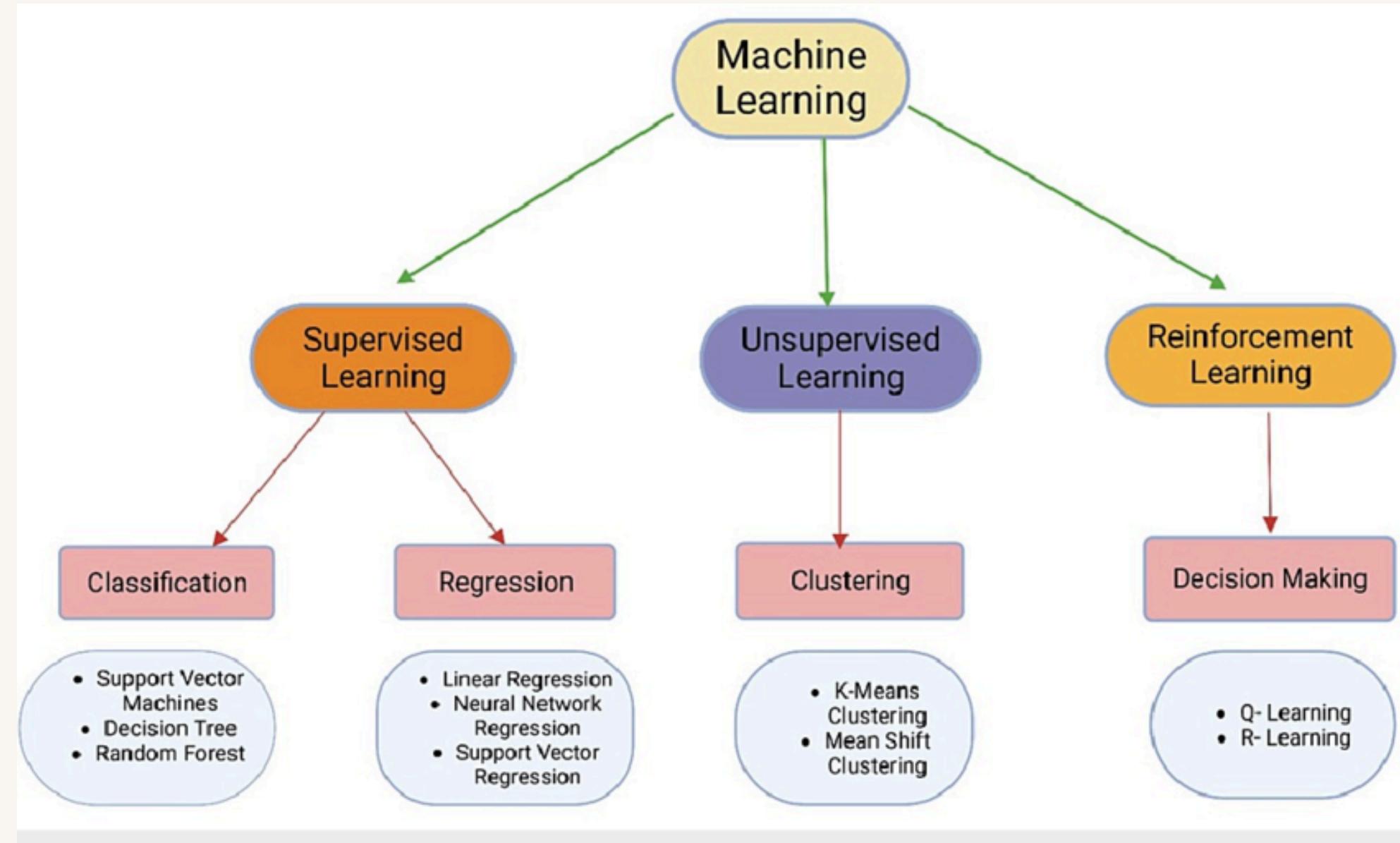
Artificial Intelligence & Machine Learning Certificate Course

Reinforcement Learning

Sehan Arandara
Associate Software Engineer
SE Graduated @ SLIIT
AI ML Instructor

03.10.2025
Day 9

Recap : Types of ML



Supervised, Unsupervised, and Reinforcement Learning

- **Supervised Learning:** Think of this as learning with a strict teacher. You provide the algorithm with labeled data (e.g., pictures of cats labeled "cat," or house features labeled with their correct price). The algorithm learns to map inputs to correct outputs. Its goal is to predict or classify. Each piece of data is often assumed to be independent.
 - Example: Training an email filter with emails already marked as "spam" or "not spam."

Supervised, Unsupervised, and Reinforcement Learning

- **Unsupervised Learning:** Here, there's no teacher and no labels. The algorithm's job is to explore unlabeled data and find hidden patterns or structures. We just finished this with K-Means, where we found customer segments without knowing them beforehand. The goal is to uncover patterns. Like supervised learning, it often assumes data records are independent.
 - Example: Grouping customers based on their shopping behavior without knowing their categories beforehand.

Supervised, Unsupervised, and Reinforcement Learning

- **Reinforcement Learning (RL):** This is different. RL doesn't use labeled data for correct behavior, nor does it just find hidden patterns. Instead, an agent learns by "doing"—interacting with an environment through trial-and-error. It receives rewards for good actions and sometimes penalties for bad ones
 - Example: Training a dog or teaching a robot to walk.

What is Reinforcement Learning (RL)?

- Imagine you're training a dog. 🐶
- When your dog performs a trick correctly, you give it a treat (a reward).
- When it misbehaves, you might ignore it or give a firm "No" (a penalty or no reward).
- Over time, the dog learns which actions lead to treats and which don't. It's learning to maximize its rewards

What is Reinforcement Learning (RL)?

- Technically, Reinforcement Learning is a framework for learning how to interact with an environment from experience.
- It's all about an "agent" figuring out the best sequence of actions to take in a given situation to maximize cumulative rewards over time.
- This concept is deeply inspired by how animals (and humans!) learn through interaction and feedback.

Real-World Examples

- Robotics: Training a robot to walk or manipulate objects.
- Games: Mastering complex strategy games like Chess or Go (e.g., AlphaGo).
- Self-Driving Cars: Making navigation decisions in real-time traffic.



The Core Concepts Inside RL: The Building Blocks

What is Reinforcement Learning (RL)?

- Agent : The learner and decision-maker. The "brain" of the operation.
 - Example: The self-driving car's AI.
- Environment : The world the agent interacts with.
 - Example: The city streets, traffic, and pedestrians.
- State (s): The agent's current situation or observation.
 - Example: The car's current speed, location, and sensor readings.

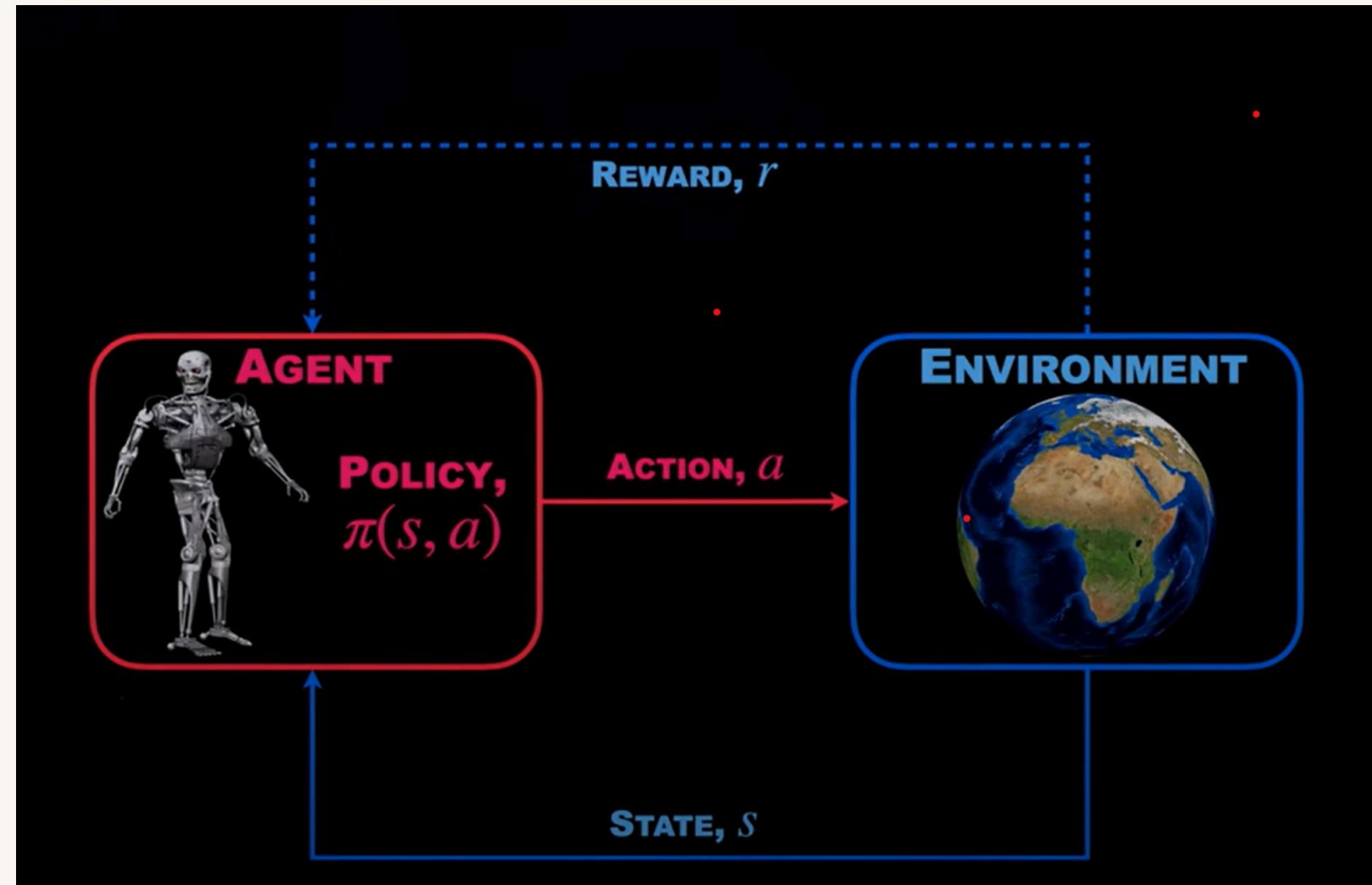
What is Reinforcement Learning (RL)?

- Policy (π): Agent's strategy; mapping states to actions
 - e.g., if in this state, take this action.
- Value Function (V): Predicts future rewards from a state (how good is this state?).

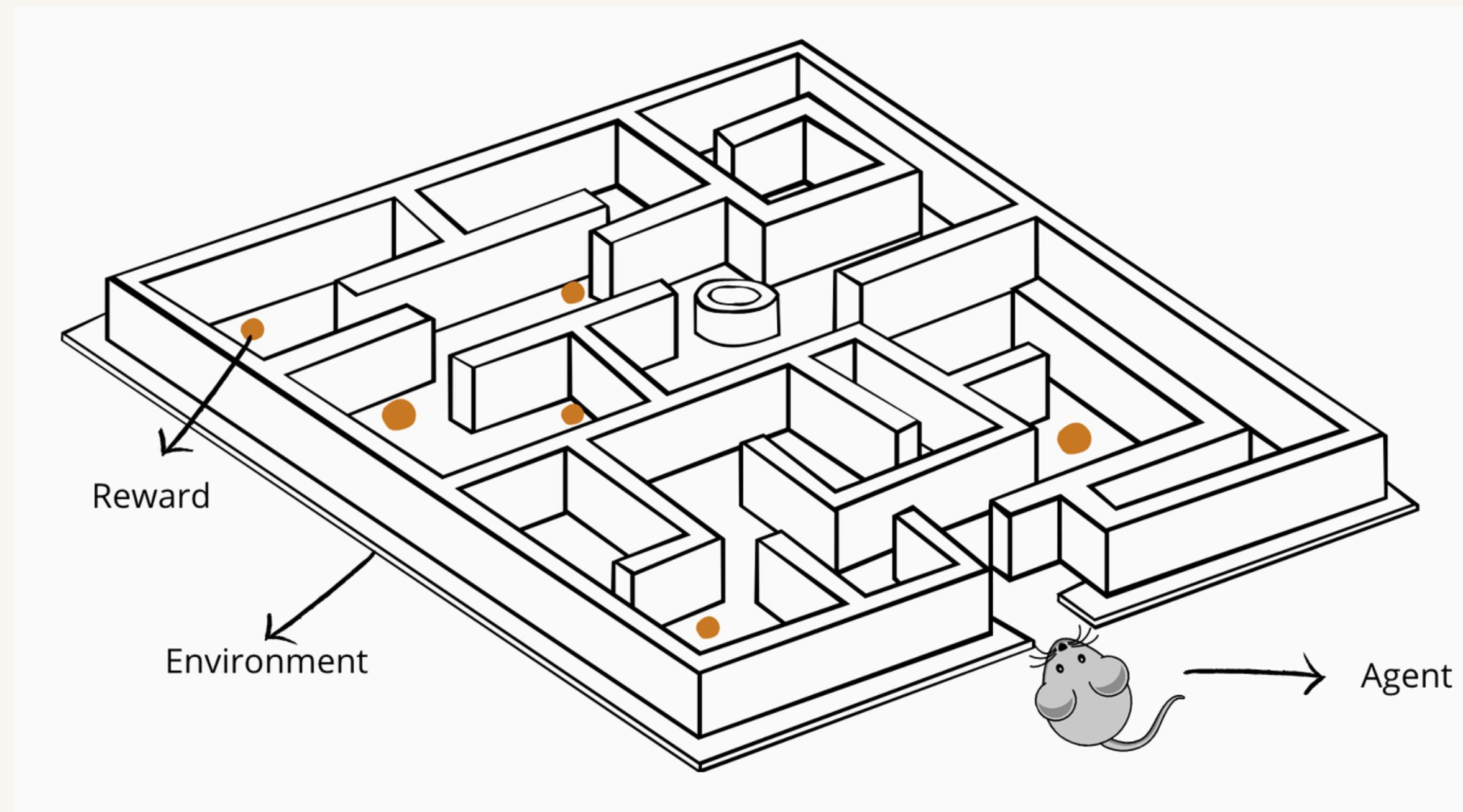
The Reinforcement Learning Loop

- Agent observes State (s) from Environment.
- Agent chooses Action (a) based on its Policy (π).
- Environment receives Action (a).
- Environment transitions to new State (s') and gives Reward (r) to Agent.
- Agent learns from Reward to update its Policy.

The Reinforcement Learning Loop



The Reinforcement Learning Loop



Main Concepts and Challenges in RL

- The ultimate goal of an RL agent is to optimize its policy (π) to get the maximum cumulative future rewards. It's not just about the immediate reward, but the total long-term gain.

Designing Policies and Value Functions

- The Policy (π):
 - The policy is the agent's strategy. It tells the agent what action to take in a given state.
 - Often, policies are probabilistic (e.g., 80% chance to move left, 20% chance to move right). This allows for exploration (trying new things) even when the agent thinks it knows the best path.

Designing Policies and Value Functions

- The Value Function (V):
 - The value function, $V(s | \pi)$, estimates how good it is to be in a particular state (s) and follow a specific policy (π). It's the expected total future reward the agent will receive from that state onwards.
 - Discount Rate (γ): When calculating future rewards, we use a discount rate (γ , between 0 and 1). This means immediate rewards are valued more highly than future rewards. This helps the agent prioritize and is a practical way to deal with infinite horizons.

The Core Challenges in RL

- Sparse Rewards:
 - Rewards can be rare or only given at the very end of a task (e.g., you only get a reward for "winning" a chess game, not for each good move).
 - This makes learning very slow and difficult.
- The Credit Assignment Problem:
 - If you win a long game, which of your hundreds of moves was the "winning" move?
 - It's extremely hard to assign credit to specific actions when the reward is delayed.
This is the central challenge in RL.

The Exploration vs. Exploitation

- Every RL agent faces a constant choice:
 - Exploitation:
 - What it is: Stick to what you know works best. Go to your favorite restaurant.
 - Goal: Maximize immediate, known rewards.
 - Exploration:
 - What it is: Try something new and unknown. Go to a new restaurant you've never been to.
 - Goal: Discover potentially better strategies for the future.
 - A good agent must balance both to achieve long-term success.

