

In this module, you will explore foundational machine learning concepts that prepare you for hands-on modeling with Python. You will explain the relevance of Python and scikit-learn in machine learning, summarize the IBM AI Engineering certification path, and classify common types of learning algorithms. You'll outline the stages of the machine learning model lifecycle and describe what a typical day looks like for a machine learning engineer. You will also compare key roles in the AI field, identify widely used open-source tools, and learn to utilize scikit-learn to build and evaluate simple models.

## Learning Objectives

---

- Classify different types of machine learning algorithms and their applications.
- Explain the significance of Python and scikit-learn in machine learning as introduced in the course.
- Summarise the structure and goals of the IBM AI Engineering Professional Certificate program.
- Outline the stages involved in the machine learning model lifecycle.
- Describe the daily responsibilities and tasks of a machine learning engineer.
- Compare the roles and skill sets of data scientists and AI engineers.
- Identify common tools and libraries used in machine learning workflows.
- Utilize scikit-learn's modules to implement basic machine learning models.

Here are **structured notes** from your course introduction 

---

# Machine Learning with Python – Course Introduction

## Why Python for ML?

- **Most widely used language** for ML, data science, analytics.
  - Popular libraries: **scikit-learn**, **NumPy**, **Pandas**.
  - Python dominates the **machine learning ecosystem**.
-

## Course Level & Prerequisites

### What you will learn



Describe the role of machine learning in various career paths



Articulate the various stages of the machine-learning lifecycle



Discuss how machine-learning models work



Implement machine-learning models using Python and scikit-learn



Solve data-related problems using machine-learning methods

- **Level:** Intermediate.
- **Recommended Skills:**
  - Python programming
  - Pandas & NumPy
  - Data preparation & analysis
- **Helpful:** Basic familiarity with machine learning.

## What You'll Learn

### 1. Foundations

- Role of ML in various career paths.
- ML lifecycle & model workflow.
- Basics of supervised vs. unsupervised learning.

## 2. Modeling Techniques

- **Classification**
- **Regression**
- **Clustering**

## 3. Advanced Intro

- Reinforcement Learning
- Deep Learning
- Artificial Intelligence

## 4. Hands-On Skills

- Build, assess, and validate models.
  - Work with real-world datasets.
  - Use **Python + Pandas + NumPy + scikit-learn**.
- 

## Course Activities

- **Instructional videos** → Learn concepts.
  - **Hands-on labs** → Apply what you learned.
  - **Practice & graded quizzes** → Assess understanding.
  - **Glossary** → Strengthen mastery of terms.
  - **Final Project** → Apply all concepts in a real scenario.
- 

## Topics Covered

- Multiple Linear Regression
  - Logistic Regression
  - Prediction
  - Fraud Detection
  - K-Nearest Neighbors (KNN)
  - Support Vector Machines (SVM)
- 

## Outcome

- Build a **solid foundation** in applied ML.
  - Prepare for **advanced opportunities** in ML, Data Science & AI.
  - Gain both **theory & practical experience**.
- 

## Support

- **Discussion Forum** for queries.
  - Course team encourages engagement & practice.
- 

 *This course is your first step towards becoming proficient in Machine Learning with Python.*

---

Got it  Here's a **structured set of notes** from the IBM Professional Certificate Overview video you shared.

---



# IBM Professional Certificate (PC)

## Overview

### ◆ Programs Included

This course is part of two IBM Professional Certificates:

## About this program

### IBM Data Science Professional Certificate

- For: Data science careers
- Covers: Data cleaning, analysis, and predictive modeling
- Tools: Python, SQL, Pandas, NumPy

### IBM AI Engineering Professional Certificate

- For: Data scientists, ML engineers, software engineers
- Covers: How to build, train, deploy models and LLMs
- Tools: Python, SciPy, Keras, PyTorch, TensorFlow

#### 1. IBM Data Science PC

- Focus: Data Science skills
- Covers: Data cleaning, analysis, predictive modeling
- Tools: **Python, SQL, Pandas, NumPy**
- Outcome: Build a job-ready portfolio
- **Beginner-friendly**, no prerequisites

#### 2. IBM AI Engineering PC

- Focus: AI engineering careers

- Prepares: Data Scientists, ML Engineers, Software Engineers
  - Covers: Build, train, deploy models (incl. LLMs)
  - Tools/Libraries: **SciPy, scikit-learn, Keras, PyTorch, TensorFlow**
  - **Prerequisites:** Python, data analysis, basic ML/AI knowledge
- 

## ◆ Learning Approach

- **Hands-on learning:** Projects & labs for real-world experience
  - **Self-paced courses:** 2–6 modules each
  - **Final capstone projects** to apply skills
  - **Portfolio-ready** projects for career preparation
- 

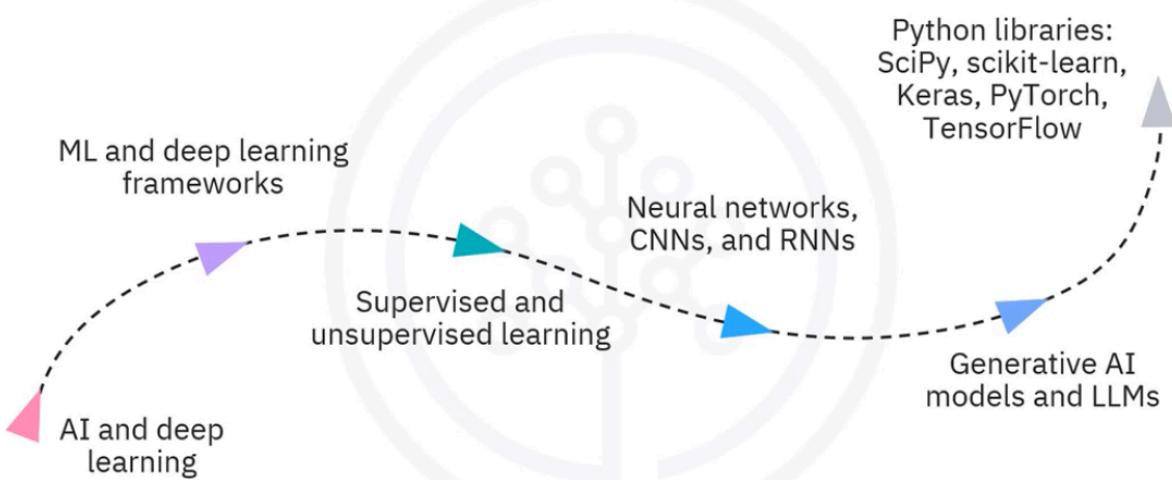
## ◆ Core Topics Covered



### Machine Learning Foundations

## AI Engineering PC: Program overview

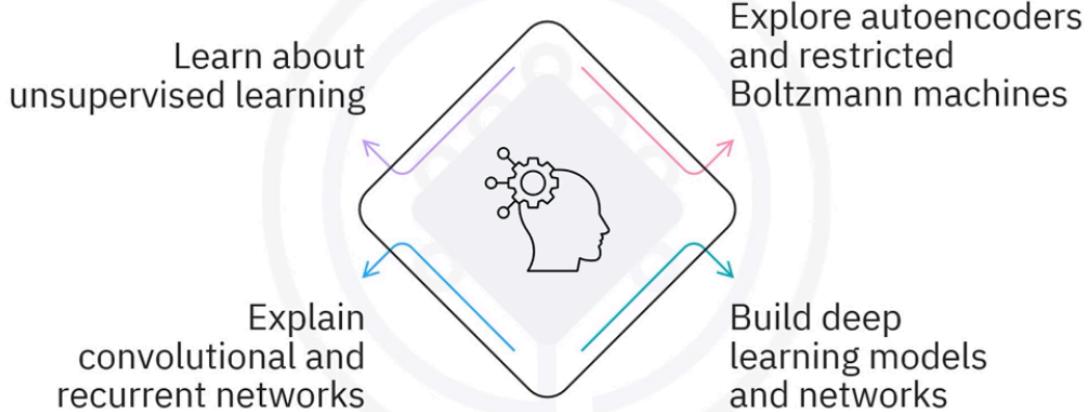
---



- ML lifecycle & supervised vs. unsupervised learning
- Algorithms: Linear regression, logistic regression, decision trees

## Deep Learning Basics

# Deep learning & neural networks with Keras

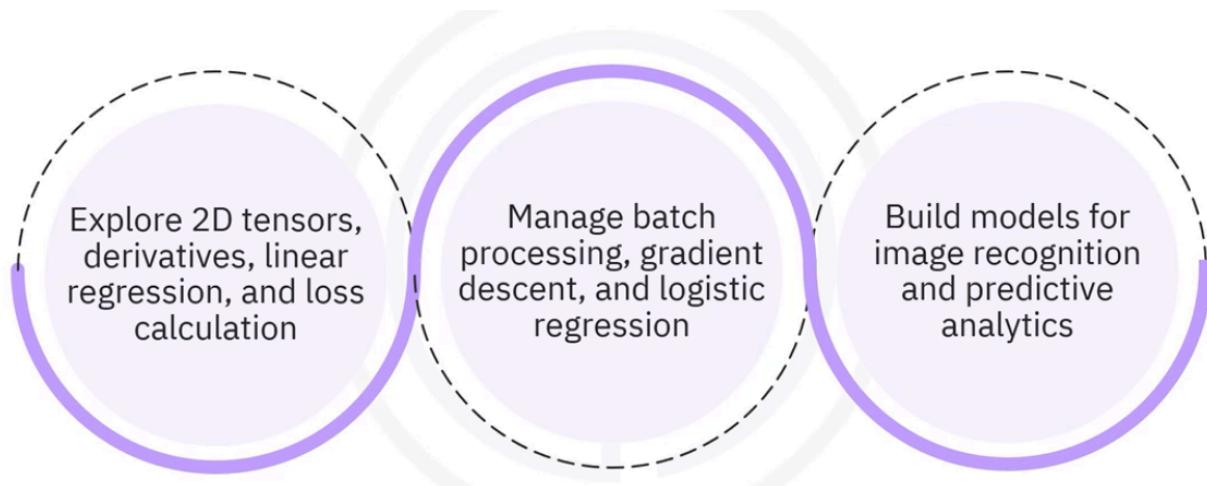


- Neural network architectures
- Unsupervised models: Autoencoders, RBMs
- CNNs & RNNs with **Keras & TensorFlow**
- Advanced: Custom layers, models, and training loops