Q n A

Artificial Intelligence & Machine Learning Certificate Course

Deep Learning

Sehan Arandara

Associate Software Engineer

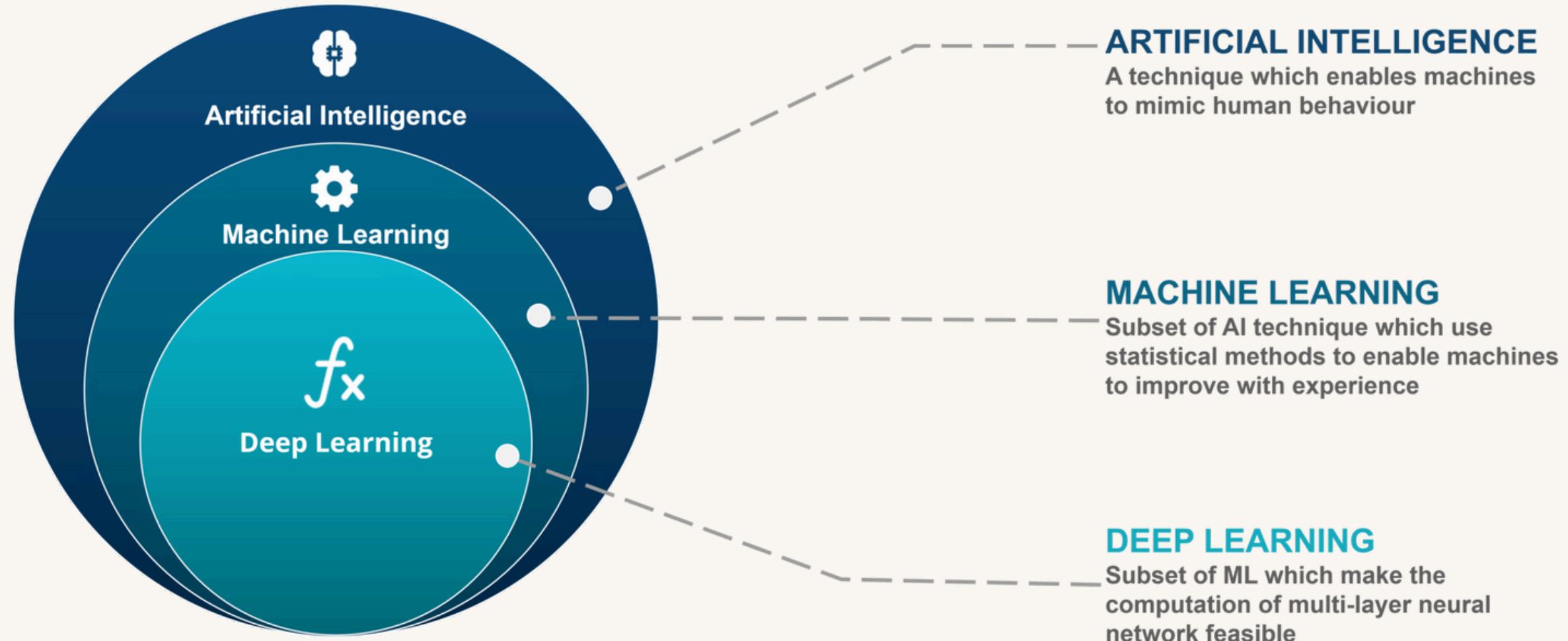
SE Graduated @ SLIIT

AI ML Instructor

12.10.2025

Day 11

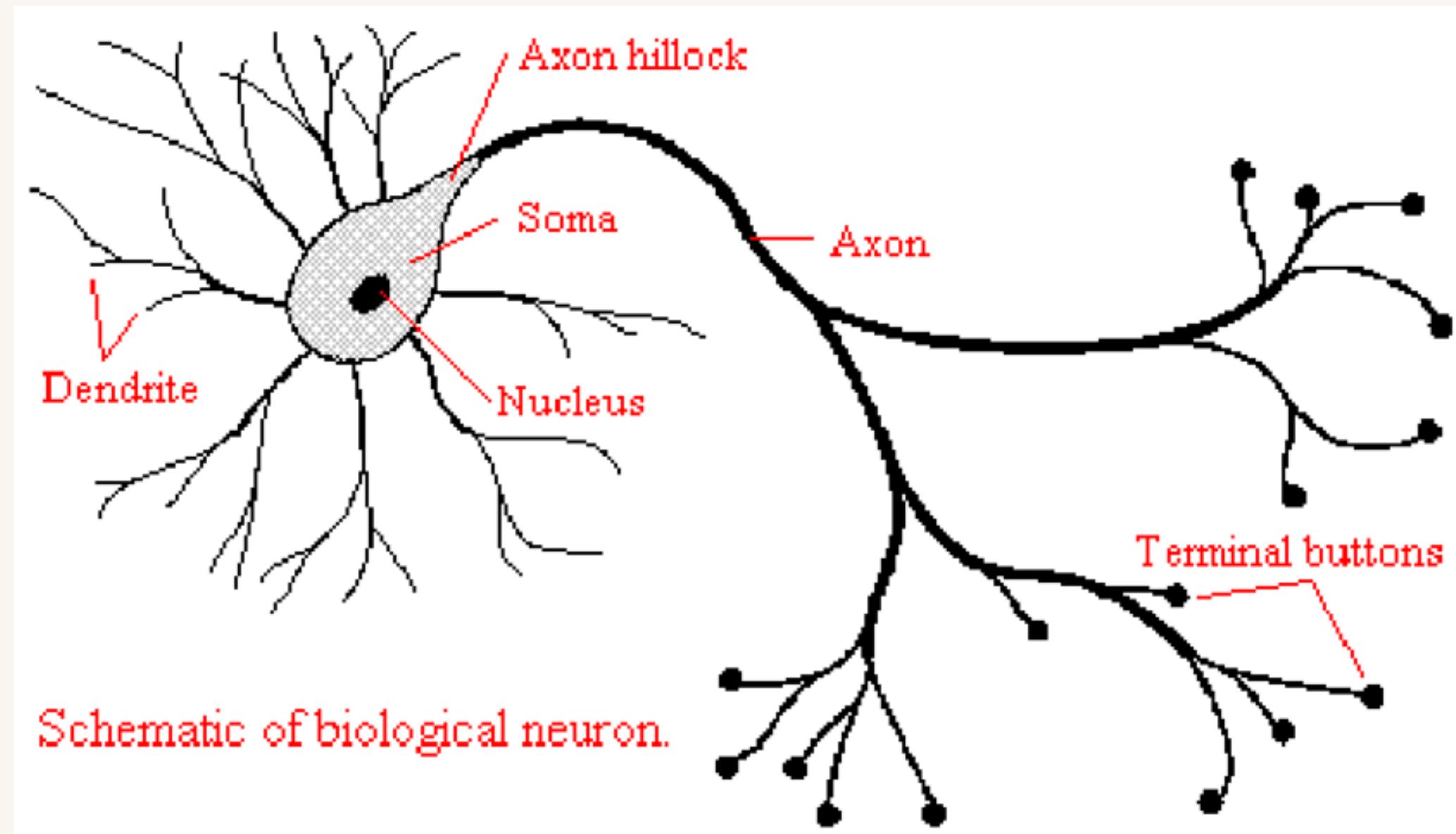
AI vs ML vs DL



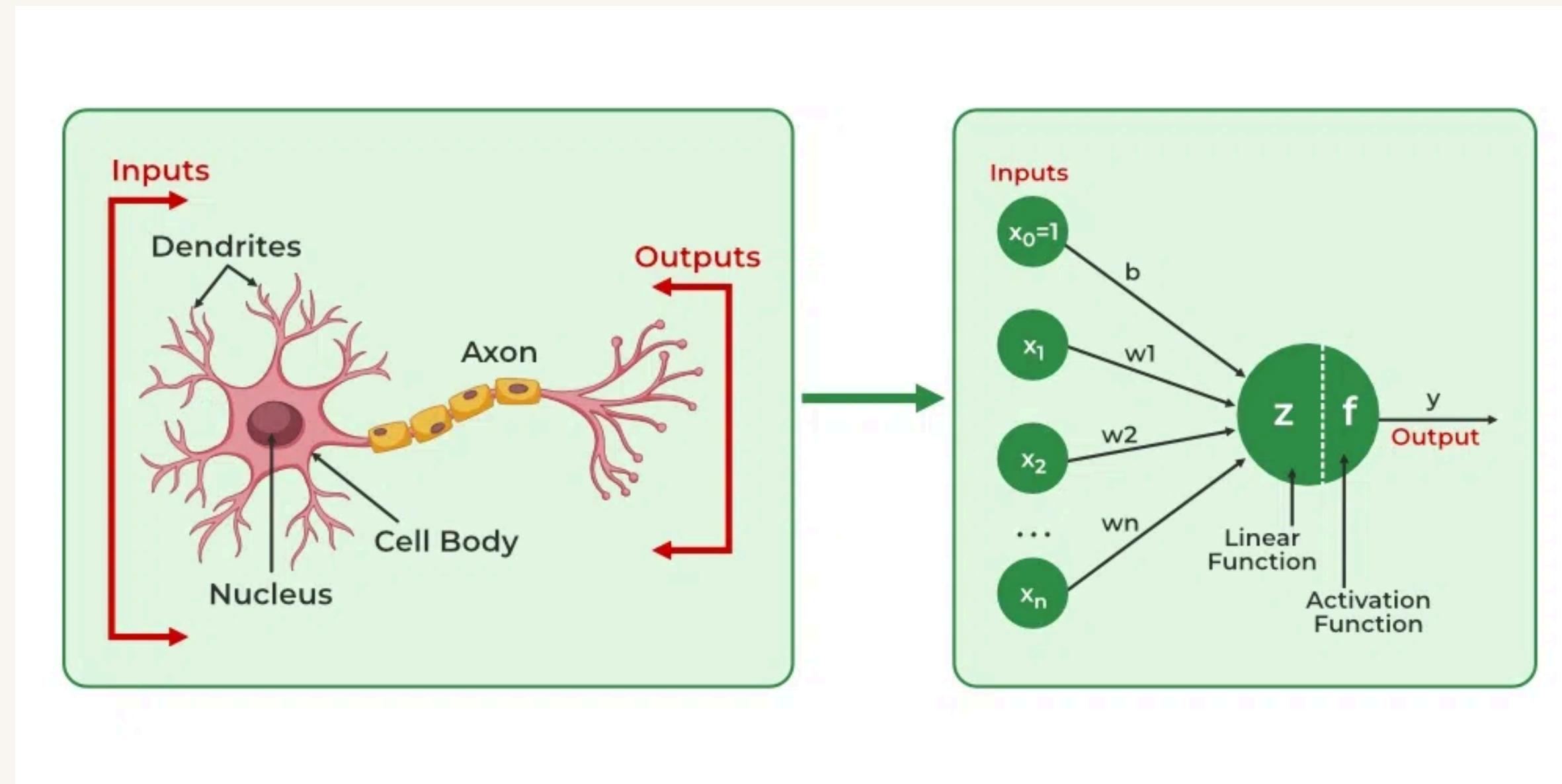
ML vs DL

Feature	Traditional Machine Learning	Deep Learning
Data Needs	Works well with small to medium datasets.	Requires large datasets ("Big Data").
Feature Engineering	Manual: A human expert must select and create features.	Automatic: The network learns the best features from the data itself.
Data Type	Best for structured data (tables, spreadsheets).	Excels at unstructured data (images, text, audio).
"Intelligence"	Learns from the features you give it.	Learns features and patterns on its own.

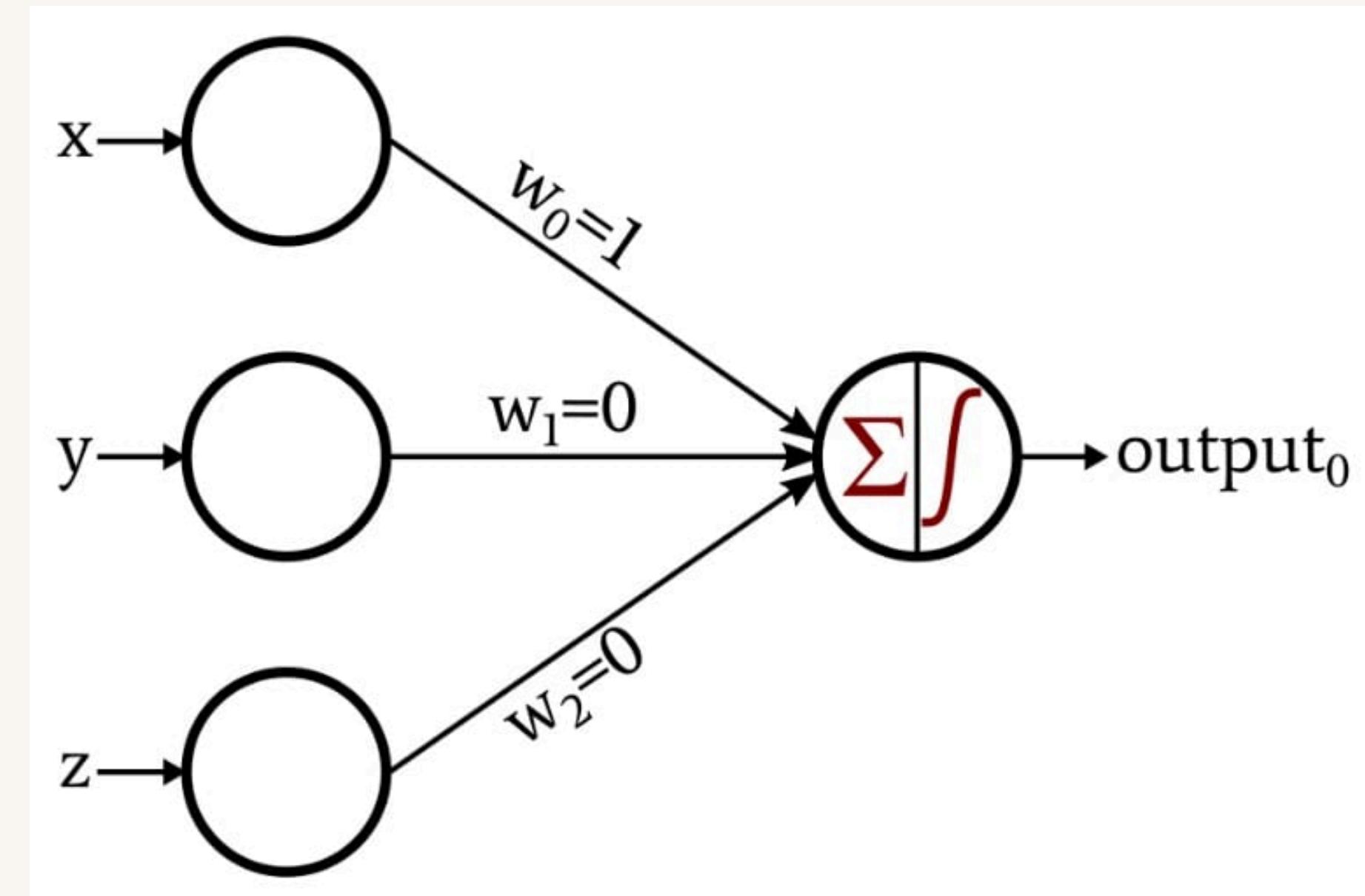
The Biological Brain & The Artificial Neuron



The Biological Brain & The Artificial Neuron



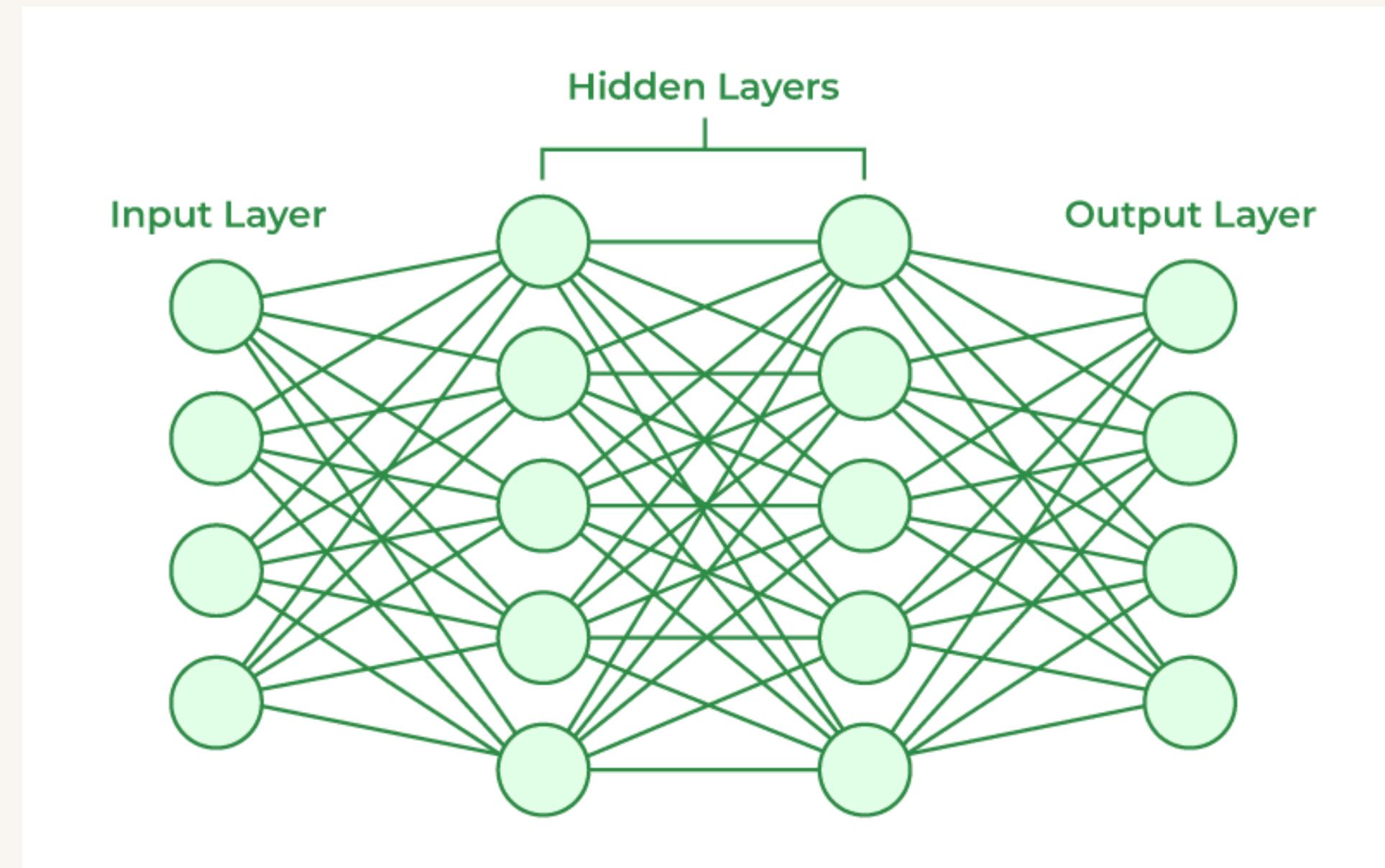
The Biological Brain & The Artificial Neuron



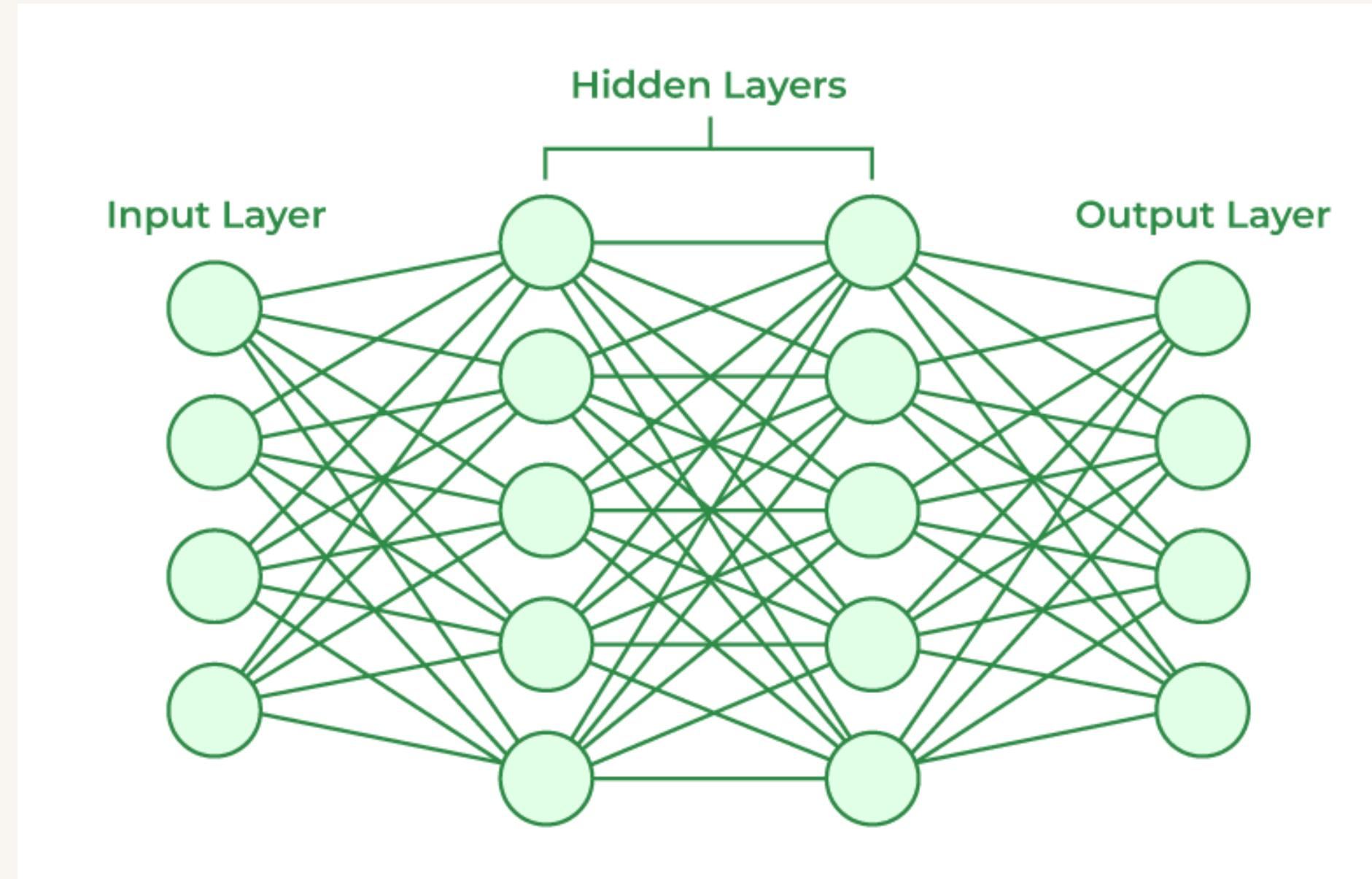
The Biological Brain & The Artificial Neuron

- Inputs: Signals coming from other neurons (or our data).
- Weights: How important is each signal? A high weight means a signal has a big effect on the dimmer.
- Bias: This is like a baseline setting. Even with no input, should the dimmer be slightly on or off?
- Activation Function: After adding up all the weighted signals, the neuron needs to 'decide' how bright to shine. The activation function takes the total and squashes it into a useful output, like a value between 0 (off) and 1 (fully on)."

What is an Artificial Neural Network (ANN)



What is an Artificial Neural Network (ANN)

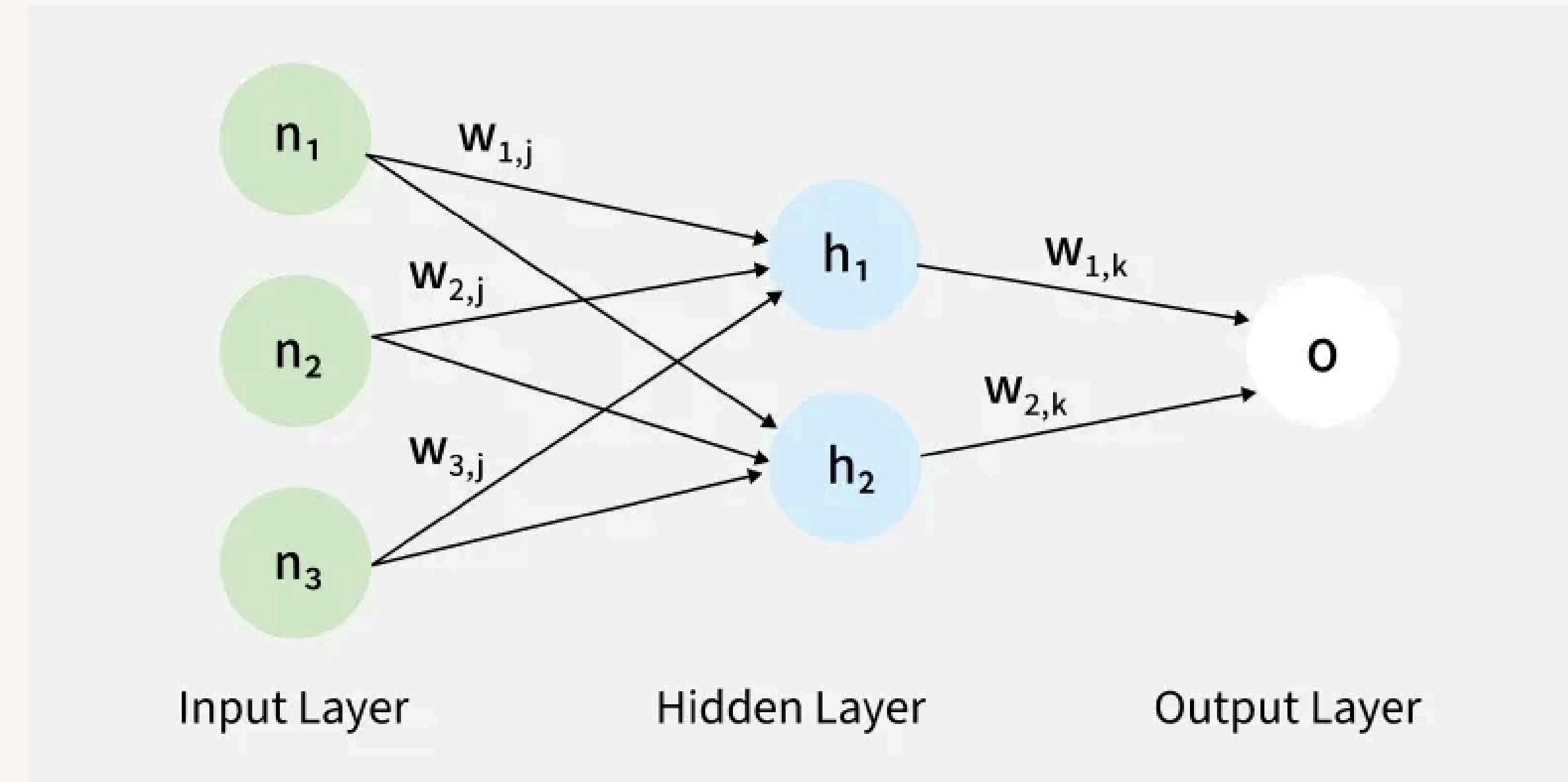


- **Input Layer:** Receives the initial data.
- **Hidden Layer(s):** The "thinking" part of the brain. The more hidden layers, the "deeper" the network.
- **Output Layer:** Produces the final prediction.

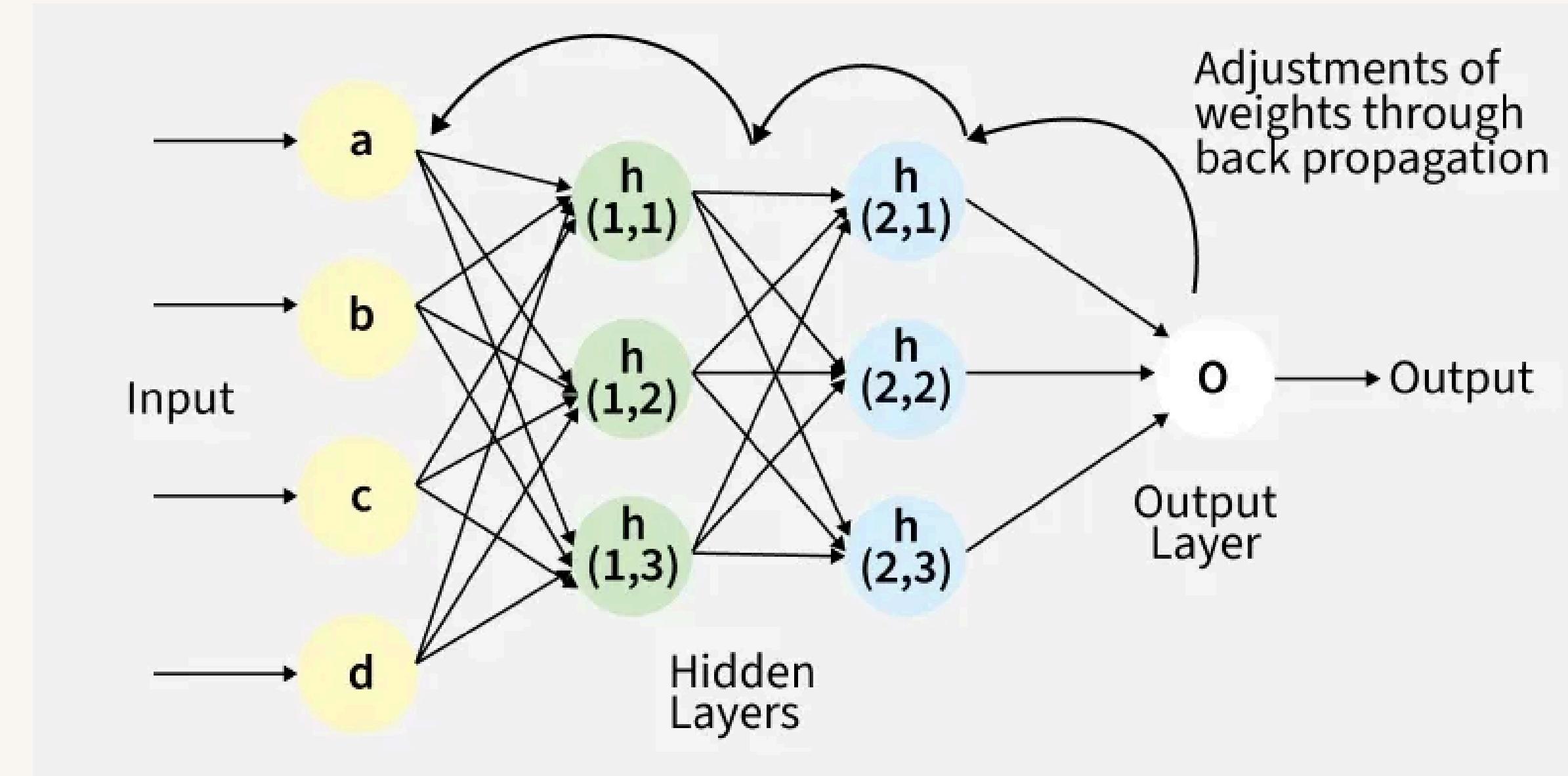
ANN Parts in a Real-World Scenario

- Problem: Predicting if a house is "expensive" or "not expensive."
- Inputs: Area (sq. ft.), Number of Bedrooms, Location Quality (1-10).
- Weights: How important is each input? Location might have a higher weight than the number of bedrooms.
- Hidden Layers: Neurons here might learn abstract concepts on their own, like "family-sized" (high bedrooms, large area) or "prime real estate" (high location quality).
- Bias: A baseline assumption, perhaps that houses are generally "not expensive."
- Output: A single neuron with a Sigmoid activation, giving a probability between 0 (not expensive) and 1 (expensive).

How an ANN Works: Forward & Backward Propagation



How an ANN Works: Forward & Backward Propagation



How an ANN Works: Forward & Backward Propagation



- Forward Propagation
- Calculate the Loss : The number of wrong answers is your Loss or "error."
- Backward Propagation: This is the most important part of learning. You go back through the test and figure out why you made mistakes. You assign blame: 'I got this question wrong because I misunderstood Topic A and forgot the formula for Topic B.' This process of assigning blame to your earlier steps is Backpropagation.
- Update Weights (Gradient Descent): Based on your mistakes, you adjust your study plan. You decide to study Topic A and B more. This is the network updating its weights to perform better on the next practice test.

This cycle of Forward -> Check Loss -> Backward -> Update is repeated thousands of times until the network becomes very accurate."

Activation Functions

- Simple graphs of the three main types:
 - Sigmoid: Squashes numbers to a range between 0 and 1 (good for output probabilities).
 - Tanh: Squashes numbers to a range between -1 and 1.
 - ReLU (Rectified Linear Unit): The most popular choice for hidden layers. It's simple: if the input is negative, the output is 0; otherwise, the output is the input.

Loss Functions

- A function that measures the "error" or "how wrong" the network's predictions are.
- Goal: The goal of training is to minimize the loss.
- Common Types:
 - Mean Squared Error (MSE): For regression tasks (predicting a number, like a house price).
 - Binary Cross-Entropy: For binary classification tasks (predicting one of two classes, like "cat" or "dog").
 - Categorical Cross-Entropy: For multi-class classification tasks (predicting one of many classes, like handwritten digits 0-9).

Applications of Deep Learning

- A visually appealing slide with icons and examples.
- Computer Vision: Image recognition (tagging photos on Facebook), self-driving cars, medical imaging analysis.
- Natural Language Processing (NLP): Language translation (Google Translate), voice assistants (Siri, Alexa), sentiment analysis.
- Recommendation Engines: Recommending movies on Netflix or products on Amazon.
- Generative AI: Creating new images (DALL-E) or text (ChatGPT).

Types of Deep Learning Networks

- ANN (Artificial Neural Networks): The standard network we've discussed. Great for structured data in tables.
- CNN (Convolutional Neural Networks): The master of images. They are specially designed to recognize patterns, objects, and features in visual data.
- RNN (Recurrent Neural Networks): The specialist for sequential data, like text, speech, or time series. They have a form of 'memory' that allows them to understand order and context."

Q n A

Artificial Intelligence & Machine Learning Certificate Course

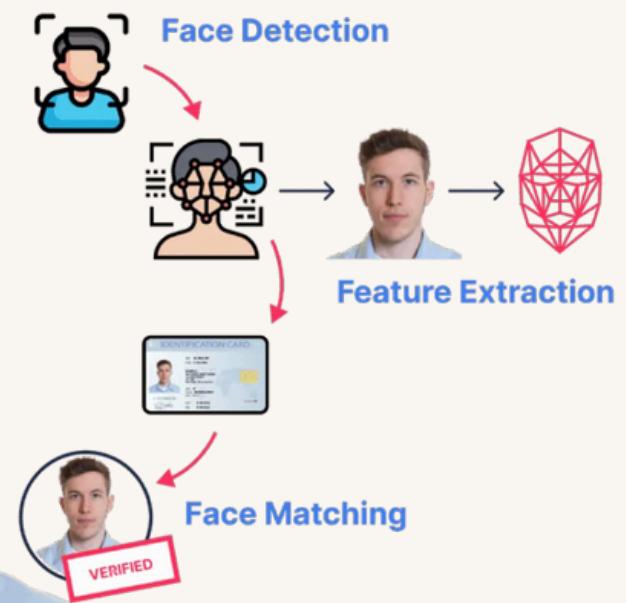
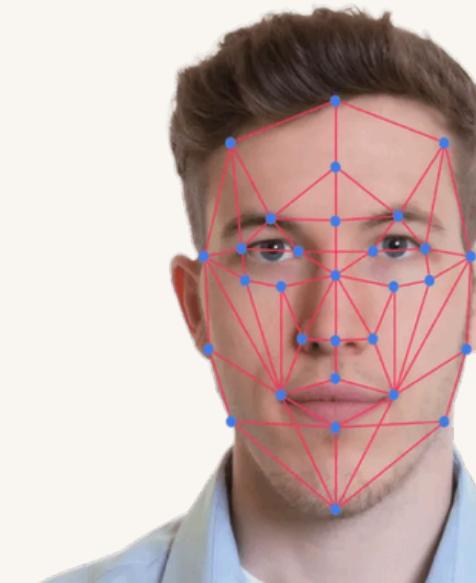
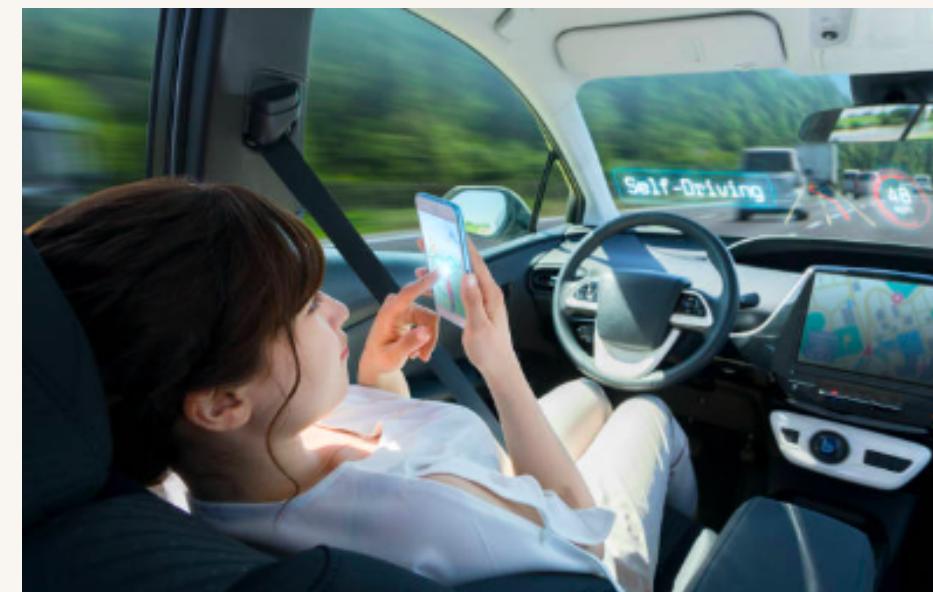
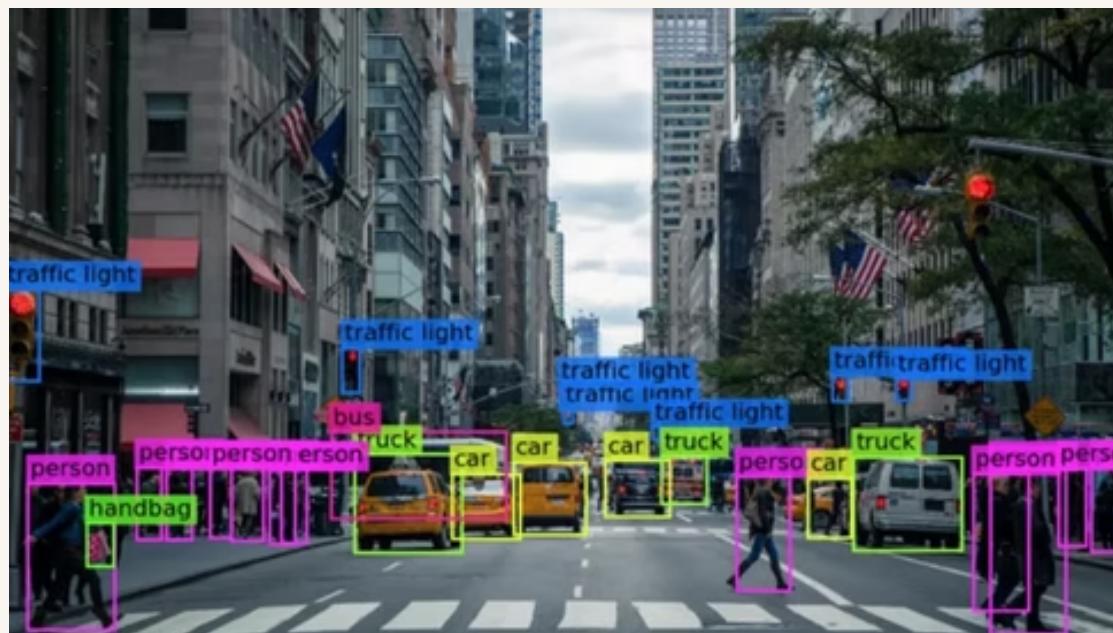
Deep Learning -2 (CNN)

Sehan Arandara
Associate Software Engineer
SE Graduated @ SLIIT
AI ML Instructor

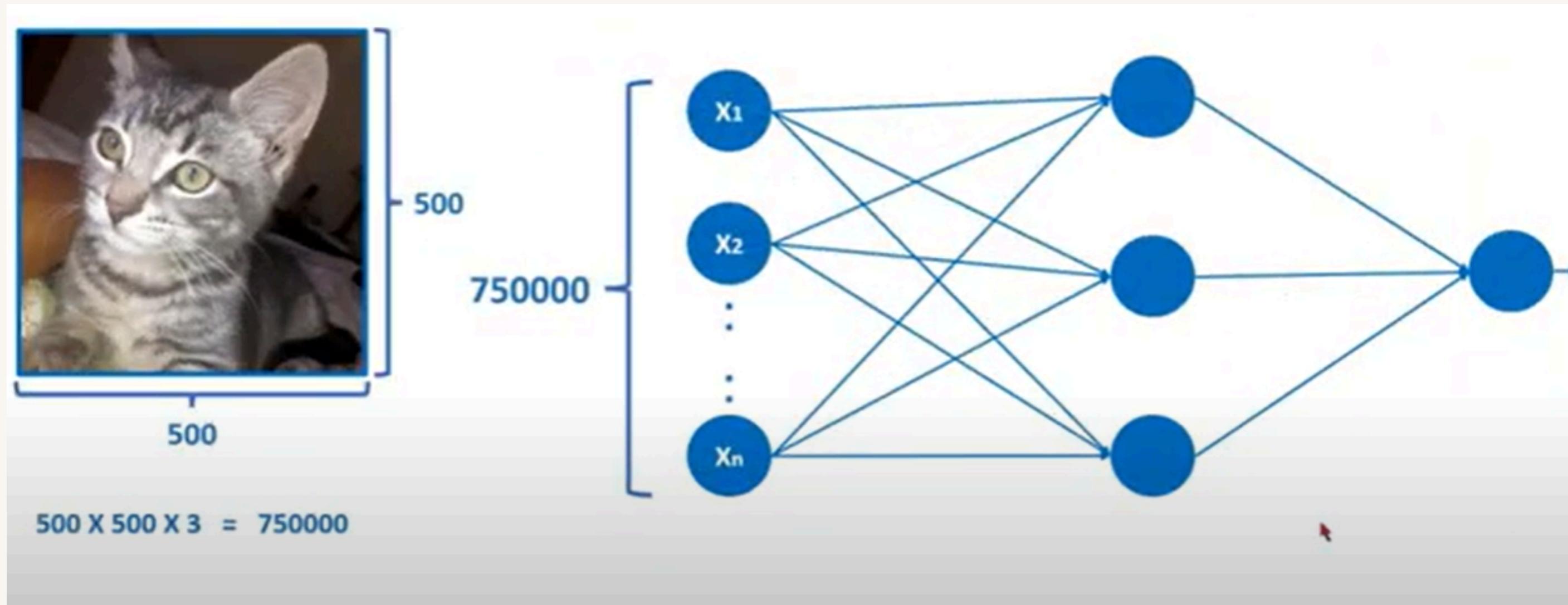
12.10.2025
Day 12

What is Computer Vision?

- A field of AI that trains computers to "see," interpret, and understand information from digital images and videos.
- To enable machines to perform tasks that the human visual system can do.



How a Computer Sees: Images are Just Numbers!



How a Computer Sees: Images are Just Numbers!