# Deep learning with Keras and TensorFlow



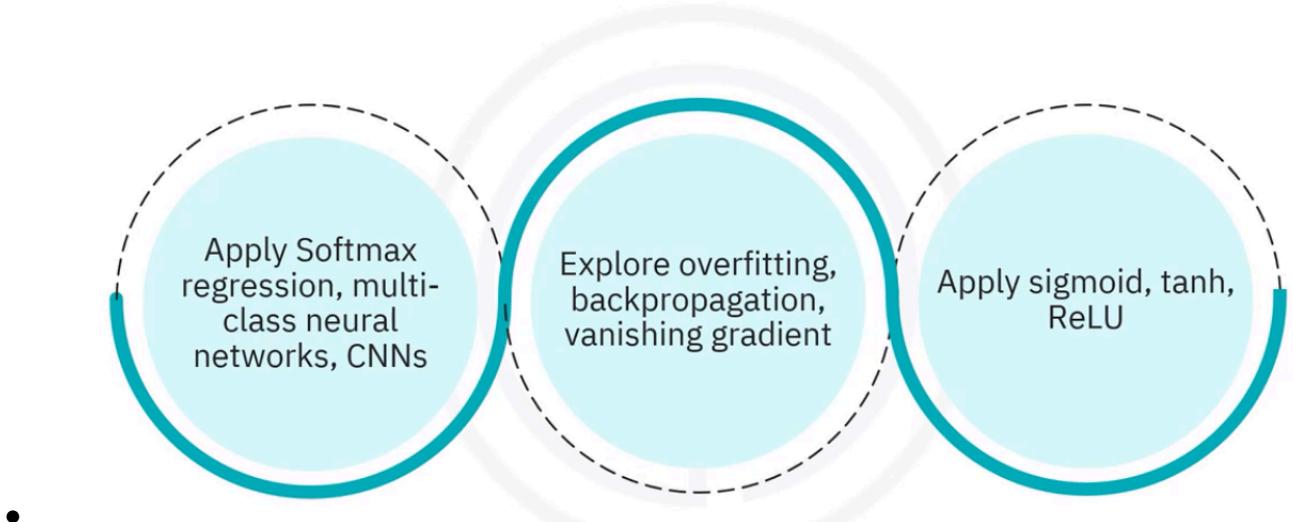
## PyTorch Skills



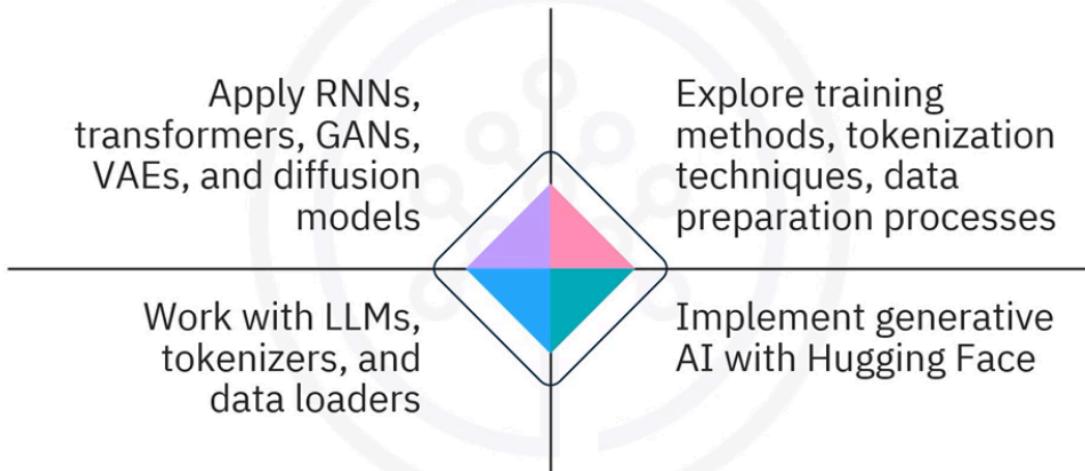
- Tensors, derivatives, loss functions, batch processing
- Logistic regression, gradient descent
- Multi-class NN, Softmax regression
- Handling overfitting, vanishing gradients, and backpropagation
- Activation functions: sigmoid, tanh, ReLU

## Deep learning with PyTorch

---



# Generative AI and LLMs



- CNNs for vision, transformers for sequential/time-series data
- Reinforcement Learning: Deep Q-Networks (DQNs)
- Generative AI models:
  - RNNs, Transformers, GANs, VAEs, Diffusion models
  - GPT & BERT (LLMs) with Hugging Face
- NLP basics: one-hot encoding, embeddings, Word2Vec
- Sequence-to-sequence models (translation, text generation)
- Evaluation: BLEU scores

# Fundamentals of AI agents using RAG and LangChain



- Fine-tuning techniques: **PEFT, LoRA, QLoRA, Prompting**
- RLHF (Reinforcement Learning with Human Feedback)
- PPO, DPO (optimization strategies)
- Instruction tuning, reward modeling

## 🔍 RAG & LangChain

### Generative AI application with RAG and LangChain



- Retrieval Augmented Generation (RAG) concepts
  - Prompt engineering best practices
  - LangChain tools & components
  - Real-world apps: QA bots, vector databases, Gradio interfaces
  - IBM Watsonx integration
- 

## ◆ Projects & Assessments

- **Applied projects:**
    - Preprocessing, building, training, testing ML/DL models
    - Generative AI projects (translation, QA bot, text generation)
  - **Final projects:**
    - Capstone in AI Engineering PC → Solve real-world problem
    - Showcase in portfolio (interview-ready)
  - **Quizzes + labs** for regular assessment
- 

## 🎯 Outcomes

- **IBM Data Science PC** → Entry-level Data Science skills, portfolio projects
  - **IBM AI Engineering PC** → AI/ML engineering role readiness, advanced projects
  - Gain **career-ready expertise** in ML, DL, NLP, LLMs, RAG, and deployment
-

✨ By the end, you'll have both theory and hands-on experience in Machine Learning, Deep Learning, and Generative AI — with projects to showcase in interviews.

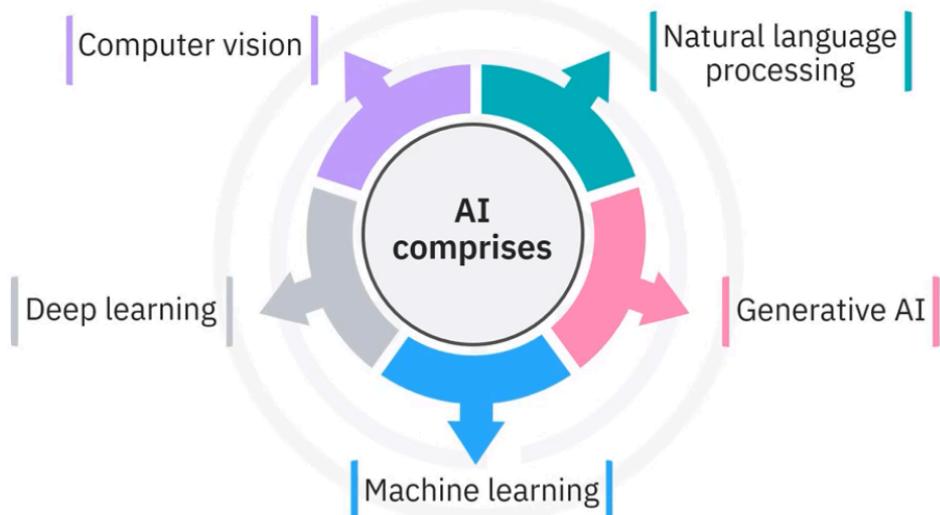
---

## Machine Learning Techniques and Applications – Notes

### 1. What is AI & ML?

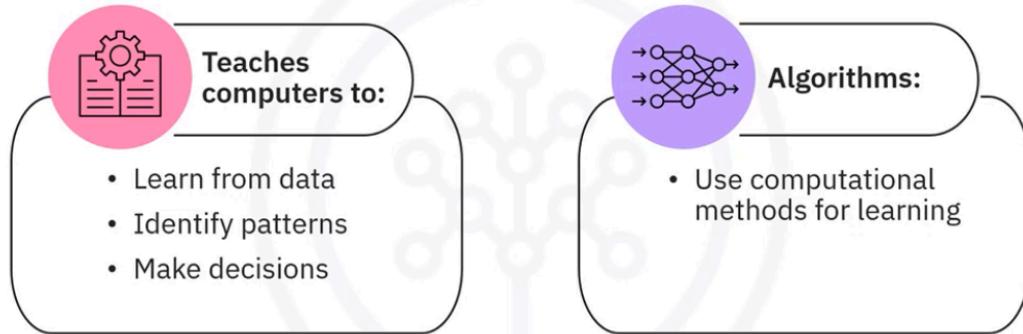
- **Artificial Intelligence (AI):** Makes computers simulate human cognitive abilities.