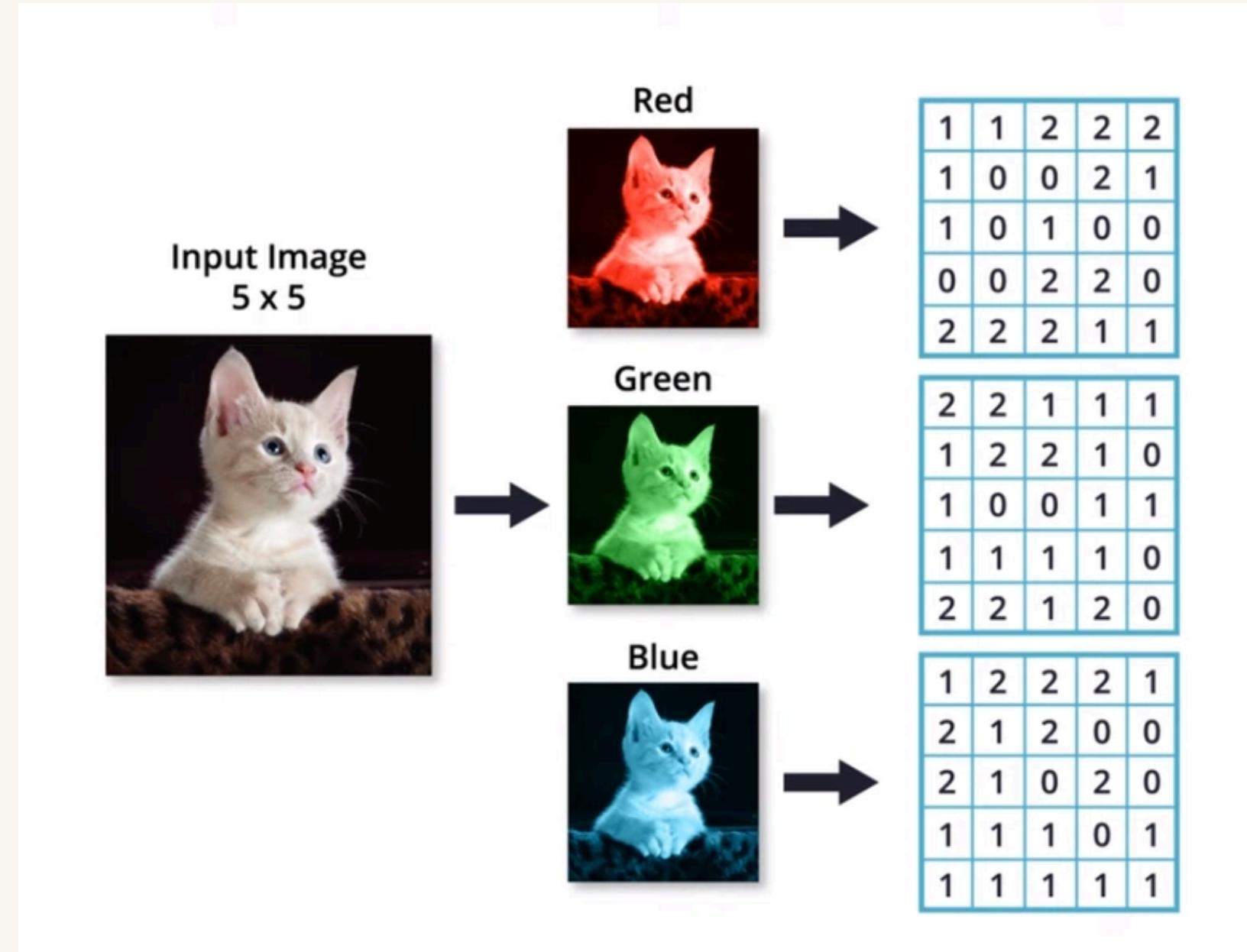
What Computer Sees

0	2	15	0	0	11	19	0	0	0	0	9	9	0	0	0	
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	115	238	255	244	245	243	250	249	255	222	185	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	123	215	251	238	255	49	
13	217	243	255	155	32	228	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0	
0	13	115	255	255	245	255	182	181	248	252	242	206	36	0	19	
1	0	6	115	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	4	86	251	255	246	254	253	255	120	11	0	1	0	
0	0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0	
0	111	255	242	255	152	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3	
0	173	255	255	161	9	20	0	13	3	13	182	251	245	61	0	
0	100	251	241	255	230	68	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	255	249	255	240	255	127	0	5	
0	0	23	115	215	255	250	248	255	255	248	248	118	34	12	0	
0	0	0	6	1	0	82	153	233	255	252	147	37	0	0	1	
0	0	0	5	5	0	0	0	0	14	1	0	6	6	0	0	

0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	6	117	251	255	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	111	255	242	255	158	24	0	0	4	39	255	232	230	56	0
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0
0	0	6	1	0	52	153	233	254	252	147	37	0	0	4	1
0	0	5	5	0	0	0	0	0	14	1	0	6	6	0	0

A picture of a simple grayscale image (e.g., a handwritten '8') on the left. On the right, show a grid of numbers (pixels) representing that same image, with 0 for black, 255 for white, and shades of gray in between.

How a Computer Sees: Images are Just Numbers!



It's just three grids stacked on top of each other one for Red, one for Green, and one for Blue (RGB). By combining the intensity values from these three grids for each pixel, we can create any color.

Key Tasks in Computer Vision

- Image classification
 - It involves analyzing an image and assigning it a specific label or category based on its content such as identifying whether an image contains a cat, dog or car.
- Object Detection
 - It involves identifying and locating objects within an image by drawing bounding boxes around them.
- Image Segmentation
 - It involves partitioning an image into distinct regions or segments to identify objects or boundaries at a pixel level.



What is a Convolutional Neural Network (CNN)?

- A specialized type of deep neural network designed specifically for processing grid-like data, such as images.
- It uses a special mathematical operation called a convolution to learn features from the input image.
- It's a neural network that has been given "eyes" to see patterns in images.
- A Convolutional Neural Network, or CNN, is a special kind of neural network that is the master of understanding images. While a standard ANN treats all inputs the same, a CNN is built with the understanding that in an image, pixels that are close to each other are related. It's designed from the ground up to find patterns in this spatial grid of numbers

Why Use CNNs Instead of ANNs for Images?

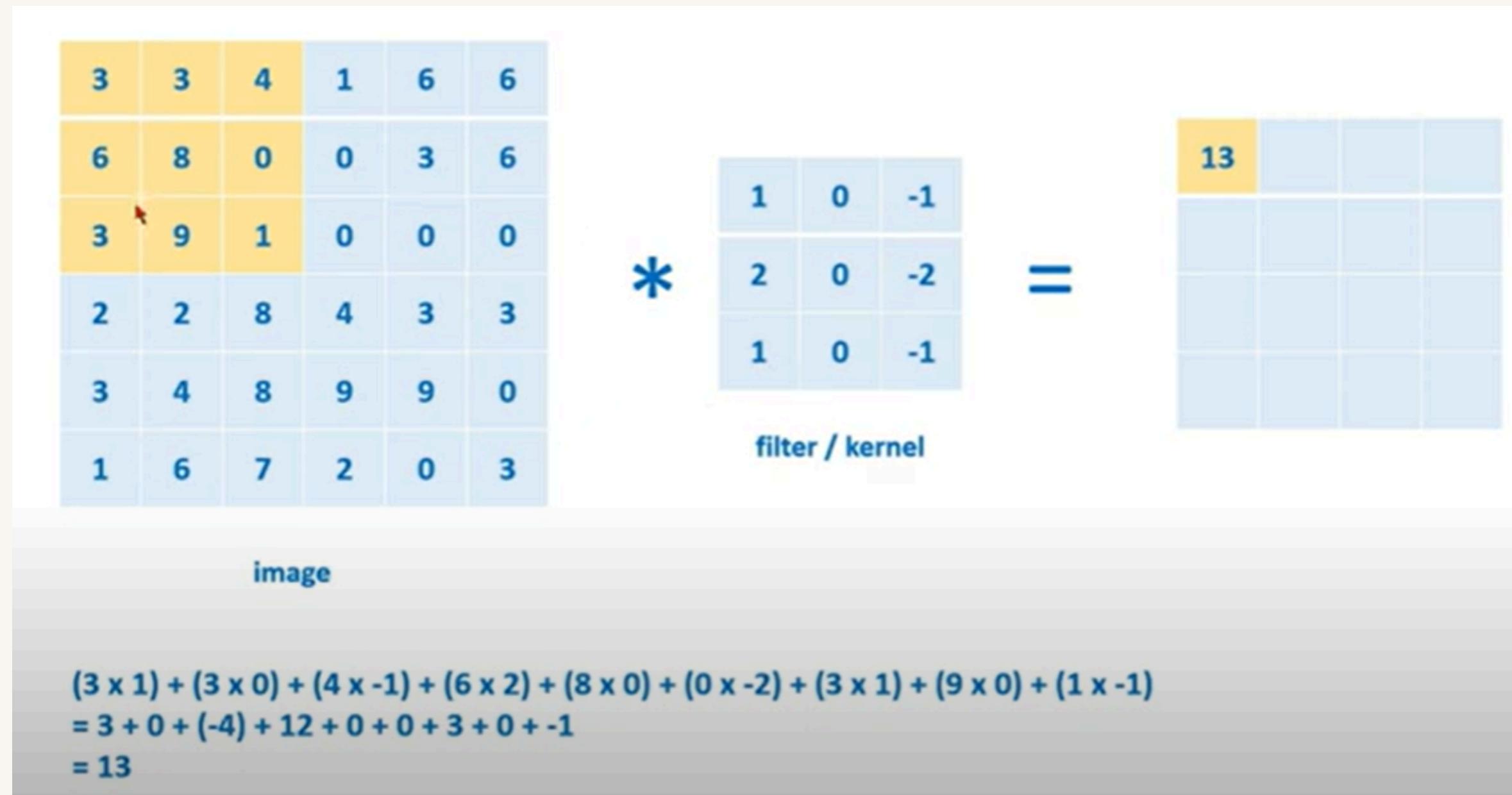
- The Problem with Using a Simple ANN for Images
 - Too Many Parameters (The Curse of Dimensionality)
 - A tiny 32x32 color image has $32 * 32 * 3 = 3,072$ input features.
 - If the first hidden layer has just 100 neurons, that's $3,072 * 100 = 307,200$ weights for the first layer alone!
 - This is computationally expensive and leads to overfitting.
 - Loss of Spatial Information : It loses all information about the spatial structure. It doesn't know which pixels were next to each other

The Core of a CNN: The Convolution Operator

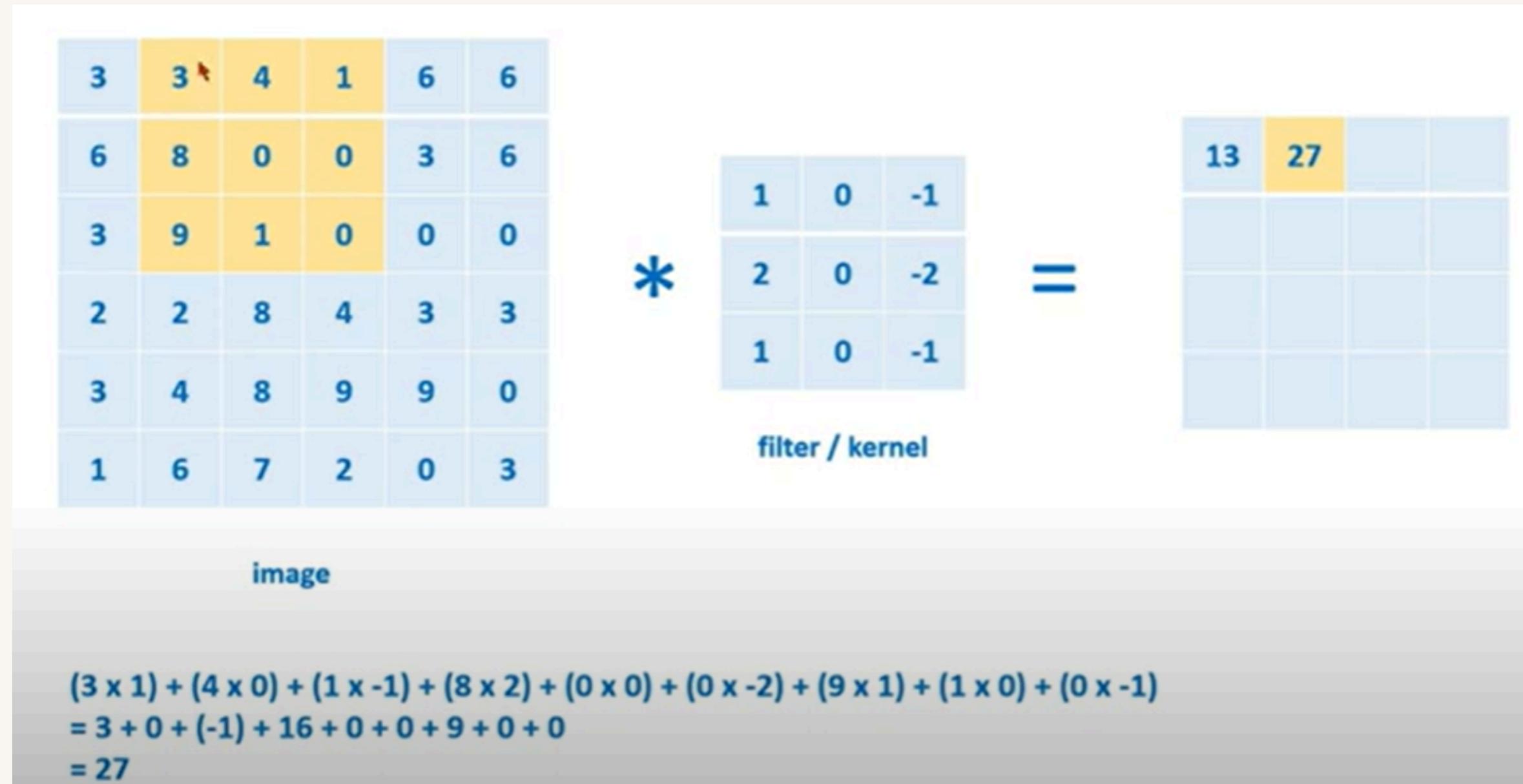
The Convolution Operator: A Pattern Detector

- Key Components:
 - Kernel / Filter: A small matrix of weights (e.g., 3x3).
 - Feature Map: The output of the convolution.
- The kernel is like a small magnifying glass that slides over the image, looking for a specific pattern (like an edge, a curve, or a color).

The Convolution Operator: A Pattern Detector



The Convolution Operator: A Pattern Detector



The Convolution Operator: A Pattern Detector

$$\begin{array}{c}
 \begin{matrix}
 3 & 3 & 4 & 1 & 6 & 6 \\
 6 & 8 & 0 & 0 & 3 & 6 \\
 3 & 9 & 1 & 0 & 0 & 0 \\
 2 & 2 & 8 & 4 & 3 & 3 \\
 3 & 4 & 8 & 9 & 9 & 0 \\
 1 & 6 & 7 & 2 & 0 & 3
 \end{matrix} \\
 \text{*} \quad \begin{matrix}
 1 & 0 & -1 \\
 2 & 0 & -2 \\
 1 & 0 & -1
 \end{matrix} \\
 = \quad \begin{matrix}
 13 & 27 & -7 & -17 \\
 4 & 24 & 4 & -5 \\
 -15 & 0 & 10 & 11 \\
 -22 & -8 & 10 & 18
 \end{matrix}
 \end{array}$$

6×6
 $n \times n$
 3×3
 $f \times f$
 4×4
 $(n-f+1) \times (n-f+1)$

Padding, Stride, Pooling & Augmentation

Padding:

- What: Adding extra pixels (usually zeros) around the border of the input image.
- Why: 1. Prevents the feature map from shrinking too quickly. 2. Allows the kernel to focus more on the pixels at the edges.
- Visual: An image grid with a border of zeros added around it.

Padding, Stride, Pooling & Augmentation

Padding, Stride, Pooling & Augmentation

$$\begin{array}{|c|c|c|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline
 0 & 3 & 3 & 4 & 1 & 6 & 6 & 0 & \\ \hline
 0 & 6 & 8 & 0 & 0 & 3 & 6 & 0 & \\ \hline
 0 & 3 & 9 & 1 & 0 & 0 & 0 & 0 & \\ \hline
 0 & 2 & 2 & 8 & 4 & 3 & 3 & 0 & \\ \hline
 0 & 3 & 4 & 8 & 9 & 9 & 0 & 0 & \\ \hline
 0 & 1 & 6 & 7 & 2 & 0 & 3 & 0 & \\ \hline
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \\ \hline
 \end{array}
 \times
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 2 & 0 & -2 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}
 =
 \begin{array}{|c|} \hline
 -14 \\ \hline
 \end{array}$$

$$\begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 3 & 3 & 4 & 1 & 6 & 6 & 0 & \\ \hline 0 & 6 & 8 & 0 & 0 & 3 & 6 & 0 & \\ \hline 0 & 3 & 9 & 1 & 0 & 0 & 0 & 0 & \\ \hline 0 & 2 & 2 & 8 & 4 & 3 & 3 & 0 & \\ \hline 0 & 3 & 4 & 8 & 9 & 9 & 0 & 0 & \\ \hline 0 & 1 & 6 & 7 & 2 & 0 & 3 & 0 & \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline -14 & 4 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|c|c|c|c|c|c|} \hline & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline & 0 & 3 & 3 & 4 & 1 & 6 & 6 & 0 \\ \hline & 0 & 6 & 8 & 0 & 0 & 3 & 6 & 0 \\ \hline & 0 & 3 & 9 & 1 & 0 & 0 & 0 & 0 \\ \hline & 0 & 2 & 2 & 8 & 4 & 3 & 3 & 0 \\ \hline & 0 & 3 & 4 & 8 & 9 & 9 & 0 & 0 \\ \hline & 0 & 1 & 6 & 7 & 2 & 0 & 3 & 0 \\ \hline & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} \quad * \quad
 \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline -14 & 4 & 12 & -7 & -16 & 15 \\ \hline -28 & 13 & 27 & -7 & -17 & 12 \\ \hline -28 & 4 & 24 & 4 & 5 & 6 \\ \hline -17 & -15 & 0 & 10 & 11 & 15 \\ \hline -16 & -22 & -8 & 10 & 18 & 21 \\ \hline -16 & -17 & 3 & 13 & 7 & 9 \\ \hline \end{array}$$

3×3
 $f \times f$

6×6

$(n+2p-f+1) \times (n+2p-f+1)$

Padding, Stride, Pooling & Augmentation

Stride:

- What: The number of pixels the kernel moves at a time as it slides across the image.
- Why: A larger stride (e.g., 2) makes the model smaller and faster by producing a smaller feature map.
- Visual: An animation showing a kernel jumping 2 pixels at a time.