### AI, machine learning, and deep learning



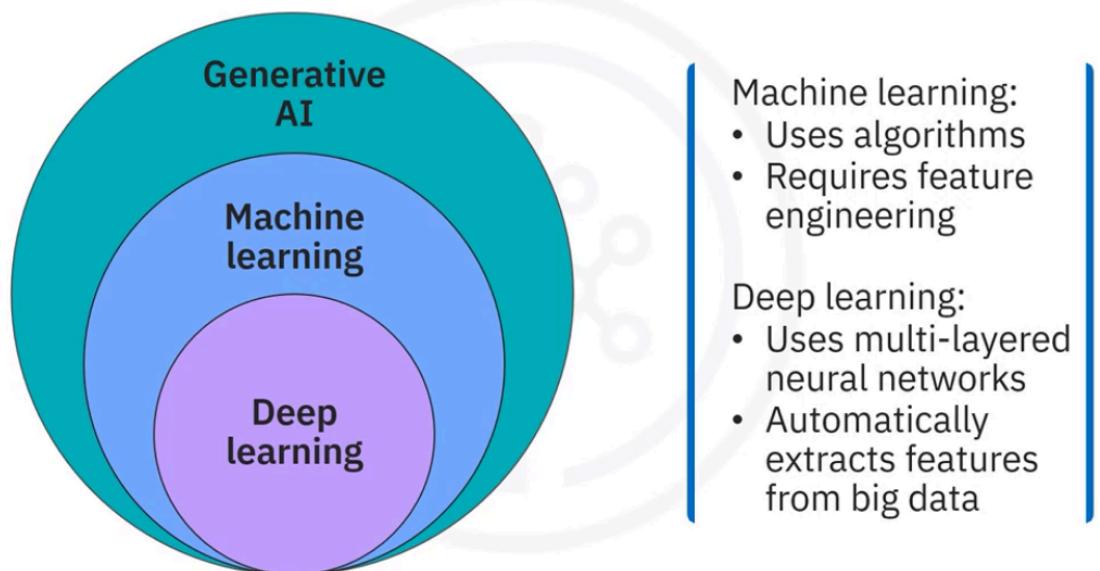
- Includes: Computer Vision, NLP, Generative AI, Machine Learning, and Deep Learning.
- **Machine Learning (ML):**

# How machine learning works



- Subset of AI.
- Uses algorithms + feature engineering.
- Learns patterns from data → makes predictions/decisions **without explicit programming**.
- Deep Learning (DL):

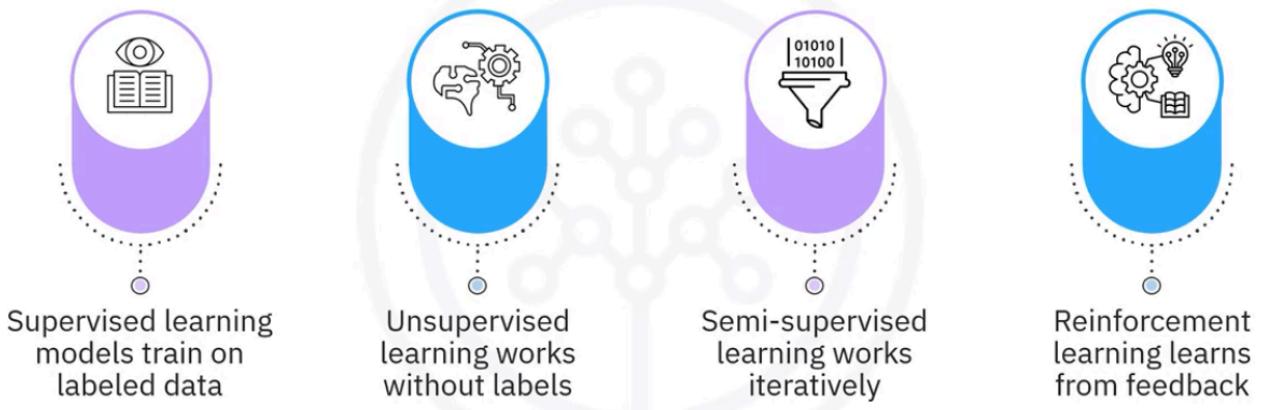
## AI, machine learning, and deep learning



- Subset of ML.
  - Uses multi-layered neural networks.
  - Automatically extracts features from complex, unstructured data (e.g., images, text, video).
  -
- 

## ◆ Types of Learning

# Machine learning paradigms

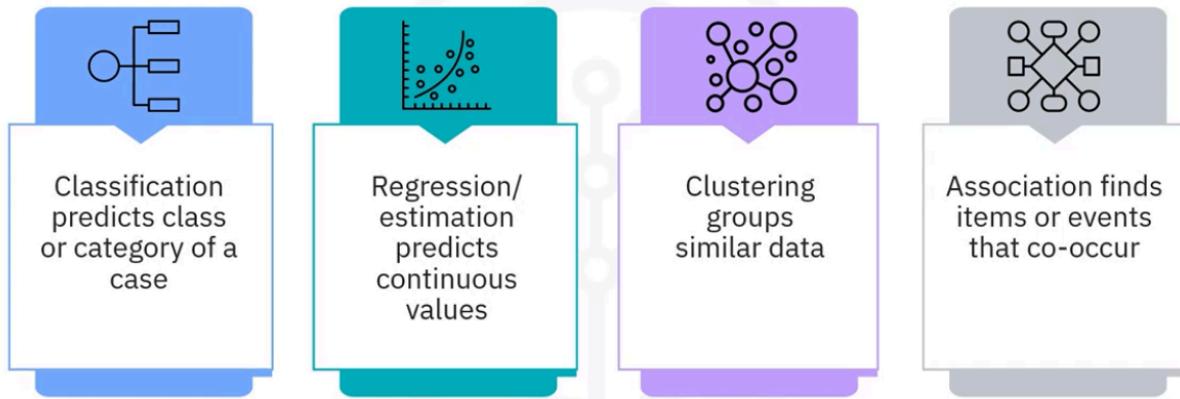


- 1. Supervised Learning** → Uses *labelled data* to make predictions.
  - Example: Predict tumour type (benign/malignant).
- 2. Unsupervised Learning** → Works with *unlabeled data*, finds hidden patterns.
  - Example: Customer segmentation.
- 3. Semi-supervised Learning** → Mix of labeled + unlabeled data; model self-labels high-confidence data.
- 4. Reinforcement Learning (RL)** → Agent interacts with environment, learns from rewards/penalties.

- Example: Game-playing AI, self-driving cars.
- 

## ◆ Key ML Techniques

# Machine learning techniques



**Classification** – predicts class/category.

- Example: Benign vs. malignant cell, customer churn.

**Regression (Estimation)** – predicts continuous values.

- Example: House price prediction, CO<sub>2</sub> emission.

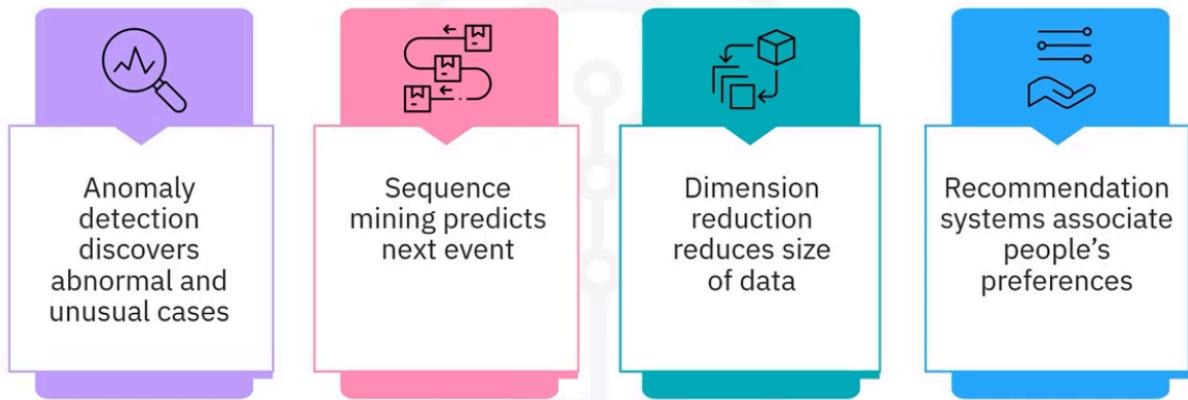
**Clustering** – groups similar data points.

- Example: Patient segmentation, customer segmentation.

**Association** – finds items/events that co-occur.

- Example: Market basket analysis (grocery items bought together).

# Machine learning techniques



**Classification** – predicts class/category.

- Example: Benign vs. malignant cell, customer churn.

**Regression (Estimation)** – predicts continuous values.

- Example: House price prediction, CO<sub>2</sub> emission.