Padding, Stride, Pooling & Augmentation

3	3	4	1	6
6	8	0	0	3
3	9	1	0	0
2	2	8	4	3
3	4	8	9	9

image

$*$ filter / kernel =

1	0	-1
2	0	-2
1	0	-1

3	3	4	1	6
6	8	0	0	3
3	9	1	0	0
2	2	8	4	3
3	4	8	9	9

image

$*$ filter / kernel =

1	0	-1
2	0	-2
1	0	-1

stride = 1

Padding, Stride, Pooling & Augmentation

3	3	4	1	6
6	8	0	0	3
3	9	1	0	0
2	2	8	4	3
3	4	8	9	9

image

$*$ filter / kernel =

1	0	-1
2	0	-2
1	0	-1

3	3	4	1	6
6	8	0	0	3
3	9	1	0	0
2	2	8	4	3
3	4	8	9	9

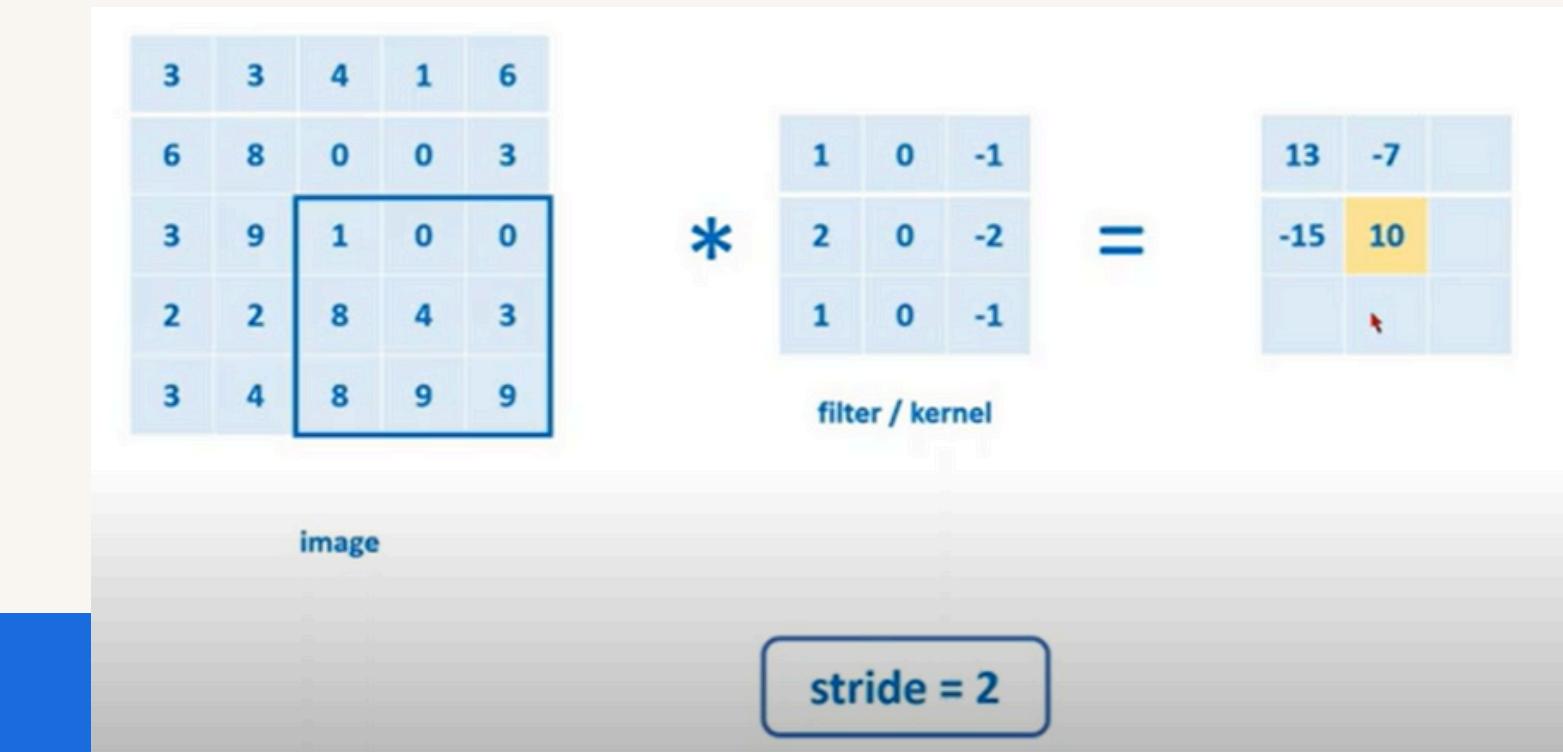
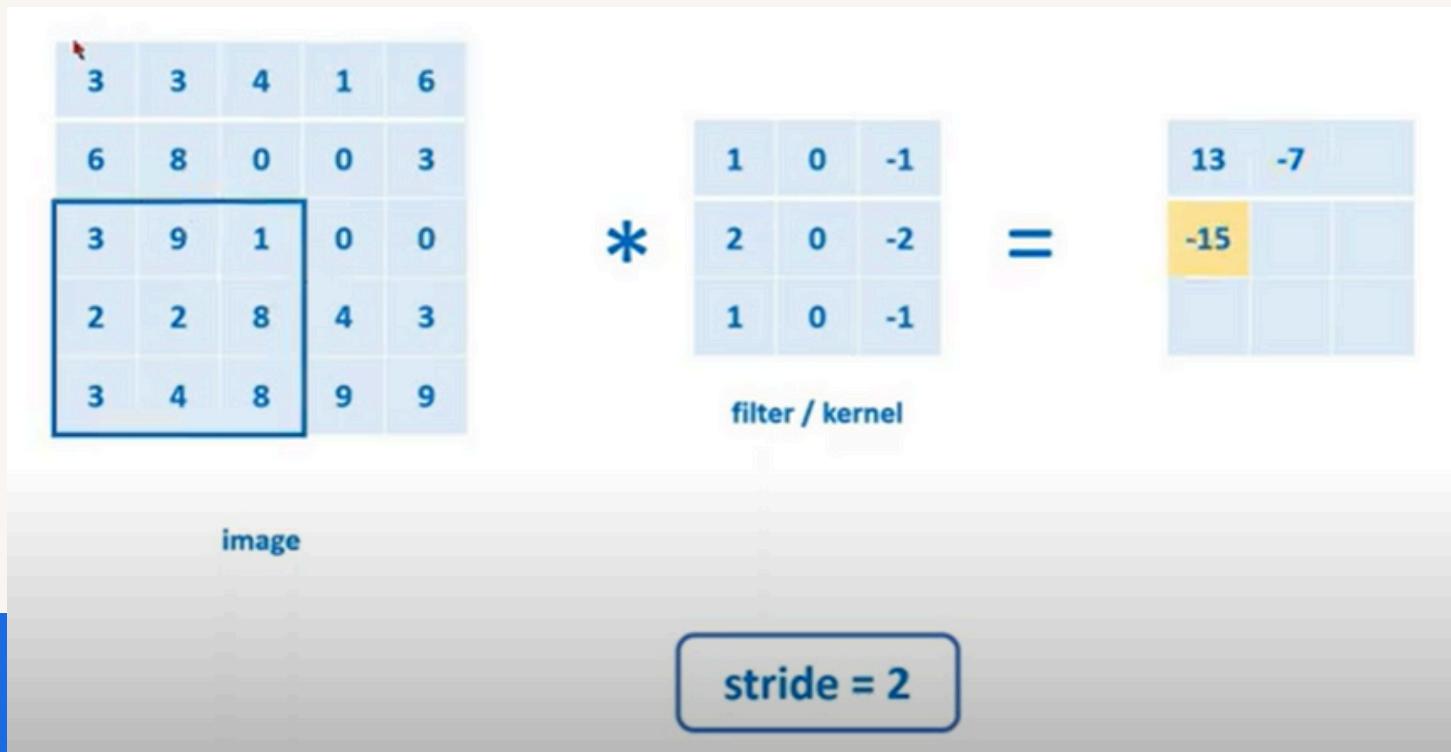
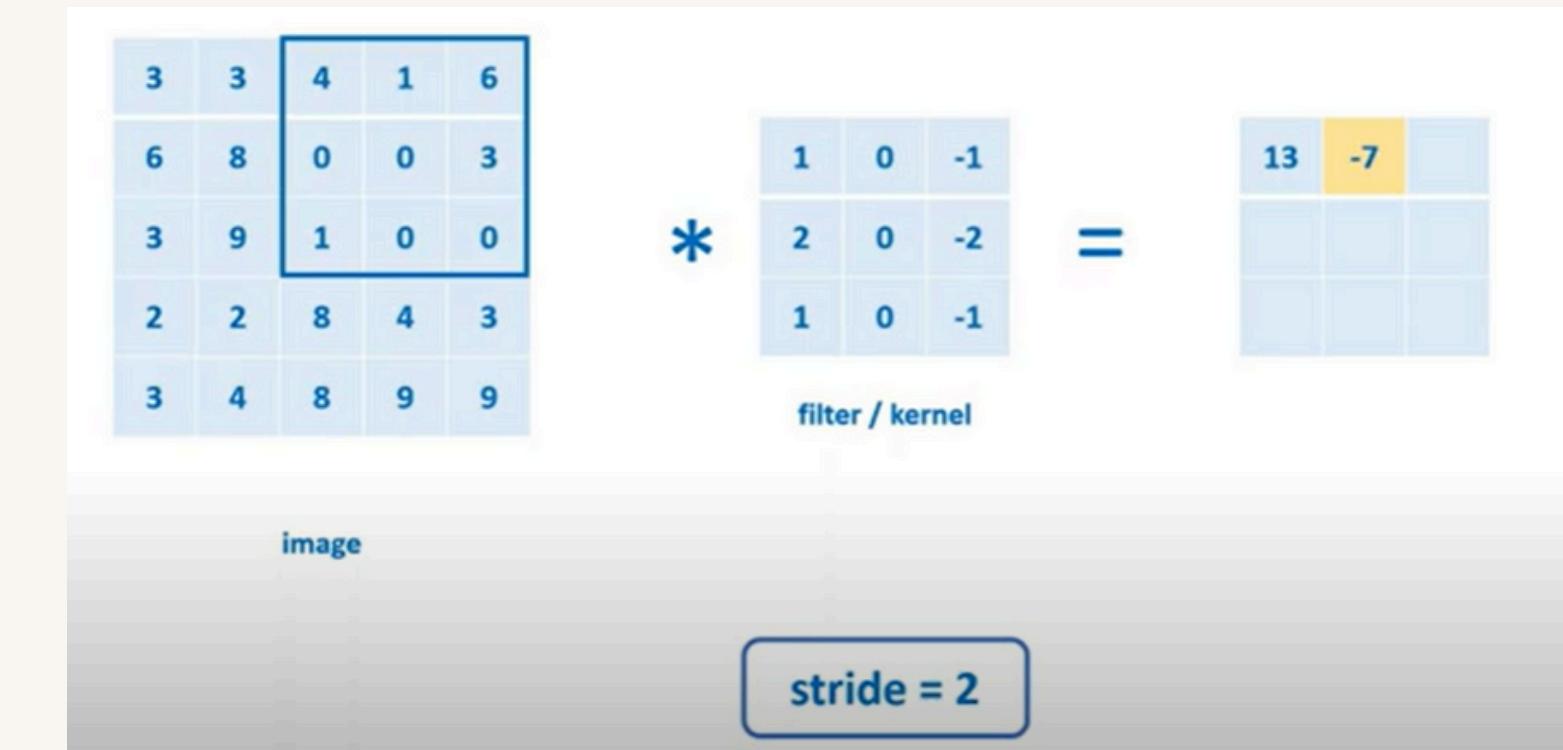
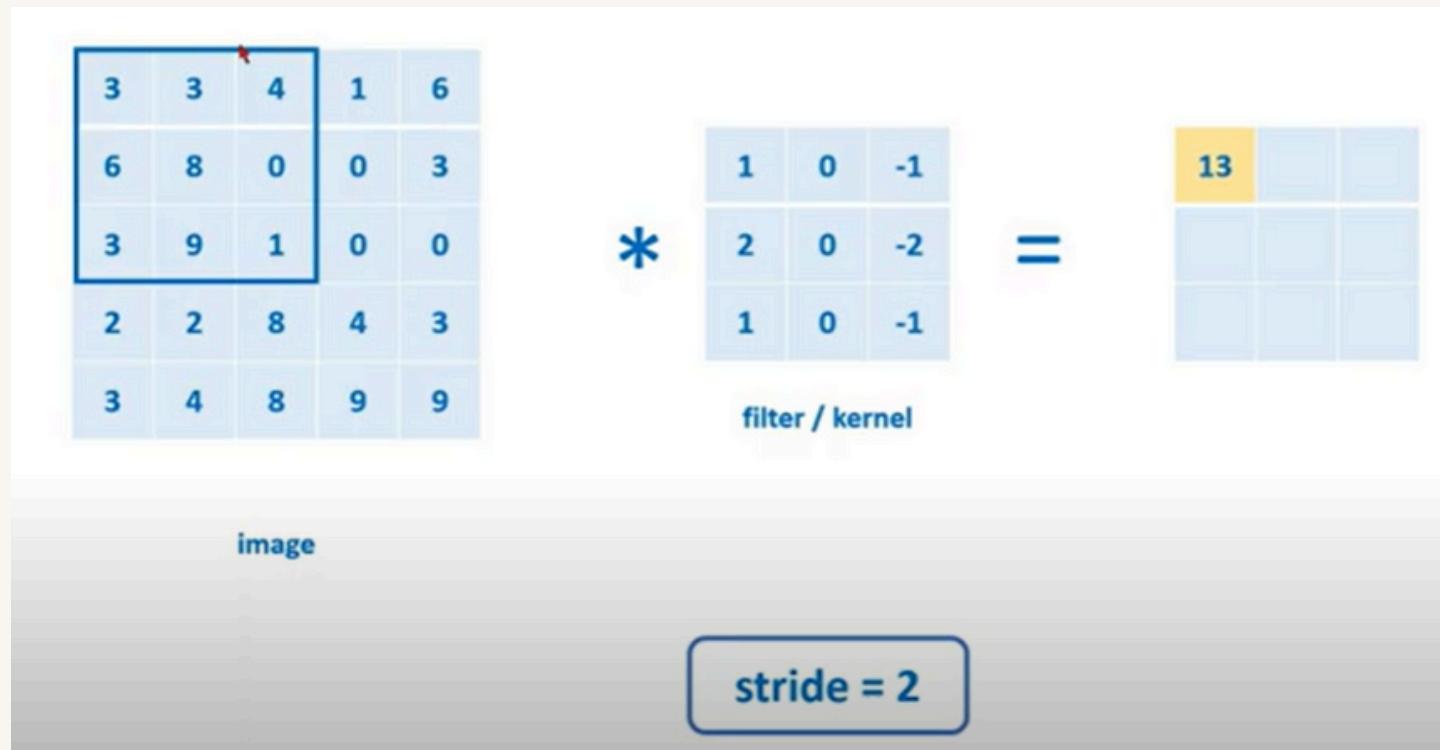
image

$*$ filter / kernel =

1	0	-1
2	0	-2
1	0	-1

stride = 1

Padding, Stride, Pooling & Augmentation



Padding, Stride, Pooling & Augmentation

Pooling Layers (e.g., Max Pooling,Avg Pooling):

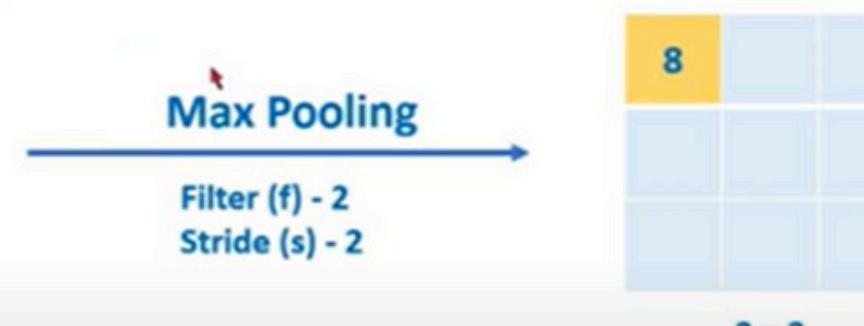
- What: A downsampling operation. It takes a small window (e.g., 2x2) of a feature map and replaces it with a single value (usually the maximum).
- Why: 1. Reduces the size of the feature maps, making the network faster. 2. Makes the model more robust by capturing the most important feature in a region, regardless of its exact position.
- Visual: A 4x4 grid being reduced to a 2x2 grid by taking the max value from each 2x2 block.

Padding, Stride, Pooling & Augmentation

Pooling Layer : Max Pooling

3	3	4	1	6	6
6	8	0	0	3	6
3	9	1	0	0	0
2	2	8	4	3	3
3	4	8	9	9	0
1	6	7	2	0	3

6×6



$$\left[\frac{n+2p-f+1}{s} \right] \times \left[\frac{n+2p-f+1}{s} \right]$$

$$\left[\frac{6+0-2+1}{2} \right] \times \left[\frac{6+0-2+1}{2} \right]$$

Pooling Layer : Max Pooling

3	3	4	1	6	6
6	8	0	0	3	6
3	9	1	0	0	0
2	2	8	4	3	3
3	4	8	9	9	0
1	6	7	2	0	3

6×6

Max Pooling
 Filter (f) - 2
 Stride (s) - 2

$$\left[\frac{n+2p-f+1}{s} \right] \times \left[\frac{n+2p-f+1}{s} \right]$$

$$\left[\frac{6+0-2+1}{2} \right] \times \left[\frac{6+0-2+1}{2} \right]$$

Pooling Layer : Max Pooling

3	3	4	1	6	6
6	8	0	0	3	6
3	9	1	0	0	0
2	2	8	4	3	3
3	4	8	9	9	0
1	6	7	2	0	3

6×6

Max Pooling
 Filter (f) - 2
 Stride (s) - 2

$$\left[\frac{n+2p-f+1}{s} \right] \times \left[\frac{n+2p-f+1}{s} \right]$$

$$\left[\frac{6+0-2+1}{2} \right] \times \left[\frac{6+0-2+1}{2} \right]$$

Padding, Stride, Pooling & Augmentation

Data Augmentation:

- What: Artificially creating more training data by applying random transformations to your existing images.
- Why: To prevent overfitting and help the model generalize better. Deep Learning needs a LOT of data, and this is a way to create more.
- Visual: A single image of a cat on the left, and multiple transformed versions on the right (rotated, zoomed, flipped horizontally, shifted).

Q n A

Artificial Intelligence & Machine Learning Certificate Course

LLM and Gen AI

Sehan Arandara

Associate Software Engineer

SE Graduated @ SLIIT

AI ML Instructor

23.10.2025

Day 14

What is Generative AI?

- A type of Artificial Intelligence that can create new and original content, such as text, images, music, and code.
- It learns patterns and structures from a massive amount of training data and then uses that knowledge to generate new, similar content.

What is a Large Language Model (LLM)?

- The Engine of Generative AI: The Large Language Model (LLM)
- An LLM is a very large, deep learning model that has been pre-trained on a vast amount of text data.
- Its fundamental job is to predict the next word in a sequence.
- Key Characteristics:
 - Large: They have billions of parameters (weights and biases).
 - Language: They are experts in understanding and generating human language.
 - Model: They are a type of neural network.

Why LLMs are So Powerful Today ?

- The Old Way (RNNs):
 - Had "memory" but struggled with long-term context.
 - Processed words one by one, making them slow.
 - Analogy: Like reading a book by covering up all the words except the one you're on.
You'd quickly forget the beginning of the chapter.
- The New Way (The Transformer Architecture):
 - The Big Idea: It can process all words in a sentence at the same time.
 - Key Innovation: A mechanism called "Self-Attention."

Why LLMs are So Powerful Today ?

- For a long time, we used RNNs for text. They were good, but they had a major problem: they forgot things easily. If you gave them a long paragraph, by the end, they would have forgotten the details from the beginning.
- Then, in 2017, a revolutionary new architecture called the Transformer was invented. It was a game-changer because it could look at an entire sentence at once and figure out which words were most important to each other. This is the architecture that powers all modern LLMs like Gemini, GPT, and LLaMA.

Key capabilities of LLM

- Text generation: Creating human-like text for articles, code, or creative writing.
- Translation: Translating text from one language to another.
- Summarization: Condensing large documents into shorter summaries.
- Question answering: Providing answers to questions based on its training data.
- Text completion: Completing sentences or paragraphs based on a given prompt.

The Magic of the Transformer: Self-Attention

- For a long time, we used RNNs for text. They were good, but they had a major problem: they forgot things easily. If you gave them a long paragraph, by the end, they would have forgotten the details from the beginning.
- Then, in 2017, a revolutionary new architecture called the Transformer was invented. It was a game-changer because it could look at an entire sentence at once and figure out which words were most important to each other. This is the architecture that powers all modern LLMs like Gemini, GPT, and LLaMA.

The Core Idea of the Transformer

- The Transformer architecture, introduced in a 2017 paper called "Attention Is All You Need," proposed a revolutionary idea: What if we process all the words in a sentence at the same time?
- To do this, it needed a way to understand the relationships between words without reading them in order. It solved this with two key innovations:
 - Self-Attention: To understand context and the relationships between words.
 - Positional Encodings: To understand the original word order.

Self-Attention

- Self-Attention is the most important concept. It's a mechanism that allows the model, for every word it processes, to look at all the other words in the sentence and determine which ones are most important to understanding that word's meaning.
- "The rocket launched successfully, but it was damaged during reentry."
 - To understand what "it" means, your brain instantly scans the sentence for nouns. You pay a lot of attention to "rocket" and very little attention to "reentry" or "launch." You create a new, richer understanding of "it" that is blended with the meaning of "rocket."

Positional Encodings

- This helps transformers to understand the relative or absolute position of tokens which is important for differentiating between words in different positions and capturing the structure of a sentence
- Without positional encoding, transformers would struggle to process sequential data effectively.
- example :
 - "The cat sat on the mat."
 - Before the sentence is fed into the Transformer model it gets tokenized where each word is converted into a token
 -

The Transformer Blueprint: The Encoder-Decoder

- The full Transformer architecture is typically made of two main parts, working together like an expert translation team.
 - The Encoder (The Reader)
 - The Decoder (The Writer)

The Encoder

- The encoder is a fundamental component of the Transformer model, responsible for processing input sequences and generating context-aware representations.
- The Encoder's job is to read and understand the input sentence.
- It uses Self-Attention and Positional Encodings to create a rich, numerical representation that captures the sentence's full meaning and context.

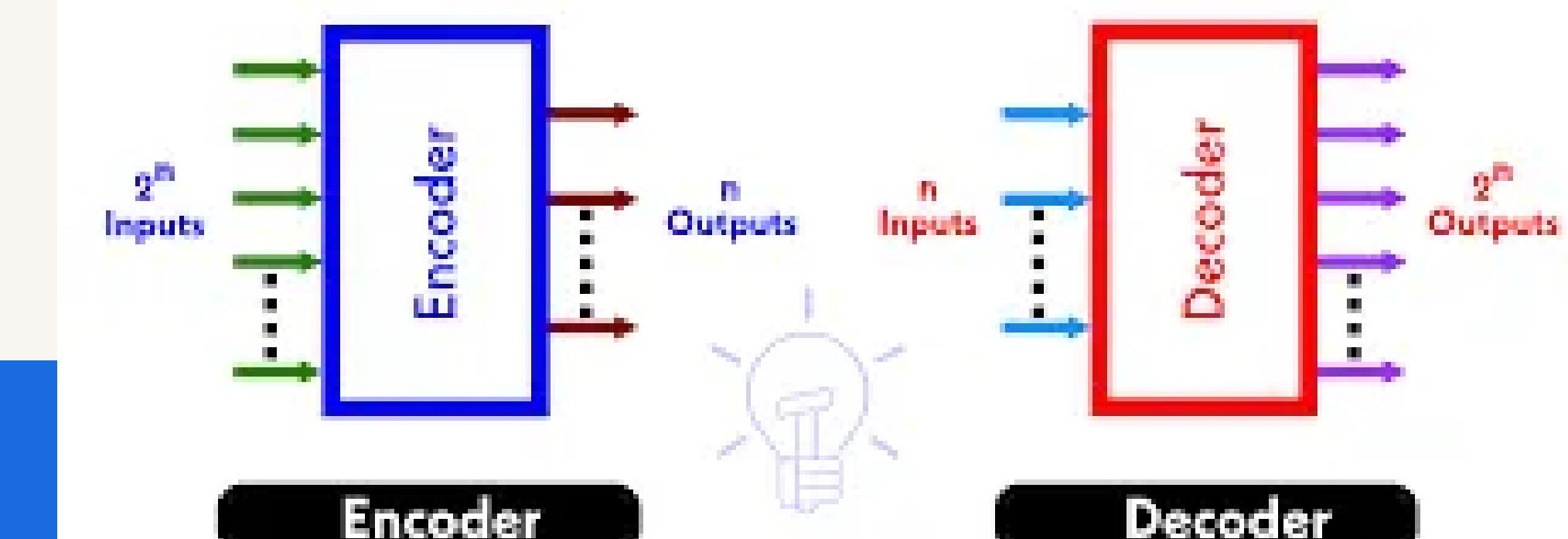
The Decoder

- The decoder is the counterpart to the encoder, designed to generate output sequences from the encoded representations.
- It generates the output one word at a time (e.g., "comment allez-vous").
- Takes the encoded input representation and generates the output sequence, one token at a time.

Encoder vs Decoder

- **Encoder:** Processes and encodes the input sequence into context-aware representations.
- **Decoder:** Generates the output sequence from the encoded representations, ensuring each token is generated based on previous tokens.
- **Encoder:** Uses self-attention to capture relationships within the input sequence.
- **Decoder:** Uses masked self-attention to ensure causal dependency in output generation and encoder-decoder attention to incorporate input context.

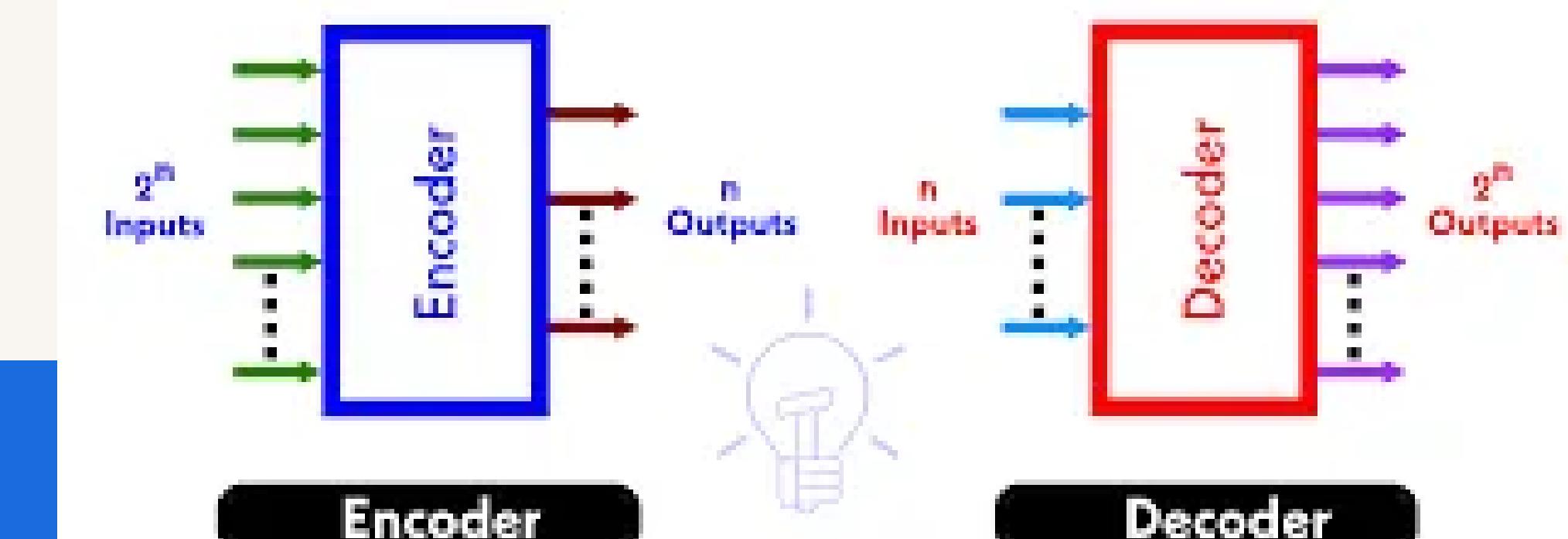
Difference Between Encoder and Decoder



Encoder vs Decoder

- **Encoder:** Processes and encodes the input sequence into context-aware representations.
- **Decoder:** Generates the output sequence from the encoded representations, ensuring each token is generated based on previous tokens.
- **Encoder:** Uses self-attention to capture relationships within the input sequence.
- **Decoder:** Uses masked self-attention to ensure causal dependency in output generation and encoder-decoder attention to incorporate input context.

Difference Between Encoder and Decoder



Large Language Model (LLM) works

- The Input Stage: Understanding Your Prompt
 - Since computers only understand numbers, the first step is to convert your words into a format it can process.
 - Tokenization (Breaking Down Words):
 - The LLM first breaks your sentence down into smaller pieces called "tokens."
 - "What is an LLM?" ---> ["What", "is", "an", "LLM", "?"]
 - Embedding (Creating a "Meaning Map")
 - The model then converts each token into a long list of numbers called a vector or an embedding. This isn't just a random list; the numbers represent the word's meaning, context, and relationships to other words.

Large Language Model (LLM) works

- The "Thinking" Stage: The Transformer Architecture
 - Self-Attention (Finding the Context)
 - Feed-Forward Network (Deep "Thinking")
 - If self-attention is about understanding the relationships between words, the feed-forward network is about "thinking" deeply about the concepts those words represent. It's like a team of experts that processes the context-rich information and prepares it for the next layer or for generating a response.

Large Language Model (LLM) works

- The Output Stage: Generating the Response
 - Predicting the First Word
 - Based on its understanding of your prompt, the model calculates the probability for every word in its vast vocabulary of being the best first word for the answer.
 - Predicting the Next Word
 - Now, the model takes your original prompt plus the first word it just generated and feeds that entire sequence back into the Transformer network. It then predicts the most likely second word
 - repeats this over and over: Prompt -> "The" -> Prompt + "The" -> "cat" -> Prompt + "The cat" -> "sat" -> and so on

Q n A

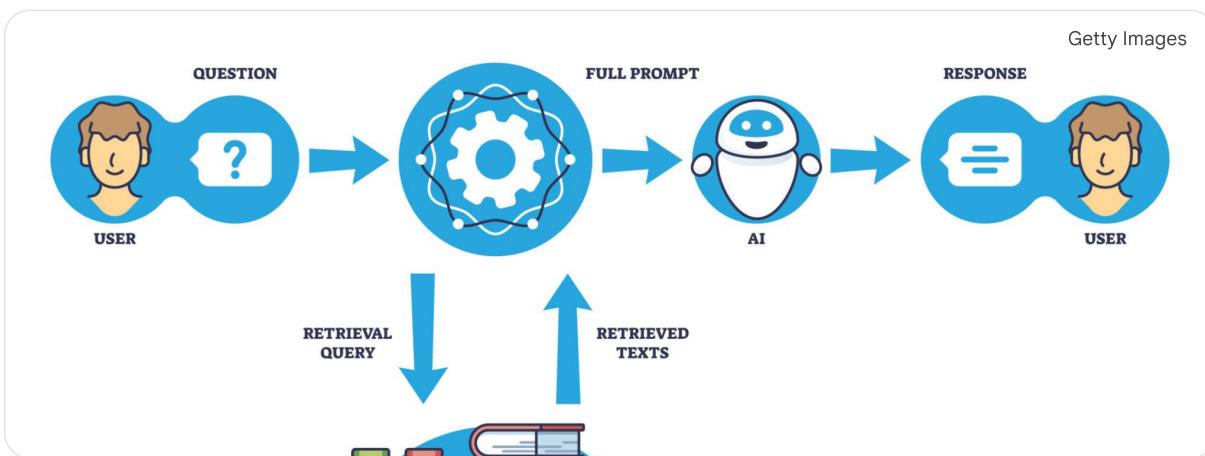
Building a "Knowledge Base" Chatbot with Gemini

This notebook teaches you the fundamentals of **Retrieval-Augmented Generation (RAG)**.

The Goal: To create a chatbot that can answer questions *only* based on a document you provide. This prevents the AI from "making things up" (hallucinating) and limits its knowledge to your specific data.

The Process (RAG):

1. **Retrieve:** When you ask a "question"... we don't send the question straight to the AI. First, we search your document for the most relevant pieces of text (the "context").
2. **Augment:** We take your original "question" and *add* (augment) the "context" we found.
3. **Generate:** We give this combined prompt (Context + Question) to the AI and tell it: "Answer this question *only* using this context."



Cell 1: Setup and Installations

Why this is important: We need to install the Python libraries that let us talk to the Gemini API (`google-generativeai`), read PDF files (`pypdf`), and perform the math for our vector search (`numpy`). The `-q` just makes the output "quiet."

```
!pip install -q google-generativeai pypdf numpy
```

Cell 2: Import Libraries & Configure API Key

Why this is important: This is where we import the libraries we just installed.

The **most important** step is setting your API key. In Colab, the *safest* way to do this is using the "Secrets" manager (look for the  icon on the left-hand side).

1. Click the .
2. Add a new secret named `GEMINI_API_KEY`.
3. Paste your API key as the value.
4. Make sure to toggle "Notebook access" ON.

This code will then securely read that key without you having to paste it directly in the cell.

```

import google.generativeai as genai
import pypdf
import numpy as np
import os
import textwrap # Used for nice text formatting

from google.colab import userdata, files

# --- Get the API Key from Colab Secrets ---
try:
    # This line securely gets your key from the Colab Secrets Manager
    os.environ['GEMINI_API_KEY'] = userdata.get('GEMINI_API_KEY')
except KeyError:
    print("GEMINI_API_KEY not found in Colab Secrets.")
    print("Please go to 'Secrets' (key icon on the left) and add your API key.")
    # This will raise an error and stop execution if the key isn't set
    raise

# --- Configure the API ---
genai.configure(api_key=os.environ['GEMINI_API_KEY'])

print("Successfully configured Gemini API.")

```

Cell 3: Initialize the Gemini Models

Why this is important: This is a key concept. We are using two different AI models for this task:

1. `generation_model` : This is a chat model (`gemini-2.5-flash`). Its job is to *understand language and write answers*.
2. `embedding_model` : This is a specialized model (`text-embedding-004`). Its *only* job is to turn text into a list of numbers (a "vector") that represents its meaning.

```

# Model for generating text (answering questions)
# THIS IS THE "System Instruction" that controls the AI's behavior
SYSTEM_INSTRUCTION = """
You are a helpful assistant. Your task is to answer the user's question based *only* or
- Do not use any external knowledge.
- If the answer is not found in the context, you *must* state: "I'm sorry, but that inf
- Be concise and directly answer the question using only the information from the conte
"""

generation_model = genai.GenerativeModel(
    model_name='gemini-2.5-flash-preview-09-2025',
    system_instruction=SYSTEM_INSTRUCTION # Set the instruction here
)

# Model for creating embeddings (vectors)
embedding_model = genai.GenerativeModel(model_name='models/text-embedding-004')

print("Models initialized.")

```

Cell 4: Upload Your PDF

Why this is important: This cell uses Colab's built-in file uploader to get your "knowledge base" (the PDF) into the notebook's environment. After you run it, a "Choose Files" button will appear.

```
print("Please upload your PDF file.")
uploaded_files = files.upload()

if not uploaded_files:
    print("\nNo file uploaded. Please run the cell again and select a file.")
else:
    # Save the file name for later
    pdf_file_name = list(uploaded_files.keys())[0]
    print(f"\nSuccessfully uploaded: {pdf_file_name}")
```

Cell 5: Read and "Chunk" the PDF

Why this is important: We can't feed a 100-page PDF to the AI at once (it's too large and inefficient). We must break it into smaller, manageable "chunks."

This code does two things:

1. Uses `pypdf` to read all the text from your uploaded PDF.
2. Splits the text into a list of chunks. We're using a simple method: splitting by blank lines (`\n\n`), which often separate paragraphs.

```
def extract_and_chunk_pdf(file_name):
    print(f"Reading {file_name}...")
    pdf_reader = pypdf.PdfReader(file_name)
    full_text = ""
    for page in pdf_reader.pages:
        full_text += page.extract_text() + "\n\n" # Add newlines between pages

    print("Chunking text...")
    # Simple chunking by paragraphs
    chunks = full_text.split("\n\n")
    # Filter out any very small or empty chunks
    chunks = [chunk for chunk in chunks if len(chunk.strip()) > 50]

    print(f"Created {len(chunks)} text chunks.")
    return chunks

# Run the function on our uploaded file
text_chunks = extract_and_chunk_pdf(pdf_file_name)

# --- Let's look at a chunk ---
print("\n--- Example Chunk ---")
# textwrap.fill formats the text nicely
print(textwrap.fill(text_chunks[0], width=80))
print("-----")
```

Cell 6: Create the "Vector Database" (Embedding)

Why this is important: This is the core "AI" part of the setup. We will loop through every single text chunk and use our `embedding_model` to turn it into a vector.

We store these vectors, along with their original text, in a simple list called `knowledge_base`. This is our in-memory "Vector Database."

- `task_type="RETRIEVAL_DOCUMENT"` : This is an important hint to the model. We're telling it to create an embedding that is optimized for being *found* in a search.

(This cell may take a minute or two if your document is large)

```
def embed_text(text):
    """Embeds a single piece of text."""
    try:
        # Use the embedding model
        result = genai.embed_content(
            model="models/text-embedding-004",
            content=text,
            task_type="RETRIEVAL_DOCUMENT" # Important: optimize for search
        )
        return result['embedding']
    except Exception as e:
        print(f"Error embedding text: {e}")
        return None

# This will be our in-memory "Vector Database"
knowledge_base = []

print("Embedding chunks (this may take a moment)....")

for i, chunk in enumerate(text_chunks):
    # Print progress
    if (i+1) % 10 == 0:
        print(f"Embedding chunk {i+1}/{len(text_chunks)}....")

    embedding = embed_text(chunk)

    if embedding:
        # Store the text and its vector
        knowledge_base.append({
            'text': chunk,
            'embedding': np.array(embedding) # Store as numpy array
        })

print("\nEmbedding complete.")
print(f"Knowledge base has {len(knowledge_base)} embedded chunks.")
```

Cell 7: Create the "Search" Function (Retrieval)

Why this is important: Now we need the "R" (Retrieval) part. This cell defines the functions that will:

1. Take a user's *question* (e.g., "What is...")
2. Embed the question itself (using `task_type="RETRIEVAL_QUERY"`)
3. Use **Cosine Similarity** math to compare the question's vector to *all* the chunk vectors in our `knowledge_base` .
4. Return the `top_k` (e.g., top 3) chunks that are *mathematically closest in meaning* to the question.

```
def embed_query(query):
    """Embeds the user's question (query)."""
```

```

try:
    result = genai.embed_content(
        model="models/text-embedding-004",
        content=query,
        task_type="RETRIEVAL_QUERY" # Note: Different task_type!
    )
    return result['embedding']
except Exception as e:
    print(f"Error embedding query: {e}")
    return None

def find_most_similar_chunks(query, top_k=3):
    """Finds the most similar chunks in the knowledge base."""
    if not knowledge_base:
        return ""

    query_embedding = embed_query(query)
    if query_embedding is None:
        return "Error embedding query."

    query_vec = np.array(query_embedding)

    similarities = []
    for item in knowledge_base:
        # Calculate Cosine Similarity
        dot_product = np.dot(query_vec, item['embedding'])
        norm_query = np.linalg.norm(query_vec)
        norm_item = np.linalg.norm(item['embedding'])

        if norm_query == 0 or norm_item == 0:
            similarity = 0.0
        else:
            similarity = dot_product / (norm_query * norm_item)

        similarities.append({'score': similarity, 'text': item['text']})

    # Sort by highest score
    similarities.sort(key=lambda x: x['score'], reverse=True)

    # Return the combined text of the top_k chunks
    context = "\n---\n".join([chunk['text'] for chunk in similarities[:top_k]])
    return context

print("Helper functions for 'Retrieval' are ready.")

```

Cell 8: Create the "Answer" Function (Augmented Generation)

Why this is important: This is the "AG" (Augmented Generation) part. This function assembles the final prompt.

Notice the `system_instruction`. This is *CRITICAL*. It's a special instruction that forces the AI to follow our rules. We are strictly telling it to *only* use the context we provide and to *admit* if the answer isn't there. This is what makes the chatbot "safe" and prevents it from making up facts.

```

def get_chat_response(query, context):
    """Generates a chat response using the Gemini model, based on context."""

    # The system_instruction is now set in Cell 3, when the model was created.
    # We no longer need to define it here.

```

```

# This is the final prompt we send to the model
final_prompt = f"""

Context:
---
{context}
---

Question:
{query}

Answer:
"""

try:
    # Call the generation model
    # We REMOVED the system_instruction argument from this call
    response = generation_model.generate_content(
        final_prompt,
        # Set temperature to 0.0 for factual, non-creative answers
        generation_config=genai.types.GenerationConfig(
            temperature=0.0
        )
    )
    return response.text
except Exception as e:
    print(f"Error generating chat response: {e}")
    return "Sorry, I encountered an error while generating the answer."

print("Helper function for 'Generation' is ready.")

```

Cell 9: Run the Full RAG Pipeline!

Why this is important: This is where we put it all together. This final cell will:

1. Ask you for a question.
2. **Retrieve:** Call `find_most_similar_chunks()` to get the relevant context from your PDF.
3. **Augment & Generate:** Call `get_chat_response()` to generate an answer based *only* on that context.
4. Print the full result!

→ Change the `user_question` variable to ask your own questions!

```

def ask_question(query):
    # 1. Retrieve
    print("Retrieving relevant context...")
    # Find the matching chunks
    context = find_most_similar_chunks(query, top_k=3)

    print("\n--- Found Context ---")
    print(textwrap.fill(context, width=80))
    print("-----\n")

    # 2. Generate
    print("Generating answer...")
    # Get the answer
    answer = get_chat_response(query, context)

```

```
print("\n--- Answer ---")
print(answer)
print("-----")

# --- TRY IT OUT! ---
# Ask a question that can be answered by your PDF
user_question = "What is the main topic of this document?" # <--- CHANGE THIS

ask_question(user_question)
```

Cell 10: (Optional) Interactive Chat

Why this is important: This is just a bonus cell that runs a loop, letting you ask questions multiple times without re-running the cell.

```
print("Starting interactive chat. Type 'quit' to exit.")
print("-" * 30)

while True:
    user_question = input("Ask a question: ")

    if user_question.lower() == 'quit':
        print("Goodbye!")
        break

    ask_question(user_question)
    print("\n")
```

Act as an expert AI Data Scientist. Your task is to take the attached dataset and perform a complete end-to-end Machine Learning (ML) pipeline analysis. You must generate a complete, cell-by-cell Google Colab Python notebook. Each code cell must be preceded by a markdown cell explaining the purpose of the code, linking it back to the standard ML pipeline steps.