**Clustering** – groups similar data points.

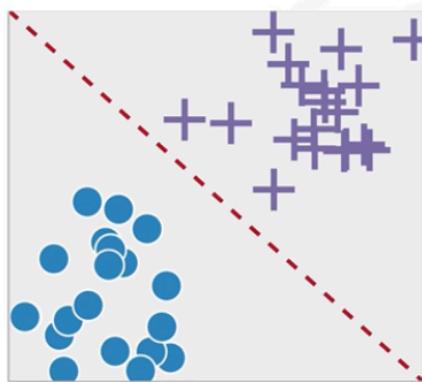
- Example: Patient segmentation, customer segmentation.

**Association** – finds items/events that co-occur.

- Example: Market basket analysis (grocery items bought together).

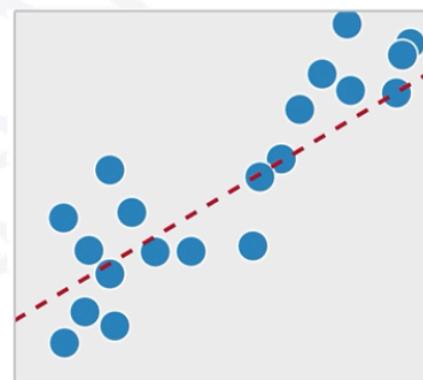
# Classification and regression

## Classification



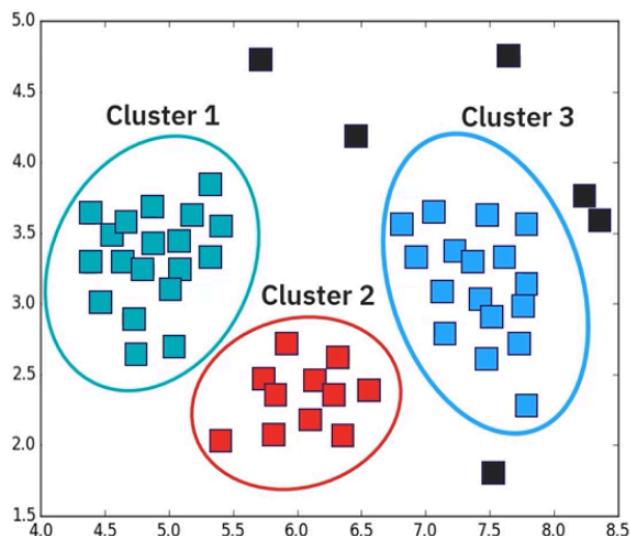
Predicts class membership of new observations

## Regression



Predicts continuous numerical values and input data

# Clustering

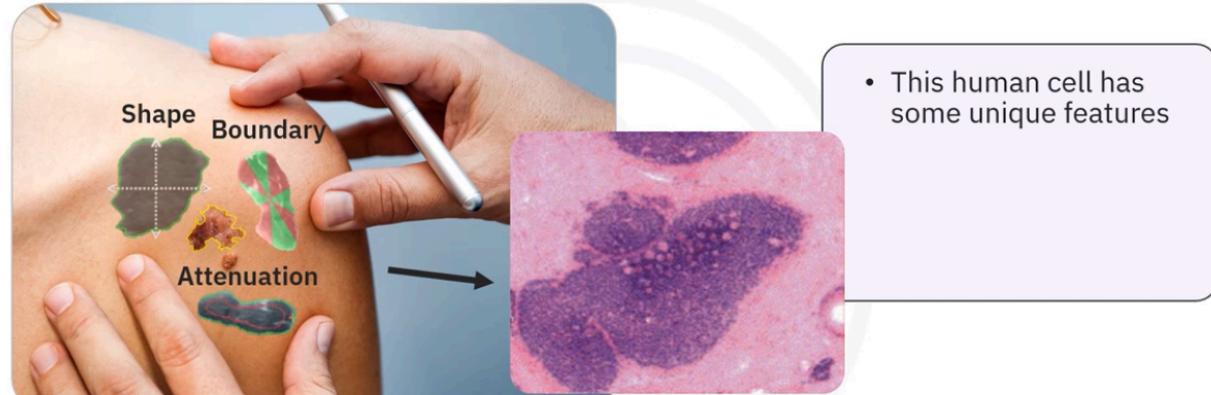


Unsupervised machine learning technique

Used to group similar data points or objects

- ◆ Applications of ML

# Applications of machine learning



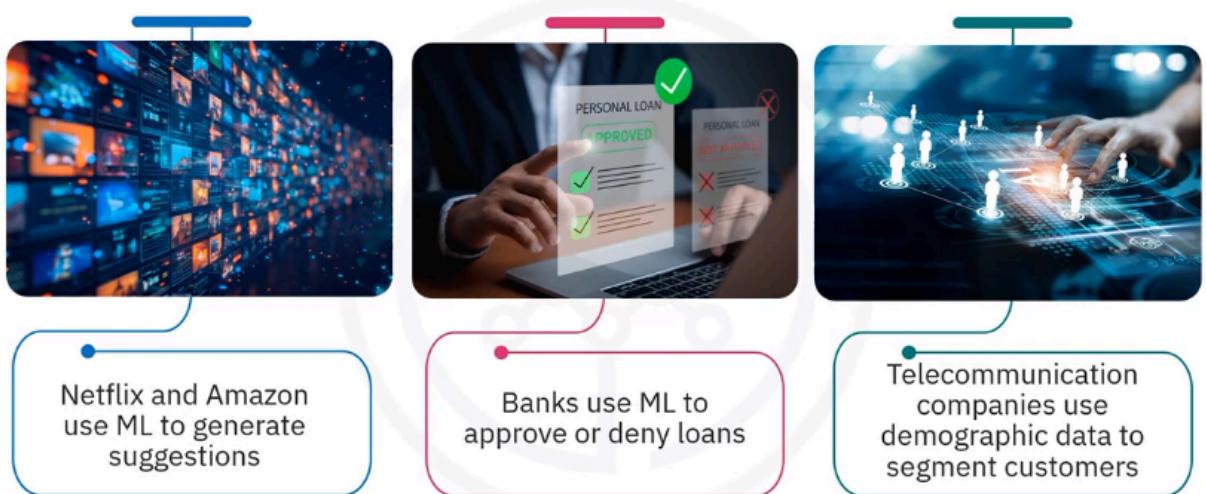
ID	Clump	UnifSize	UnifShape	MargAdh	SingEpi Size	Bare Nuc	Bland Chrom	Norm Nucl	Mit	Class
1000015	6	1	1	1	7	1	3	1	1	

# Applications of machine learning



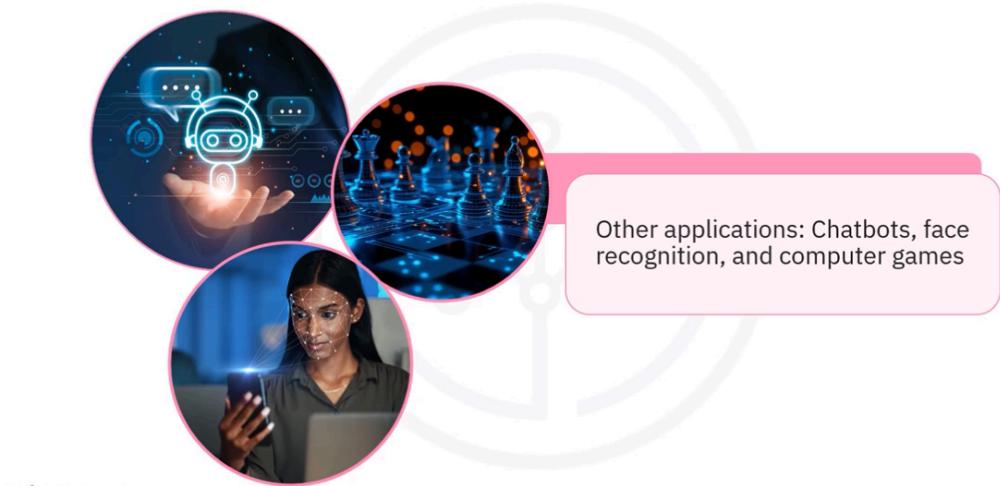
- **Healthcare:** Early disease prediction (cancer cell classification).

# Applications of machine learning



- **Finance:** Loan approval prediction, fraud detection.
- **Telecom:** Customer churn prediction.
- **Retail/E-commerce:** Recommendation engines, customer segmentation.
- **Computer Vision:** Identify objects (cats vs. dogs in images).

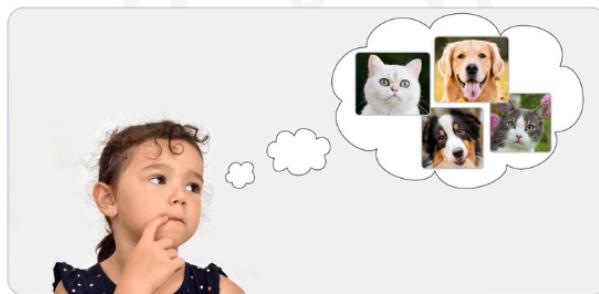
## Pairing machine learning with human intelligence



- **Everyday Life:** Virtual assistants (chatbots, Siri), face recognition, gaming (chess AI).