[USER INPUTS - PLEASE FILL THESE OUT] 1. **DATASET FILE:** [YOUR DATASET FILE] (This is the file I have attached) 2. **TARGET COLUMN:** '[TARGET COLUMN NAME]' (The column we want to predict) 3. **PROBLEM TYPE:** '[PROBLEM TYPE]' (Specify either 'classification' or 'regression') --- **[PIPELINE INSTRUCTIONS - FOLLOW THESE STEPS EXACTLY]** **Notebook Title:** Machine Learning Pipeline for [TARGET COLUMN NAME]

Prediction

Step 1: Setup & Data Collection * Import all necessary libraries (pandas, numpy, scikit-learn, matplotlib, seaborn). * Load the attached dataset into a pandas DataFrame.

* Clearly state the problem objective based on the target column and problem type.

* Display the first 5 rows of the DataFrame (`.head()`).

Step 2: Exploratory Data Analysis (EDA) * Display the DataFrame's info (`.info()`) and statistical summary (`.describe()`).

* Check for and report the sum of missing values for each column.

* Create a visualization for the distribution of the target variable.

* For **regression**, use a histogram with a KDE plot. Comment on its skewness.

* For **classification**, use a count plot. Comment on the class balance.

Step 3: Data Preprocessing & Feature Engineering * Separate the data into features (X) and the target (y).

* **CRITICAL:** If the problem is **regression** and the target variable is skewed, apply a log transformation (`np.log1p`) to `y`.

* Automatically identify categorical and numerical features from X.

* Split the data into training (80%) and testing (20%) sets. Use a `random_state` for reproducibility. For classification, use `stratify`.

* Create a `ColumnTransformer` preprocessing pipeline that:

1. Applies `StandardScaler` to all numerical features.
2. Applies `OneHotEncoder` (with `handle_unknown='ignore'`) to all categorical features.

Step 4: Model Training * Train three different models appropriate for the `PROBLEM_TYPE`. Encapsulate the preprocessor and the model in a scikit-learn `Pipeline` for each.

* **If 'classification':** 1. `LogisticRegression` (as a baseline) 2. `DecisionTreeClassifier` 3. `RandomForestClassifier`

* **If 'regression':** 1. `LinearRegression` (as a baseline) 2. `DecisionTreeRegressor` 3. `RandomForestRegressor`

Step 5: Model Evaluation * For each trained model, make predictions on the test set.

* **If 'regression':** Remember to transform the predictions back to their original scale (`np.expm1`) before evaluation.

* Calculate and display the standard evaluation metrics.

* **If 'classification':** Provide the `Accuracy`, `Precision`, `Recall`, `F1-Score`, and a `Confusion Matrix` for each model.

* **If 'regression':** Provide the `R-squared (R²)`, `Mean Squared Error (MSE)`, and `Root Mean Squared Error (RMSE)` for each model.

Step 6: Model Comparison & Selection * Create a pandas DataFrame that compares the key performance metric (F1-Score for classification, R² for regression) of all three models.

* Explicitly state which model is the best performer and why.

Step 7: Overfitting Analysis * For the best-performing model (likely Random Forest or Decision Tree), create a visualization to check for overfitting.

* **Primary Method:** Generate and display a **Learning Curve** plot. Explain how to interpret the plot in the markdown comments (i.e., the gap between training and validation scores).

* **Secondary Method:** For **regression**, also show a scatter plot of **Actual vs. Predicted values**.

Step 8: Model Saving * Save the complete pipeline of the best-performing model to a file using `joblib`. Name the file `best_model.joblib`.

Step 9: Inference * Create a final cell to demonstrate a real-world prediction.

* Load the saved `best_model.joblib` file.

* Create a new,

sample data point (a dictionary or single-row DataFrame) for prediction. * Create an inference function that takes the new data and the loaded model as input and returns a clear, human-readable prediction. * Execute the function and display the final prediction.

A Quick Guide to Hierarchical Clustering ♠

What is it & Why use it?

- **What:** Hierarchical Clustering is an algorithm that builds a tree-like hierarchy of clusters. Unlike K-Means, it doesn't produce a single set of clusters, but rather a whole range of clusterings at different levels.
 - **Why:** Its biggest advantage is that you **don't need to specify the number of clusters (K) beforehand**. The algorithm reveals the structure of the data at various scales, which can be very insightful.
-

The Two Types: Agglomerative vs. Divisive

There are two main strategies for building the cluster hierarchy:

1. **Agglomerative (Bottom-up)**
 - Starts with each data point as its own individual cluster.
 - It then progressively **merges** the closest pairs of clusters until only one large cluster remains.
 - This is the most common and widely used approach.
 2. **Divisive (Top-down)**
 - Starts with all data points in one single, giant cluster.
 - It then progressively **splits** the clusters into smaller ones until each data point is its own cluster.
 - This is less common because splitting a cluster is computationally very complex.
-

What is a Dendrogram?

A **Dendrogram** is the tree diagram that visualizes the entire hierarchical clustering process. It's the main output of the algorithm.

- The **y-axis** represents the distance or dissimilarity. Longer vertical lines mean that two more distant clusters were merged.
 - The **x-axis** represents the individual data points.
 - You can **choose the number of clusters** by making a horizontal cut across the dendrogram. The number of vertical lines your cut intersects is the number of clusters you get at that distance level.
-

How the Algorithms Work

How Agglomerative Clustering Works (Step-by-Step)

This is the "building up" method:

1. **Start:** Make each data point a single-point cluster. If you have N points, you have N clusters.
2. **Find & Merge:** Find the two **closest** clusters and merge them. Now you have $N-1$ clusters.
3. **Repeat:** Continue finding and merging the closest pairs of clusters.
4. **End:** The process stops when only one cluster (containing all data points) remains.

How Divisive Clustering Works (Step-by-Step)

This is the "breaking down" method:

1. **Start:** Begin with one giant cluster containing all data points.
2. **Find & Split:** Find the best way to split the cluster into two smaller clusters.
3. **Repeat:** Choose one of the new clusters and continue splitting it.
4. **End:** The process stops when every single data point is its own cluster

Choosing where to "cut" the dendrogram

Method 1: The Longest Vertical Line (Visual Method)

This is the most common and intuitive method for interpreting a dendrogram. The idea is to find the cut that represents the largest "jump" in distance, where two very dissimilar groups were merged.

The Logic: A long vertical line on a dendrogram means that two clusters that were far apart have been merged. This is often a "bad" merge, suggesting that the clusters were naturally distinct. We want to choose a number of clusters that exists *before* these distinct groups are forced together.

How to Do It:

1. Look for the longest vertical line in the dendrogram that is not crossed by any other horizontal lines.
2. Draw a horizontal line (your "cut") that goes through this longest vertical line.
3. Count the number of vertical lines your horizontal cut intersects. This number is your optimal number of clusters (K).

In the example above, the longest vertical jump is the blue line on the left. A horizontal cut made there intersects two vertical lines, suggesting that the optimal number of clusters is K=2. If the second-longest line was chosen, the cut would intersect three lines, suggesting K=3 is also a possibility.

Method 2: Business or Domain Knowledge (Practical Method)

Sometimes, the best number of clusters is determined by the problem you're trying to solve. The "best" clustering is the one that is most useful.

The Logic: The data analysis should serve a practical purpose. If your goal is to segment customers for a marketing campaign, and the marketing department only has the budget for three different campaign strategies, then the most useful number of clusters is K=3, regardless of what the dendrogram suggests.

How to Do It:

1. Understand the business goal. What will the clusters be used for?
2. Identify any real-world constraints. Are there limits on budget, resources, or categories?
3. Choose a K that aligns with those constraints.

K-Means Clustering: Mall Customer Segmentation

A Step-by-Step Guide for Students

Project Objective

Welcome! In this hands-on lab, we will act as data scientists for a shopping mall. Our goal is to analyze customer data to discover natural groupings or "segments." This is a core task in **unsupervised learning**, where we find patterns in data without any pre-existing labels. By the end, you'll have a model that can help the mall create smarter, more effective marketing strategies.

Tools We'll Use

- **Python** in a **Google Colab** notebook.
 - **Pandas** library for data manipulation.
 - **Scikit-learn** library for the K-Means algorithm.
 - **Matplotlib** library for data visualization.
-

Step 1: Setting Up Your Workspace and Loading the Data

 **Goal:** Import our necessary tools (libraries) and load the customer dataset into our environment.

 **Why it Matters:** Every data project begins here. We must first load our raw materials (the data) into a workable format (a Pandas DataFrame) and import the libraries that provide the functions we need to analyze and visualize it.

Your Task:

1. **Import Libraries:** Write `import` statements for `pandas` (as `pd`), `KMeans` (from `sklearn.cluster`), and `matplotlib.pyplot` (as `plt`).
 2. **Load Data:** Use the `pd.read_csv()` function to load the `Mall_Customers.csv` file into a DataFrame variable named `df`.
 3. **First Look:** Use the `.head()` method on your `df` to display the first five rows. This confirms your data has been loaded correctly.
-

Step 2: Feature Selection

 **Goal:** Select the specific columns—'Annual Income (k\$)' and 'Spending Score (1-100)’—that we will use for clustering.

 **Why it Matters:** We want to segment customers based on their financial behavior. Choosing only these two features makes our model focused on this specific question and allows us to easily visualize the results on a 2D plot.

Your Task:

1. **Select Columns:** From the `df` DataFrame, select the two columns 'Annual Income (k\$)' and 'Spending Score (1-100)’. **Hint:** To select more than one column, pass a list of names like `df[['column1', 'column2']]`.
 2. **Format for Scikit-learn:** The algorithm requires the data as a numerical array. Add `.values` to the end of your selection to extract the numbers.
 3. **Store the Data:** Save this final array in a variable called `X`.
-

Step 3: Finding the Optimal Number of Clusters (The Elbow Method)

 **Goal:** Determine the best value for 'K' (the number of clusters) using a visual technique called the Elbow Method.

 **Why it Matters:** The K-Means algorithm requires us to specify how many clusters to find. The Elbow Method helps us find the "sweet spot" by plotting the WCSS (Within-Cluster Sum of Squares) for different values of K. We look for the "elbow" in the plot—the point where the graph's steepness decreases significantly.

Your Task:

1. **Initialize:** Create an empty list named `wcss`.
 2. **Loop and Train:** Create a `for` loop that iterates from `i = 1` to `10`. Inside the loop:
 - Create a `KMeans` model with `n_clusters=i`. Set `random_state=42` for consistent results.
 - Fit the model to your data `X` using the `.fit()` method.
 - Append the model's WCSS value (accessible via `.inertia_`) to your `wcss` list.
 3. **Plot the Results:**
 - After the loop, use `plt.plot()` to graph the `range(1, 11)` on the x-axis and your `wcss` list on the y-axis.
 - Add a title and labels for clarity.
 - Use `plt.show()` to display the plot. Observe the plot to find the elbow.
-

Step 4: Building the Final K-Means Model

🎯 **Goal:** Train the definitive K-Means model using the optimal 'K' you discovered in the previous step.

💡 **Why it Matters:** The last step was an experiment. Now, we use the results of that experiment (our optimal K, which is 5) to build the one model that will actually perform the customer segmentation.

Your Task:

1. **Create the Model:** Instantiate a new `KMeans` object. Set `n_clusters` to 5 and `random_state` to 42.
 2. **Fit and Predict:** Use the `.fit_predict()` method on your data `X`. This trains the model and assigns each data point to a cluster.
 3. **Store the Assignments:** Save the returned list of cluster assignments (which will contain numbers from 0 to 4) in a variable named `y_kmeans`.
-

Step 5: Visualizing and Interpreting the Clusters

🎯 **Goal:** Create a scatter plot to visualize the customer segments and interpret what each segment represents.

💡 **Why it Matters:** This is the payoff! A visualization turns abstract numbers into an insightful story. By analyzing the characteristics of each cluster, we can create customer personas (e.g., "Target," "Careful") that are valuable for the business.

Your Task:

1. **Create the Plot:** Use `plt.figure()` to set a good size for your plot.
2. **Plot Each Cluster:** You will call `plt.scatter()` five times, once for each cluster.
 - To plot just the points belonging to cluster 0, use the condition `y_kmeans == 0`. The x-values will be `X[y_kmeans == 0, 0]` and the y-values will be `X[y_kmeans == 0, 1]`. Give this a `label` in the plot call (e.g., 'Standard').
 - Repeat this for clusters 1, 2, 3, and 4, giving each a descriptive label.
3. **Plot the Centroids:** Find the centers of the clusters using the `.cluster_centers_` attribute of your trained model. Plot these points using another `plt.scatter()` call, but make them visually distinct (e.g., use a black 'X' marker).
4. **Finalize the Chart:** Add a title, axis labels, and a legend using `plt.legend()`. Display your final chart with `plt.show()`.

Module: From Notebook to Live Web App with Streamlit & AI

Objective: In this module, you will learn the complete workflow to turn your machine learning model (from a Colab notebook) into a live, interactive web application running on your local computer. We will use Anaconda for our environment setup and an AI assistant (LLM) to help write the code.

Part 1: Setting Up Your Development Environment 🔧

Before you can build anything, you need to prepare your digital workspace.

Task 1.1: Install Anaconda

Anaconda is a free, all-in-one toolkit for data science. It includes Python, a powerful environment manager, and many common libraries.

1. Go to the official download page: [anaconda.com/download](https://www.anaconda.com/download)
2. Download the installer for your operating system (Windows/macOS/Linux).
3. Run the installer and **accept all the default recommended settings**. This is the easiest and safest option.

Task 1.2: Create Your Project Folder

On your computer (e.g., on your Desktop), create a new, empty folder. Give it a clear name, like `my_streamlit_app`.

Task 1.3: Set Up Your Virtual Environment

We will create an isolated "sandbox" for our project's libraries.

1. **Open Anaconda Prompt.** From your computer's Start Menu or Applications, find and open "Anaconda Prompt". It will open a special terminal.
2. **Navigate to Your Project Folder.** Use the `cd` (change directory) command.

- **Tip:** If your folder is on a different drive (like D:), you must first switch to that drive by typing D: and pressing Enter.
- Use quotes if your path has spaces.
-

Bash

```
# First, switch drive if needed
```

```
D:
```

```
# Then, navigate to your folder
```

```
cd "path\to\your\my_streamlit_app"
```

3.

Create the Virtual Environment. Run the following command. We will name our environment

`venv`.

Bash

```
python -m venv venv
```

4.

Activate the Environment. This "turns on" your isolated sandbox.

Bash

```
.\venv\Scripts\activate
```

5. You have succeeded when you see (`venv`) appear at the beginning of your terminal prompt.
-

Part 2: Generating Your App with an AI Assistant

Now, we'll use your Colab notebook to generate the Streamlit code.

Task 2.1: Get Your Files from Colab

You need two files from your Google Colab session:

1. The trained model file: `employee_attrition_model.joblib`. Download this from the Colab file explorer.
2. The notebook file: `Classifications.ipynb`. Download this using `File -> Download -> Download .ipynb`.

Task 2.2: The Universal Best Prompt

Go to your favorite LLM (like Gemini, ChatGPT, etc.) that allows file uploads. Upload your **.ipynb notebook file** and then copy and paste the following prompt.

Plaintext

You are an expert Python developer who specializes in building data science applications with Streamlit.

I have uploaded my Google Colab notebook (`.ipynb`) where I trained and saved a scikit-learn model pipeline named `employee_attrition_model.joblib`.

Based on the uploaded notebook, please perform the following two tasks:

****Task 1: Generate the Streamlit Application Code****

Generate a complete, single-file Streamlit application named `streamlit_app.py`. The app must:

1. Load the `employee_attrition_model.joblib` file safely using `@st.cache_resource` for efficiency.
2. Display a clear title for the application.
3. Create user-friendly input widgets (like `st.selectbox` and `st.number_input`) for all the features the model needs. Use the notebook to determine the correct feature names, categories for dropdowns, and sensible default values.
4. Include a "Predict" button.
5. When the button is clicked, collect the user inputs, create a Pandas DataFrame, and use the loaded pipeline to make a prediction.
6. Display the final prediction (e.g., "Leave" or "Stay") and the confidence probability in a clear, easy-to-read format.
7. The code must be well-commented.

****Task 2: Generate the `requirements.txt` File****

Analyze all the libraries and their versions used in my notebook and in the Streamlit app you just generated. Create the content for a `requirements.txt` file. It is critical that you specify the exact versions (e.g., `scikit-learn==1.6.1`) to ensure the model loads without any version mismatch errors. Please use this version for the `sklearn pip install scikit-learn==1.6.1` do not change this version

Part 3: Building and Running Your App Locally

Task 3.1: Create Your Project Files

1. Place the downloaded **employee_attrition_model.joblib** file inside your **my_streamlit_app** folder.
2. Create a new file in the folder named **streamlit_app.py**. Paste the Python code the LLM gave you for Task 1 into this file.

3. Create another new file named **requirements.txt**. Paste the list of libraries the LLM gave you for Task 2 into this file.

Your folder should now look like this:

```
/my_streamlit_app/  
|-- venv/  
|-- streamlit_app.py  
|-- requirements.txt  
|-- employee_attrition_model.joblib
```

Task 3.2: Install the App's Dependencies

Go back to your Anaconda Prompt, where your (**venv**) is still active. Run this command to install all the libraries your app needs:

Bash
pip install -r requirements.txt

Task 3.3: Run Your Streamlit App!

This is the final step. Run the following command in your terminal:

Bash
streamlit run streamlit_app.py

Your web browser will automatically open a new tab with your live, interactive application running on your local host. Congratulations!

Here are the key hyperparameters and how you can tune them:

1. Chunking

This is arguably the **most important** hyperparameter. Our current method (`full_text.split("\n\n")`) is simple, but not very precise.

- **The Problem:** What if a "paragraph" is 3,000 words long? That's too big to be useful. What if it's 10 words? That's too small and lacks context.
- **How to Tune:** A much better method is "Fixed-Size Chunking with Overlap."
 - **chunk_size:** The number of characters (or tokens) in each chunk (e.g., `1000`).
 - **chunk_overlap:** How many characters to *repeat* at the start of the next chunk (e.g., `100`). This overlap is crucial; it prevents a single idea from being split between two chunks.

2. Retrieval

- **top_k=3:** This is the number of chunks we retrieve to send to the AI.
 - **If you set k=1:** You risk not getting enough context for a complex answer.
 - **If you set k=10:** You risk "drowning" the AI in too much information, and some of it might be irrelevant "noise." This can confuse the model.
- **How to Tune:** Start with `k=3` or `k=5`. If you find answers are missing information, try increasing `k`. If you find answers are "confused" or pulling from irrelevant facts, try decreasing `k`.

3. Generation

- **temperature=0.0:** This controls the AI's "creativity."
 - `0.0` means it will be very factual, deterministic, and "boring"—which is exactly what we want for a RAG chatbot!
 - **If you set temperature=0.7:** The AI might start to "guess" or "infer" information that isn't *explicitly* in the context, which defeats the purpose of our system.
- **How to Tune:** For factual Q&A, you should almost always **leave this at 0.0**.

AI Product Development Capstone Project

1. ⚙️ Project Objective

The goal of this project is to simulate a real-world product development cycle. You will not just train a model; you will build a **complete, end-to-end AI product**. This involves ideation, data processing, building an ML pipeline, creating a backend API to serve your model, and building a frontend user interface to interact with it.

This project will serve as a key portfolio item, demonstrating your ability to build and deploy functional AI applications.

2.💡 Phase 1: Project Ideation & Proposal

You must decide what problem you want to solve. Your idea must be approved by your instructor before you begin development.

Brainstorming Themes:

- **Health & Wellness:** A symptom checker, a food nutrition calculator from an image, a workout pose corrector.
- **E-commerce:** A product recommendation engine, a customer review sentiment analyzer, a fake review detector.
- **Education:** A text summarizer for study notes, a simple "question-answering" bot for a specific textbook.
- **Productivity:** An email spam classifier, a meeting transcript summarizer, an "invoice data extractor" from a PDF.
- **Computer Vision:** A plant disease detector, a "mask on/off" detector, an object counter in an image.