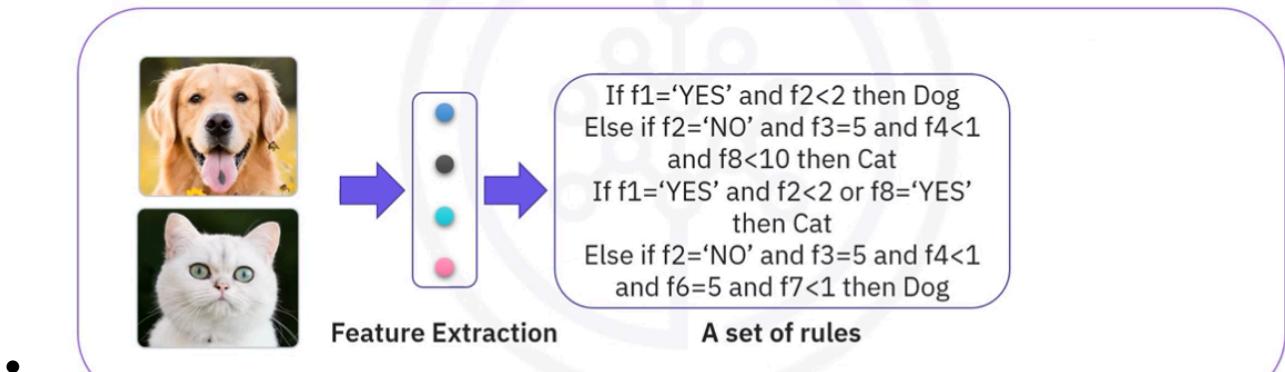
## Image recognition with machine learning

- **Data:** Images of cats and dogs



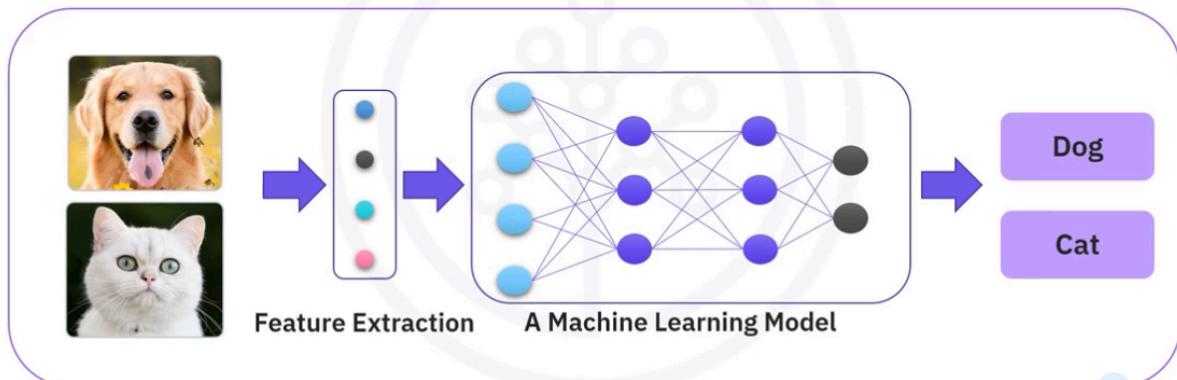
# Image recognition with machine learning

- **Data:** Images of cats and dogs
- **Before ML:** Create rules to detect the animals



# Image recognition with machine learning

- **Data:** Images of cats and dogs
- **Before ML:** Create rules to detect the animals
- **With ML:** Build a model to infer the animal type



## ◆ Key Takeaways

# Pairing machine learning with human intelligence



Increasing ML applications

Humans are still in the mix

Example: If ML algorithm refuses a loan, a banker wants to know why

- ML enables systems to learn directly from **data** without fixed rules.
- Different ML techniques are chosen depending on **data type, problem, and resources**.
- ML impacts multiple industries: healthcare, banking, telecom, retail, entertainment.
- **Humans remain in the loop** → Interpret results, ensure fairness, accountability, and explainability.

---

Here's a **well-structured note** on the **Machine Learning Model Lifecycle** from your text 👇

## Machine Learning Model Lifecycle – Notes

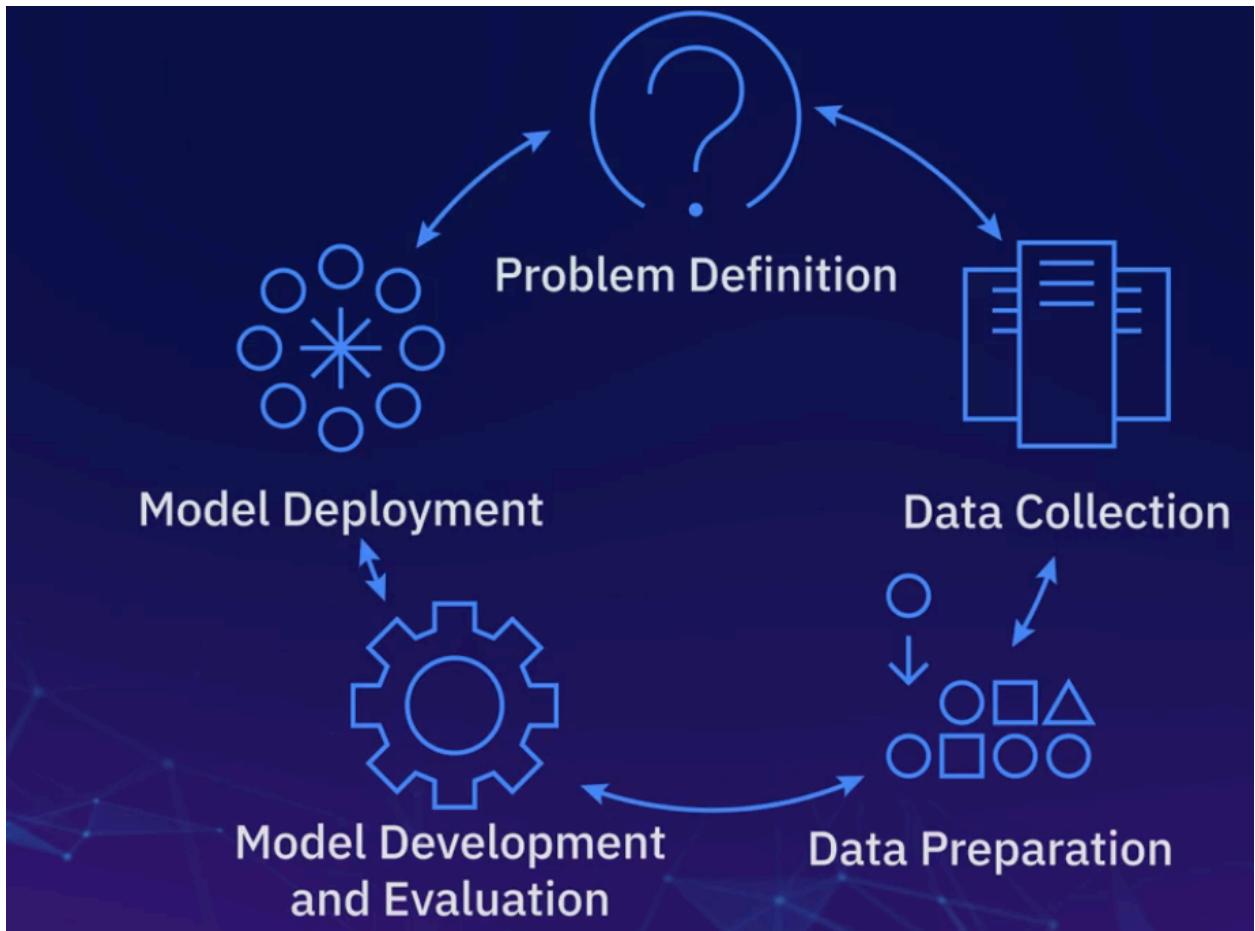
# Process Lifecycle



## 1. Overview

- A **Machine Learning (ML) model lifecycle** describes the steps taken from identifying a problem to deploying a working ML solution.
  - The process is **iterative** → steps are revisited if issues occur in production.
- 

## 2. Key Stages in the ML Lifecycle



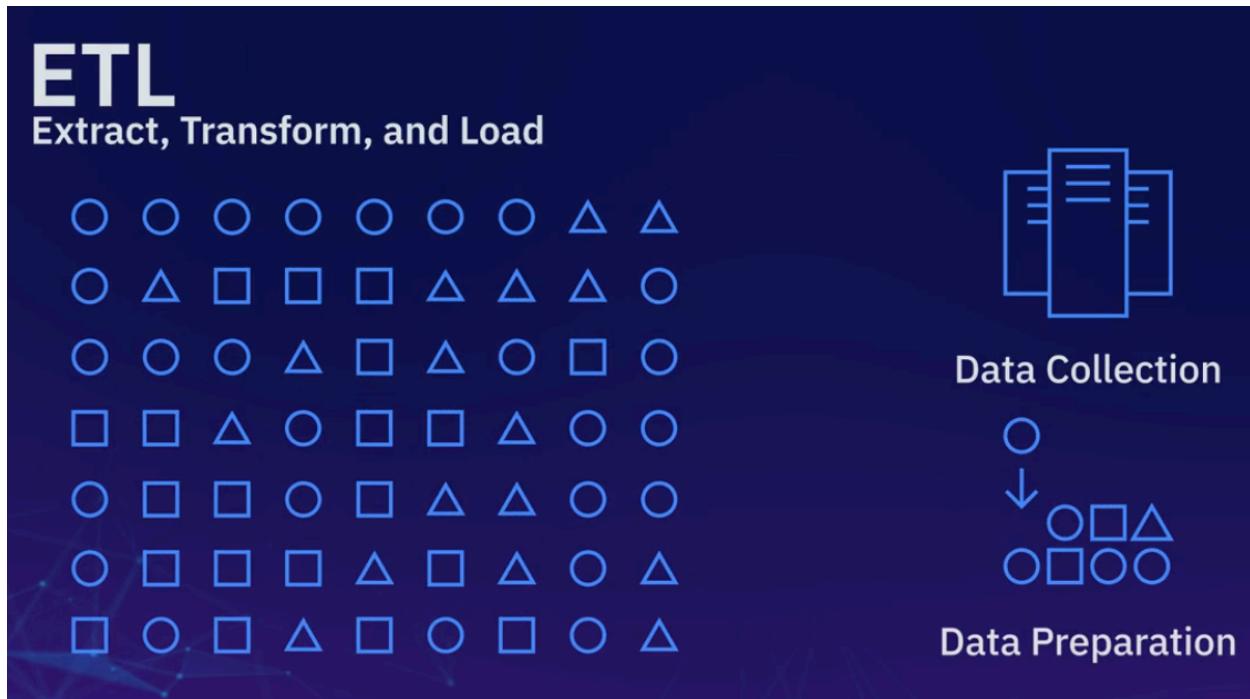
#### ◆ 1. Problem Definition

- Clearly state the **problem/situation** to solve.
- Aligns ML goals with business objectives.

#### ◆ 2. Data Collection

- Gather raw data from multiple sources.
- Must ensure **relevance, completeness, and quality**.

#### ◆ 3. Data Preparation (ETL Process)



- **ETL = Extract, Transform, Load**

- **Extract:** Collect data from sources.
- **Transform:** Clean, format, and structure the data.
- **Load:** Store into a single accessible location (e.g., data warehouse).

- Prepares data for model building.

- ◆ **4. Model Development & Evaluation**

- Select suitable ML algorithm(s).
- Train and tune the model using prepared data.
- Evaluate model performance (accuracy, precision, recall, etc.).

- ◆ **5. Model Deployment**

- Deploy trained model into a production environment.

- Continuously monitor model performance.
  - Retrain or re-iterate lifecycle if performance drops.
- 

### 3. Iterative Nature

- Lifecycle is **not linear**.
  - Example:
    - If a deployed model performs poorly → revisit **data preparation** or even **problem definition**.
- 

### 4. Key Takeaway

- ML Lifecycle = **Problem Definition** → **Data Collection** → **Data Preparation (ETL)** → **Model Development & Evaluation** → **Model Deployment**.
  - It is **cyclical**, allowing continuous improvement.
- 

Here's a **clear and structured note** from *A Day in the Life of a Machine Learning Engineer* 

---

## A Day in the Life of a Machine Learning Engineer

### Learning Goals

- Understand the **importance and requirements** of each process in the ML lifecycle.

- Identify which processes are **more time-consuming**.
- 

- ◆ **Case Example: Beauty Product Recommendation System**

- **Business Goal:** Increase revenue by recommending products to customers based on their purchase history.
- **Pain Point (Problem Definition):**

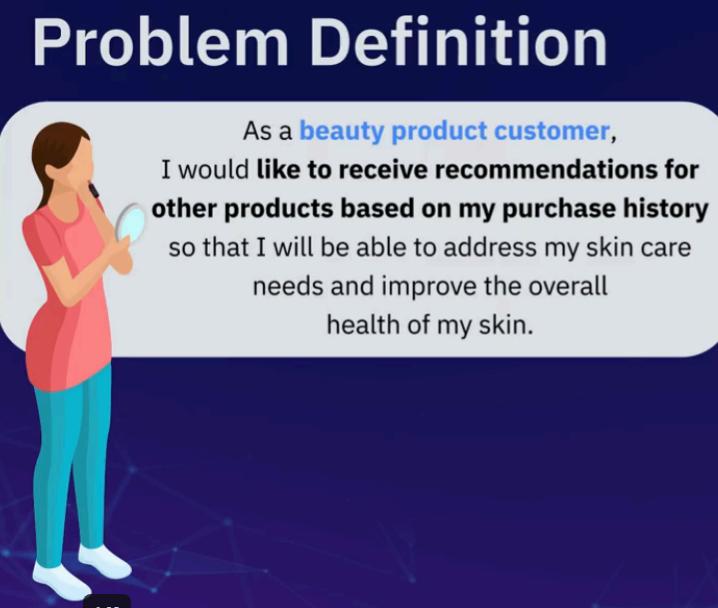
“As a beauty product customer, I would like to receive recommendations for other products based on my purchase history so I can address skincare needs and improve skin health.”

---

- ◆ **Lifecycle Steps in Action**

### 1. Problem Definition

## Problem Definition



As a **beauty product customer**,  
I would **like to receive recommendations for  
other products based on my purchase history**  
so that I will be able to address my skin care  
needs and improve the overall  
health of my skin.



Purchased      Recommended

- Align ML solution with client's needs.
  - Ensure the problem statement is **clear, measurable, and relevant**.
- 

## 2. Data Collection

# Data Collection



**User Data**  
demographics, purchase history, etc.

**Product Data**  
inventory, ingredients, popularity

**Other Data**  
saved products, liked products, search history

- Sources:
  - User data → demographics, purchase history, transactions.
  - Product data → ingredients, ratings, popularity, inventory.
  - Behavioral data → saved/liked products, search history, most visited items.
- Aim: **Gather comprehensive data** to understand both users & products.



---

### 3. Data Preparation & EDA

## Data Preparation

The background features a blue gradient with a network of white dots and lines. In the foreground, there are three blue 3D cubes of increasing size, stacked on a grey base.

- Clean data to remove irrelevant data
- Remove extreme values
- Missing values should be removed or replaced
- Each data column should be in the proper format
- Create new features
- Exploratory Data Analysis
- Identify how to split up the data into training and test sets

- Handle **errors, formatting issues, missing values.**
- Clean irrelevant/extreme data.
- Feature engineering: e.g.
  - Avg. duration between user transactions.
  - Most frequently purchased products.
  - Mapping skin issues targeted by products.

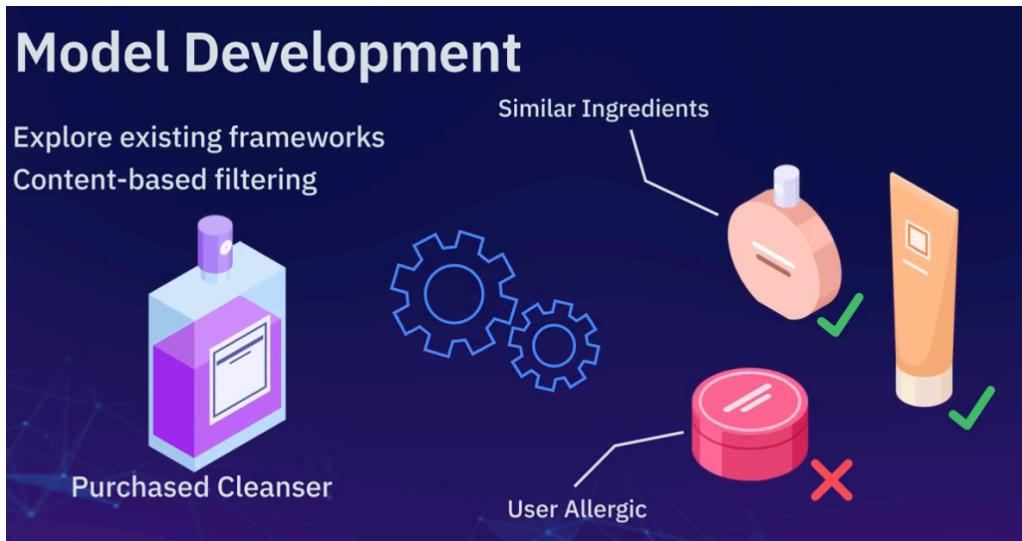
- **Exploratory Data Analysis (EDA):**
  - Use plots & correlation analysis.
  - Validate patterns with domain experts.
- **Train/Test split:** Example → use most recent transaction as test set.

⌚ Note: Data collection + preparation = **most time-consuming**.

---

## 4. Model Development

- **Techniques Used:**
  - *Content-Based Filtering*: Recommend products similar to what a user purchased.



- *Collaborative Filtering*: Recommend based on similar users' preferences/ratings.