Project Proposal (To be submitted for approval): You must submit a 1-page document outlining:

1. **Problem Statement:** What problem are you solving and who is it for?
2. **Dataset:** Where will you get your data? (e.g., Kaggle, a public API, or will you create it?)
3. **ML Model:** What type of model do you plan to use? (e.g., "A CNN model using TensorFlow," "A text classifier using Scikit-learn").
4. **Proposed Tech Stack:** What tools do you plan to use? (e.g., "Python/Flask for backend, React for frontend").

3.🛠️ Phase 2: Core Technical Requirements

Your final product **must** include these three distinct components that work together.

A. The ML Pipeline (The "Brain")

This is the Python script (e.g., a Jupyter Notebook or `.py` file) where you process data and train your model.

- **Data Ingestion:** Load your dataset.
- **Data Preprocessing:** Clean the data (handle missing values, normalize, etc.). For text or images, this includes tokenization, resizing, etc.
- **Feature Engineering:** Create relevant features for your model.
- **Model Training:** Train one or more models. You should experiment with different parameters.
- **Model Evaluation:** Evaluate your model fairly on a "test set" and record its accuracy, precision, F1-score, or other relevant metrics.
- **Model Serialization:** Save your final, trained model to a file (e.g., `model.pkl`, `model.h5`, or using `joblib`).

B. The Backend (The "Plumbing")

This is a server that exposes your model to the world via a REST API.

- **Technology:** Use a simple backend framework like **Flask** or **FastAPI** (Python).
- **Function:** Create a server that can:
 1. Load your saved model file (`model.pkl`) when the server starts.
 2. Define an API endpoint (e.g., `/predict`).
 3. This endpoint must accept data from the frontend (e.g., as JSON or an image upload).
 4. It should use the loaded model to make a prediction.
 5. It must send the prediction back to the frontend (e.g., as JSON).

C. The Frontend (The "Face")

This is the user interface (UI) that a user interacts with.

- **Technology:** Use a web framework like **React**, **Vue**, or even simple **HTML/CSS/JavaScript**.
- **Function:** Your UI must:
 1. Provide a simple, clean interface for a user to input data (e.g., a text box, a form, an "upload image" button).
 2. Have a "Submit" or "Predict" button.
 3. When clicked, it must send the user's data to your backend API endpoint (`/predict`).
 4. It must wait for the response from the backend.
 5. It must display the prediction clearly to the user (e.g., "This review is: **Positive**" or "Prediction: **This is a 'Rose'**").

4. 🎯 Phase 3: Final Submission Deliverables

You must submit a single link (e.g., to Google Drive or a GitHub repository) containing the following:

1. **Source Code:** All your code, organized into folders (e.g., `/frontend`, `/backend`, `/ml_pipeline`).
2. **README.md File:** A documentation file that explains:
 - A brief description of your project.
 - Instructions on how to run it (e.g., `pip install -r requirements.txt`, `python app.py`).
 - A link to your dataset.
3. **Trained Model File:** The `model.pkl` or `model.h5` file. (If it's too large for Google Drive, include a link to it).
4. **Project Report & Screenshots:** A short report (PDF or Google Doc) that includes:
 - Your project idea and final model performance.
 - **Screenshots of your final product in action.** (This is what you asked about).
 - **Required Screenshots:**
 - A screenshot of your frontend UI before prediction.
 - A screenshot of your frontend UI *after* showing a successful prediction.
 - A screenshot of your ML pipeline's evaluation metrics (e.g., the confusion matrix or accuracy score).
5. **(Optional but Recommended)** A short 2-5 minute video demonstrating your project working.

Understanding Transfer Learning and Hugging Face

1. The Big Problem: "How Do I Build an AI?"

Imagine your job is to create an app that can read any movie review and instantly tell if it's "positive" or "negative." This is called **Sentiment Analysis**.

The Old Way (Building from Scratch): If you tried to do this 10 years ago, you would need:

- A team of Ph.D.s in mathematics and linguistics.
- A massive "dataset" of billions of sentences.
- Millions of dollars to rent supercomputers for months or even years.
- All this, just to teach the AI what a "word" is, what "grammar" is, and what "context" means.

This was impossible for almost everyone.

The New Way (Transfer Learning): Today, you can "download" a giant, pre-trained AI "brain" that *already* understands language, and then just "tweak" it for your specific task. This process is fast, cheap, and can be done by a single person in an afternoon.

This "new way" is called **Transfer Learning**, and the "tweak" is called **Fine-Tuning**. The place where everyone gets these "brains" is **Hugging Face**.

2. Key Concept: What is a "Transformer"?

A **Transformer** is the "brain" we've been talking about. It's a specific, powerful type of AI architecture (you've heard of its famous members: **GPT** and **BERT**).

A Transformer's superpower is **understanding context**.

It doesn't just read word-by-word. It reads a whole sentence at once and understands the relationships between all the words. It knows that the word "bank" means something different in "river bank" vs. "money in the bank."

Analogy: A Transformer is like a "master of language." It understands grammar, nuance, and relationships. It even understands abstract concepts, like the famous example: "King" - "Man" + "Woman" = "Queen".

3. Key Concept: Pre-training vs. Fine-tuning

This is the single most important process in modern AI.

Part 1: Pre-training (The "General Education")

This is the expensive, difficult part that a big company (like Google or OpenAI) does *once*.

- **What it is:** A giant, "blank" Transformer model is fed all of Wikipedia, a library of millions of books, and huge parts of the internet.
- **The Goal:** It is **NOT** taught a specific task. It's just learning the *rules of language*—grammar, context, facts, and how ideas relate.
- **The Analogy:** This is a person getting a Ph.D. in Literature. They graduate as an expert in language. They understand nuance and context, but they don't have a specific *job*. This "Ph.D." is the **pre-trained model**.

Part 2: Fine-tuning (The "Specific Job Training")

This is the fast, easy part that *you* get to do.

- **What it is:** We take that pre-trained "Ph.D." model and give it a *new*, small, specific dataset. For our project, we would give it a list of 10,000 movie reviews and a label for each one ("positive" or "negative").
- **The Goal:** To teach this expert in language a *new, specific job*.
- **The Analogy:** This is the 1-week of company training. We tell the Ph.D. graduate, "Welcome to the job. Forget all that complex literature for a minute. Just read these 10,000 customer reviews and tell me if they are happy or sad."

Because the "brain" *already understands language*, it learns this new, simple task incredibly fast. It doesn't need to re-learn what a "word" is. It just needs to learn that "masterpiece" is positive and "boring" is negative.

4. How It Works (Technically): The "Body" and the "Head"

This is a simple way to visualize *how* fine-tuning works.

- **The "Body":** This is the giant, pre-trained Transformer model (like **BERT** or **DistilBERT**). Its only job is to *read text and understand its meaning*. It "digests" a sentence and outputs a list of numbers (a "vector") that represents the sentence's meaning.
- **The "Head":** This is a *tiny*, new layer that we add on top of the "body." This layer is specific to our task.
 - For our Sentiment Analysis task, we add a "**Classification Head**" with two outputs (0 for "negative", 1 for "positive").
 - If we wanted to translate, we would add a "**Translation Head**."

When we "fine-tune," we are really just training this *tiny new head*. We "freeze" the giant "body" and tell it to just pass its "understanding" to the new "head." This is why it's so fast and can be done on a normal computer (or in a free Google Colab notebook!).

5. What is Hugging Face? (The "Ecosystem")

So, where do you get these models, datasets, and tools? That's **Hugging Face**.

Hugging Face is "**The GitHub for AI.**" It's a central community and platform built around these three key things:

1. **Models (The "AI Library"):** A massive library of over 1 million pre-trained and fine-tuned models. You can find "bodies" (like `bert-base-uncased`) or fully-formed, fine-tuned models (like `distilbert-base-uncased-finetuned-sst-2-english`).
2. **Datasets (The "Textbooks"):** A library of thousands of datasets, all pre-formatted and ready to be used for fine-tuning (like the `imdb` dataset).
3. **Spaces (The "Demo Showroom"):** An "app store" or "science fair" where people build and host live demos of their AI models. It's where you go to *play* with all this technology.

6. Putting It All Together: The "Live" Tour

As you look at the `huggingface.co` website with your teacher, here is what to look for:

- **On the "Models" Page:**
 - Notice the "**Tasks**" filter on the left. This is the list of "jobs" a model can do (the "head" it has).
 - Click "**Text Classification**" (our movie review task). You are now looking at models that are ready to sort text into categories.
- **On a "Model Card" (e.g., `distilbert-base-uncased-finetuned-sst-2-english`):**
 - This is the "nutrition label" for the AI. It tells you what "body" was used (`distilbert-base-uncased`), what dataset it was fine-tuned on, and what its limitations are.
 - **The "Inference API" Widget:** This is the "wow" moment. It's a live, running demo of the model.
 - Type: `This movie was a masterpiece. I loved it.`
 - See the result: `POSITIVE (score: 0.999...)`
 - You are *using* a powerful, fine-tuned AI model directly in your browser, for free, with no code.
- **On the "Datasets" Page (e.g., `imdb`):**
 - This is the "job training" packet. You can see the raw data: a column for "text" (the review) and a column for "label" (0 or 1).
 - This is the exact "textbook" that was used to fine-tune the model you just tested.

7. Conclusion: Why This Matters (The "So What?")

This "pre-train / fine-tune" workflow has changed everything.

- A startup in a garage can now build a world-class AI app in a weekend.
- A hospital can fine-tune a model on its *own* X-rays to find tumors.
- A law firm can fine-tune a model on *its own* legal documents to find contracts.

You no longer need to be a giant corporation to build powerful, world-changing AI. You just need to be smart about finding the right "body" and fine-tuning it on the right "job training" data.

Student Notes: Beyond Chatbots, The Rise of AI Agents

1. The Big Problem: "A Chatbot Can't Do Anything"

In our last lesson, we learned about "fine-tuning," which lets us create specialized AI models that are experts at *one thing* (like classifying movie reviews).

But what about the real world?

Imagine you're using a chatbot and you say: "**It looks like it's going to rain today. Can you email my 2 PM meeting and move it to a coffee shop near my office?**"

A normal chatbot (like ChatGPT) will fail. It might *write* the email for you to copy and paste, but it can't:

1. **Know** your location or check the *real* weather.
2. **Access** your calendar to find the meeting attendees.
3. **Search** Google Maps for a "coffee shop near my office."
4. **Open** your email and *actually send* the message.

A chatbot can **talk**, but it can't **act**.

This is the problem AI Agents are built to solve.

2. What is an "AI Agent"?

An **AI Agent** is a system that uses a "brain" (a powerful LLM like GPT or Gemini) to **autonomously reason, plan, and take actions** to achieve a specific goal.

- A **Chatbot** *responds* to you.
- An **AI Agent** *works* for you.

It's the difference between asking someone for a recipe and having a personal chef who goes to the store, buys the ingredients, and cooks the meal for you.

3. How Agents "Think": The Core Loop

Agents work in a continuous cycle of thought and action. The most famous framework for this is called **ReAct (Reason + Act)**.

It's a simple, powerful loop:

1. **GOAL:** The user gives a complex goal (e.g., "Find me a cheap flight to Paris for next weekend and book it.")
2. **REASON:** The agent's "brain" (the LLM) thinks *step-by-step*.
 - *Thought:* "My goal is to book a flight. First, I need to know what 'next weekend' is. Then I need to search for flights."
3. **ACT:** The agent chooses a "**Tool**" to accomplish its thought.
 - *Action:* Use `calendar` to find the dates for next weekend.
4. **OBSERVE:** The agent gets the result from the tool.
 - *Observation:* "The dates are November 22nd to November 24th."
5. **REPEAT (Loop):** The agent takes this new information and starts the loop over.
 - *Thought:* "Great. Now I have the dates. My next step is to find flights."
 - *Action:* Use `Google Flights_api` to search for "flights from [My_Location] to Paris on Nov 22-24."
 - *Observation:* "I found 3 flights. Flight 1 is \$800. Flight 2 is \$950. Flight 3 is \$1200."
 - *Thought:* "The user asked for a 'cheap' flight. Flight 1 is the cheapest. I should ask the user to confirm before booking."
 - *Action:* Use `ask_user` tool: "I found a flight for \$800. Would you like me to book it?"

This loop—**Reason, Act, Observe**—continues until the final goal is achieved.

4. The Anatomy of an AI Agent (The Components)

Agents are built by combining several key parts:

1. **The Brain (The LLM):** This is the core reasoning engine. It's the "Ph.D. graduate" from our last lesson, but now it's been promoted to a "Manager." Its job isn't to *do* the work, but to *make a plan* and *delegate* tasks to its "tools."
 - *Examples:* GPT-4, Gemini, Claude 3
2. **The "Toolbelt" (The APIs & Functions):** This is the most important part! An agent with no tools is just a chatbot. Tools give the agent "superpowers" to interact with the real world.
 - `Google Search` (To look up information)
 - `weather_api` (To check the weather)
 - `calendar_api` (To manage your schedule)
 - `python_interpreter` (To run code, do math, or analyze data)
 - `email_sender` (To send messages)
3. **Memory (The "Clipboard" and "Notebook"):** The agent needs to remember things.
 - **Short-Term Memory:** The history of the conversation (what you just said).
 - **Long-Term Memory:** A special database (called a "Vector Database") that acts as the agent's permanent "notebook." This is where it can store facts, preferences ("I prefer window seats"), and lessons learned from past tasks.

5. Key Platforms (How People Build Agents)

You don't have to build all this from scratch. Companies have created "agent frameworks" or "platforms" that provide all the Lego blocks.

For Developers (The "Lego" Kits):

1. **LangChain:** This is the most famous "Swiss Army knife" for building AI apps. It's a library that gives developers all the pieces (or "chains") to connect a **Brain** (LLM) to **Tools** and **Memory**. It's powerful but complex.
2. **LlamaIndex:** This is similar to LangChain but specializes in connecting agents to *your own data*. It's the best tool for building agents that can "read" your PDFs, emails, or company documents (a process called RAG - Retrieval-Augmented Generation).

For Everyone (The "No-Code" Platforms):

You have probably already used a basic AI agent!

- **Custom GPT "Actions":** When you build a Custom GPT and use the "Actions" feature to connect it to an API (like the Weather API), you are *turning your chatbot into a true AI agent*. You are giving it a "tool" and a "brain" to use it.
- **Zapier AI / Make.com:** These platforms are "agent builders." You can visually link an AI's "thought" to real-world "actions" (like "when I get an email, have an AI read it, and if it's important, add it to my Trello board").

6. Conclusion: Why This Matters (The "So What?")

- **Phase 1 (The Past):** We had **AI Models** (like **BERT**). They were "bodies" with no "head," only useful for data scientists.
- **Phase 2 (The Present):** We have **Chatbots** (like ChatGPT). These are models with a "chat head." They are powerful "talkers" and "thinkers" that can help you with tasks.
- **Phase 3 (The Future):** We have **AI Agents**. These are models with a "chat head," "tools," and "autonomy." They are "doers" that can work for you.

This is the leap from AI as a "tool you use" to AI as an "assistant you manage" or an "autonomous teammate." It's the next major step in making AI useful in the real world.

Student Notes: Beyond Chatbots, The Rise of AI Agents

1. The Big Problem: "A Chatbot Can't Do Anything"

In our last lesson, we learned about "fine-tuning," which lets us create specialized AI models that are experts at *one thing* (like classifying movie reviews).

But what about the real world?

Imagine you're using a chatbot and you say: "**It looks like it's going to rain today. Can you email my 2 PM meeting and move it to a coffee shop near my office?**"

A normal chatbot (like ChatGPT) will fail. It might *write* the email for you to copy and paste, but it can't:

1. **Know** your location or check the *real* weather.
2. **Access** your calendar to find the meeting attendees.
3. **Search** Google Maps for a "coffee shop near my office."
4. **Open** your email and *actually send* the message.

A chatbot can **talk**, but it can't **act**.

This is the problem AI Agents are built to solve.

2. What is an "AI Agent"?

An **AI Agent** is a system that uses a "brain" (a powerful LLM like GPT or Gemini) to **autonomously reason, plan, and take actions** to achieve a specific goal.

- A **Chatbot** *responds* to you.
- An **AI Agent** *works* for you.

It's the difference between asking someone for a recipe and having a personal chef who goes to the store, buys the ingredients, and cooks the meal for you.

3. How Agents "Think": The Core Loop

Agents work in a continuous cycle of thought and action. The most famous framework for this is called **ReAct (Reason + Act)**.

It's a simple, powerful loop:

1. **GOAL:** The user gives a complex goal (e.g., "Find me a cheap flight to Paris for next weekend and book it.")
2. **REASON:** The agent's "brain" (the LLM) thinks *step-by-step*.
 - *Thought:* "My goal is to book a flight. First, I need to know what 'next weekend' is. Then I need to search for flights."
3. **ACT:** The agent chooses a "**Tool**" to accomplish its thought.
 - *Action:* Use `calendar` to find the dates for next weekend.
4. **OBSERVE:** The agent gets the result from the tool.
 - *Observation:* "The dates are November 22nd to November 24th."
5. **REPEAT (Loop):** The agent takes this new information and starts the loop over.
 - *Thought:* "Great. Now I have the dates. My next step is to find flights."
 - *Action:* Use `Google Flights_api` to search for "flights from [My_Location] to Paris on Nov 22-24."
 - *Observation:* "I found 3 flights. Flight 1 is \$800. Flight 2 is \$950. Flight 3 is \$1200."
 - *Thought:* "The user asked for a 'cheap' flight. Flight 1 is the cheapest. I should ask the user to confirm before booking."
 - *Action:* Use `ask_user` tool: "I found a flight for \$800. Would you like me to book it?"

This loop—**Reason, Act, Observe**—continues until the final goal is achieved.

4. The Anatomy of an AI Agent (The Components)

Agents are built by combining several key parts:

1. **The Brain (The LLM):** This is the core reasoning engine. It's the "Ph.D. graduate" from our last lesson, but now it's been promoted to a "Manager." Its job isn't to *do* the work, but to *make a plan* and *delegate* tasks to its "tools."
 - *Examples:* GPT-4, Gemini, Claude 3
2. **The "Toolbelt" (The APIs & Functions):** This is the most important part! An agent with no tools is just a chatbot. Tools give the agent "superpowers" to interact with the real world.
 - `Google Search` (To look up information)
 - `weather_api` (To check the weather)
 - `calendar_api` (To manage your schedule)
 - `python_interpreter` (To run code, do math, or analyze data)
 - `email_sender` (To send messages)
3. **Memory (The "Clipboard" and "Notebook"):** The agent needs to remember things.
 - **Short-Term Memory:** The history of the conversation (what you just said).
 - **Long-Term Memory:** A special database (called a "Vector Database") that acts as the agent's permanent "notebook." This is where it can store facts, preferences ("I prefer window seats"), and lessons learned from past tasks.

5. Key Platforms (How People Build Agents)

You don't have to build all this from scratch. Companies have created "agent frameworks" or "platforms" that provide all the Lego blocks.

For Developers (The "Lego" Kits):

1. **LangChain:** This is the most famous "Swiss Army knife" for building AI apps. It's a library that gives developers all the pieces (or "chains") to connect a **Brain** (LLM) to **Tools** and **Memory**. It's powerful but complex.
2. **LlamaIndex:** This is similar to LangChain but specializes in connecting agents to *your own data*. It's the best tool for building agents that can "read" your PDFs, emails, or company documents (a process called RAG - Retrieval-Augmented Generation).

For Everyone (The "No-Code" Platforms):

You have probably already used a basic AI agent!

- **Custom GPT "Actions":** When you build a Custom GPT and use the "Actions" feature to connect it to an API (like the Weather API), you are *turning your chatbot into a true AI agent*. You are giving it a "tool" and a "brain" to use it.
- **Zapier AI / Make.com:** These platforms are "agent builders." You can visually link an AI's "thought" to real-world "actions" (like "when I get an email, have an AI read it, and if it's important, add it to my Trello board").

6. Conclusion: Why This Matters (The "So What?")

- **Phase 1 (The Past):** We had **AI Models** (like **BERT**). They were "bodies" with no "head," only useful for data scientists.
- **Phase 2 (The Present):** We have **Chatbots** (like ChatGPT). These are models with a "chat head." They are powerful "talkers" and "thinkers" that can help you with tasks.
- **Phase 3 (The Future):** We have **AI Agents**. These are models with a "chat head," "tools," and "autonomy." They are "doers" that can work for you.

This is the leap from AI as a "tool you use" to AI as an "assistant you manage" or an "autonomous teammate." It's the next major step in making AI useful in the real world.