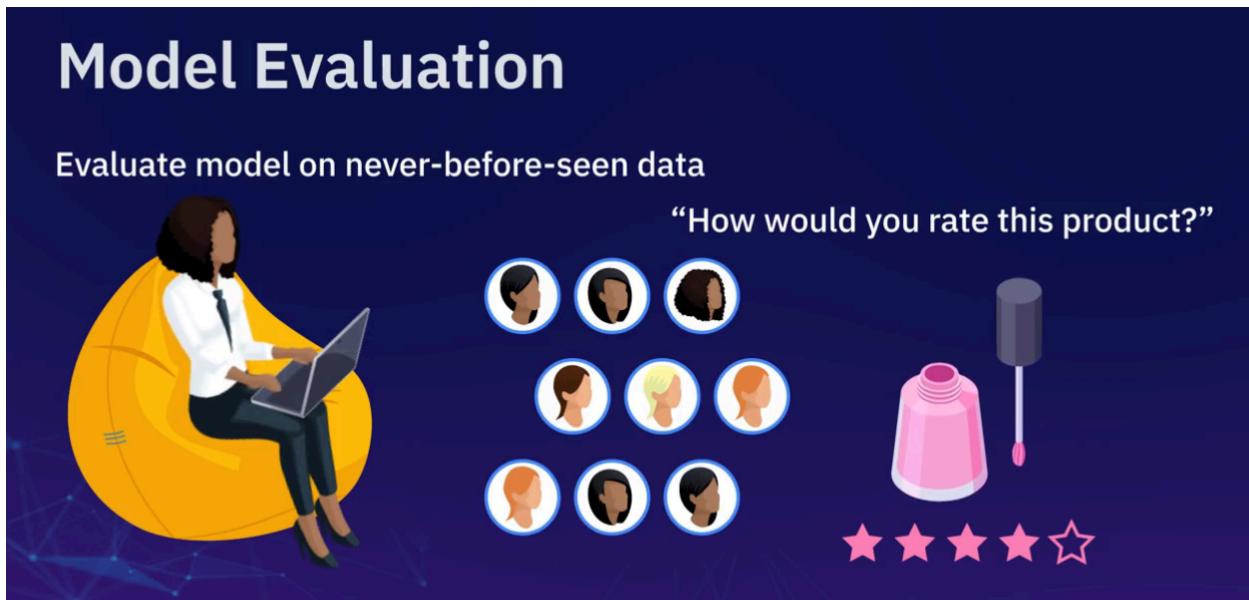
## Model Development

Collaborative filtering



- Final model = **Hybrid** of both techniques.
- 

## 5. Model Evaluation



- Initial evaluation → test dataset.
  - Tune model parameters.
  - Real-world validation → deploy to small user group for feedback.
    - Metrics: ratings, clicks, purchases, engagement.
-

## 6. Deployment



- Integrated into app & website.
- Model made available to customers in real-time.

---

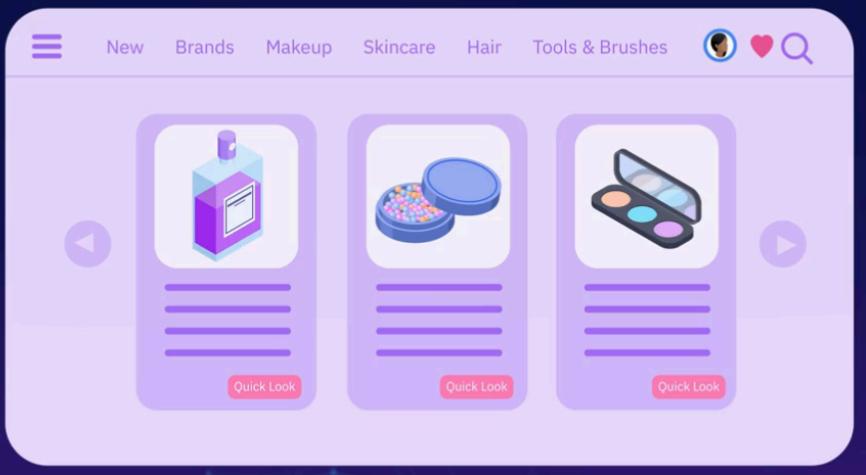
## 7. Monitoring & Iteration

# Model Deployment

Continue to track model performance



Retrain model based on new information



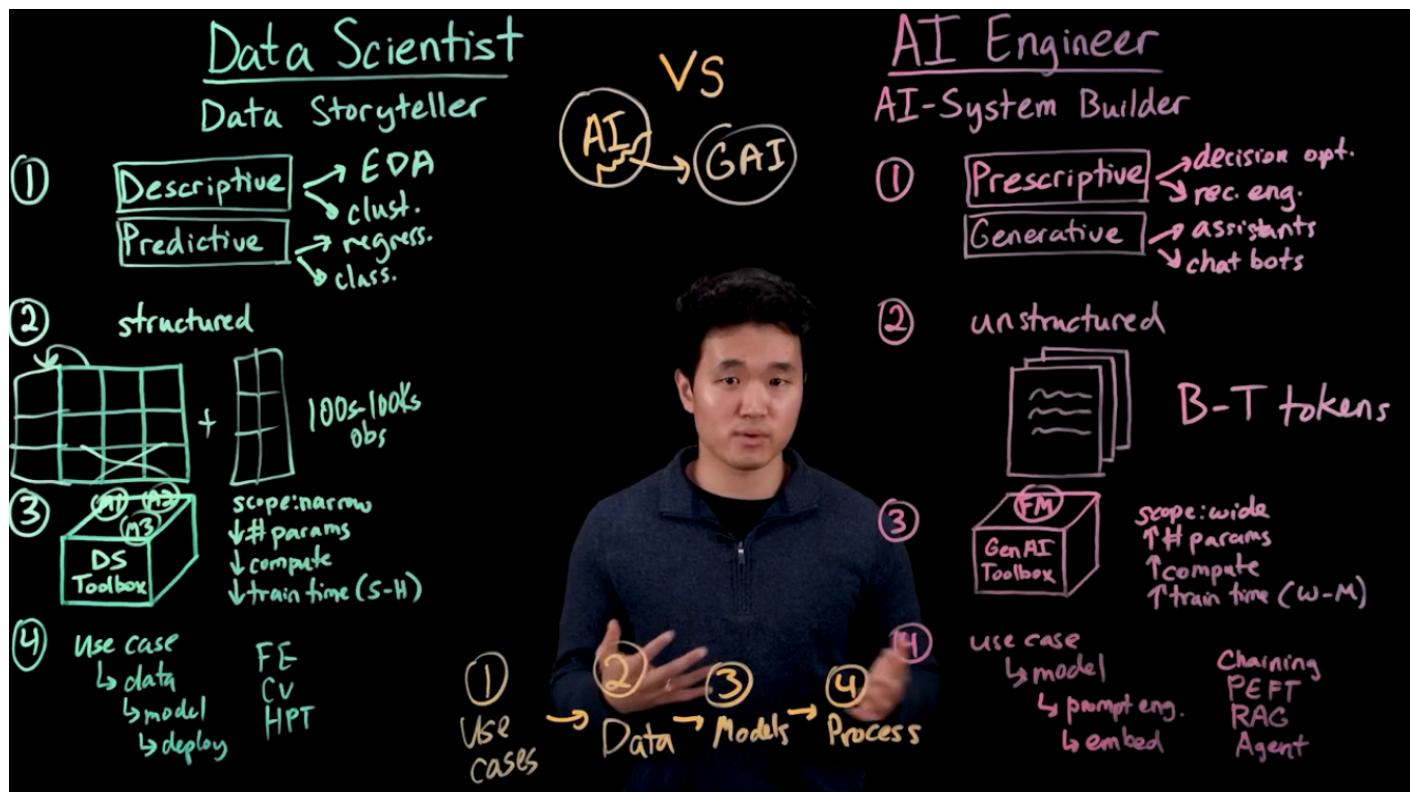
- Track ongoing performance (accuracy, click-through, conversions).
- Retrain/update with new data.
- Continuous improvement cycle.

## ◆ Key Takeaways

- **Every step** in ML lifecycle is crucial for success.
- **Data Collection & Preparation** are the **most time-consuming** processes.
- After deployment → continuous **monitoring & retraining** ensures long-term performance.

Here's a **crisp, structured note** comparing **Data Scientists vs AI Engineers (Generative AI Engineers)** ↴

# Data Scientist vs AI Engineer



## 🌐 Industry Context

- **Past:** Data Science = "sexiest job of the 21st century."
- **Now:** Rise of **AI Engineers** due to **Generative AI breakthroughs**.
- Generative AI has become distinct enough to form its own field → **AI Engineering**.

## 🔑 4 Key Differences

### 1. Use Cases

- **Data Scientist** = *Data Storyteller*

- Descriptive: Exploratory Data Analysis (EDA), Clustering (customer segmentation).
  - Predictive:
    - Regression (predict numeric values like revenue).
    - Classification (predict categories like success/failure).
  - **AI Engineer = AI System Builder**
    - Prescriptive: Decision Optimization, Recommendation Engines.
    - Generative: Chatbots, coding assistants, digital advisors, content summarization.
- 

## 2. Data

- **Data Scientist**
    - Works mainly with **structured/tabular data** (hundreds → thousands of rows).
    - Focus: Cleaning (remove outliers, joins, filtering), Feature Engineering.
  - **AI Engineer**
    - Works mainly with **unstructured data** (text, images, audio, video).
    - Example: Large Language Models (LLMs) → billions/trillions of tokens.
- 

## 3. Models

- **Data Scientist**
  - Toolbox = many specialized models (linear regression, decision trees, clustering).
  - Narrow scope → models generalize less.
  - Smaller models (seconds → hours to train, less compute).

- **AI Engineer**
    - Toolbox = **Foundation Models** (e.g., GPT, LLaMA, Stable Diffusion).
    - Wide scope → generalize across tasks.
    - Massive size (billions of parameters).
    - Training requires huge compute (hundreds–thousands GPUs, weeks–months).
- 

## 4. Processes & Techniques

- **Data Science Workflow**
  1. Define use case.
  2. Gather & prepare data.
  3. Train/validate model → Feature Engineering, Cross-validation, Hyperparameter Tuning.
  4. Deploy for real-time predictions.
- **AI Engineering Workflow**
  1. Define use case.
  2. Use **pre-trained foundation models** (skip raw training).
  3. Interact via **Prompt Engineering**.
  4. Extend with frameworks:
    - Prompt chaining.
    - Parameter-Efficient Fine-Tuning (PEFT).
    - Retrieval-Augmented Generation (RAG).
    - Autonomous agents.

5. Embed in larger system (assistants, apps, automation).
- 

## ◆ **Similarities & Overlap**

- Both fields still intersect:
    - Data Scientists can do prescriptive tasks.
    - AI Engineers can work with structured data.
- 



## **Key Takeaway**

- **Data Scientist** = translates messy data into insights & predictions.
  - **AI Engineer** = builds generative systems using foundation models to transform processes.
  - Both fields are evolving rapidly → new tools, research, and opportunities daily. Here's a **well-structured summary note** of the lecture on **Tools for Machine Learning**
-

# Tools for Machine Learning

## ◆ Importance of Data

### Data



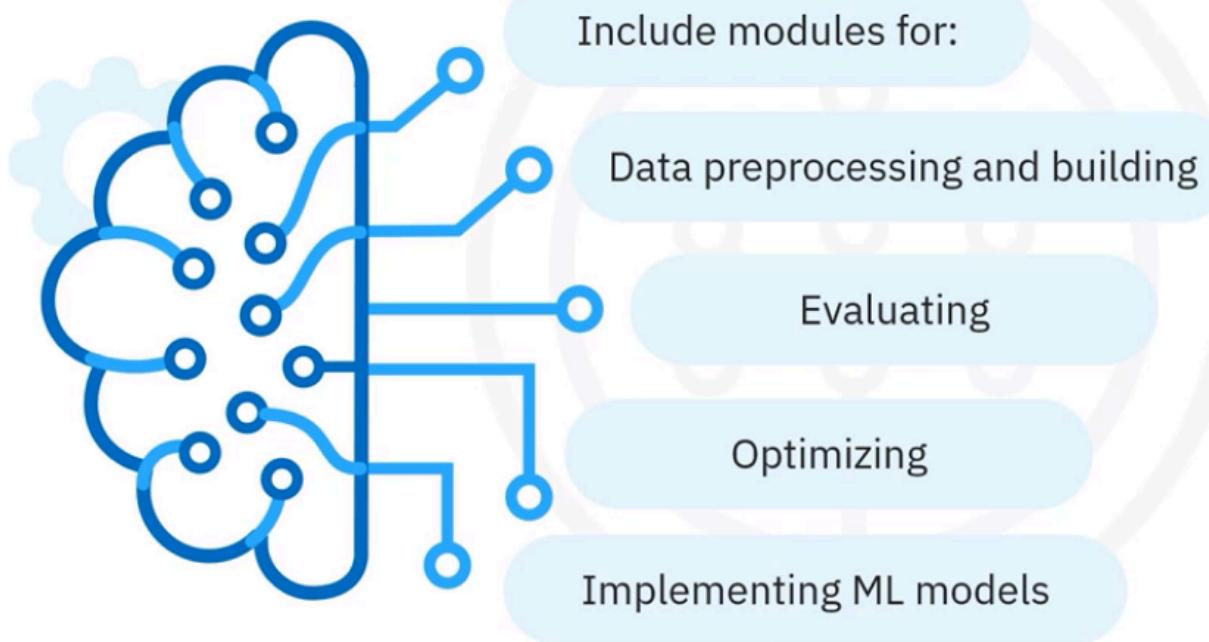
- **Data = raw facts, figures, information.**
- Central to **every ML algorithm** → patterns, predictions, insights.
- Drives all stages of the ML pipeline.

## ◆ Machine Learning Tools

- Provide functionalities for:
  - Data preprocessing
  - Model building, evaluation, optimization
  - Model deployment

# Machine learning tools

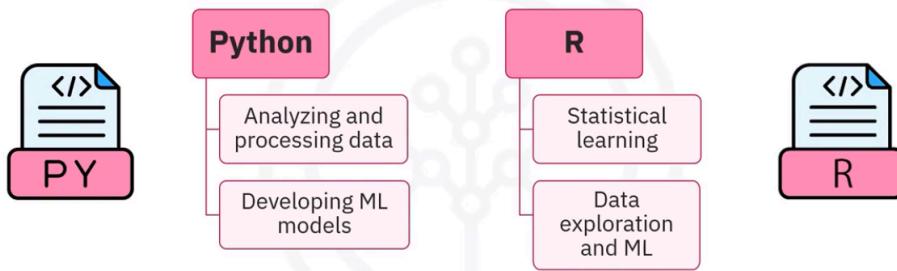
Provide functionalities for ML pipelines



- Simplify tasks like handling big data, statistical analysis, prediction.
- Example: **Pandas (data analysis), Scikit-learn (ML models)**.

## ◆ Programming Languages for ML

### Machine learning programming language



- **Python** → Most popular (rich libraries, easy to use).
- **R** → Statistical learning, visualization.

## Machine learning programming language



**Julia:** Parallel and distributed numerical computing support

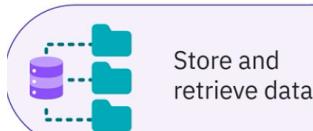
**Scala:** Processing big data and building ML pipelines

**Java:** Supporting scalable ML applications

**JavaScript:** Running ML models in web browsers

- **Julia** → High-performance, numerical computing.
- **Scala** → Big data, ML pipelines.
- **Java** → Scalable ML applications in production.
- **JavaScript** → Runs ML in browsers (client-side apps).

## Uses of machine learning tools



Store and retrieve data



Explore, visually inspect, and clean data



Work with plots, graphs, and dashboards



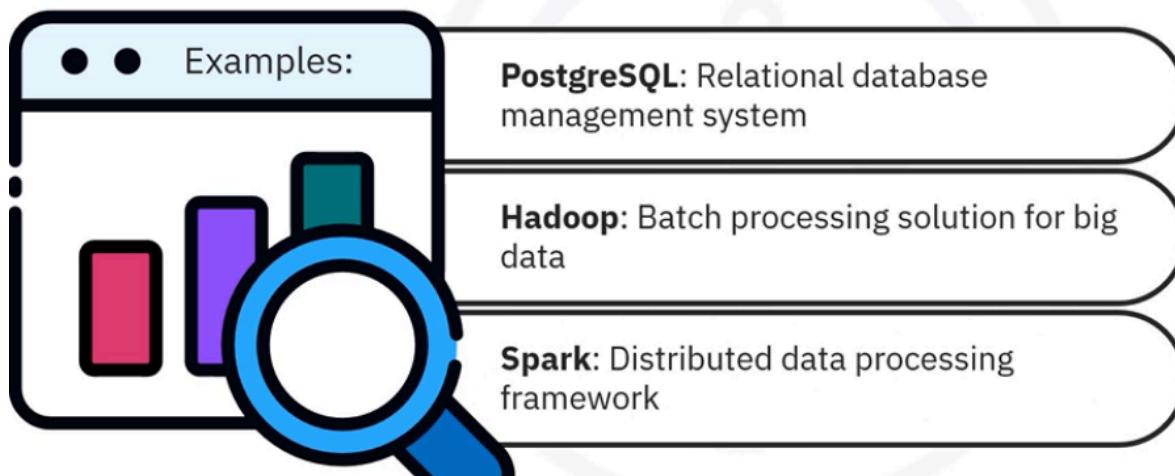
Prepare data for developing ML models

## ◆ Categories of ML Tools

### 1. Data Processing & Analytics

#### Data processing and analytics

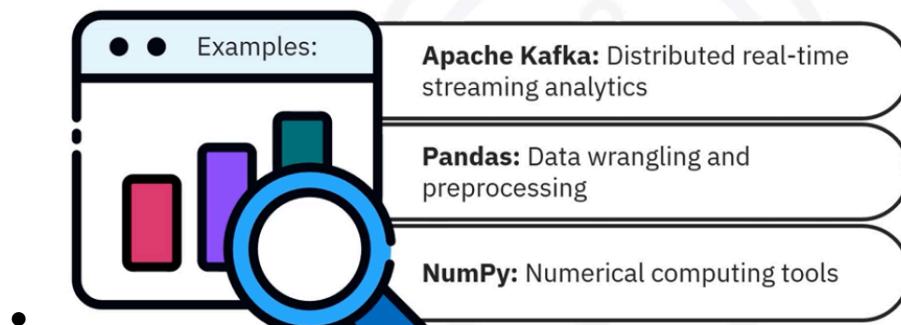
Process, store, and interact with data.



- **PostgreSQL** → SQL-based relational database.
- **Hadoop** → Disk-based, scalable batch processing for massive data.
- **Spark** → In-memory, distributed framework (faster than Hadoop).

#### Data processing and analytics

Process, store, and interact with data.



- **Kafka** → Real-time streaming & pipelines.
- **Pandas** → DataFrames, cleaning, transformation.

- **NumPy** → Numerical operations, linear algebra, GPU support.
- 

## 2. Data Visualization

### Data visualization

- Understand and visualize data structure

Examples:



**Matplotlib:** Customizable plots and interactive visualizations

**Seaborn:** Attractive statistical graphics

**ggplot2:** Building and adding elements in layers

**Tableau:** Interactive data visualization dashboards

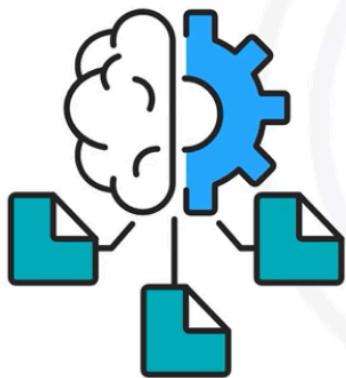
- **Matplotlib** → Foundational library for plots.
  - **Seaborn** → High-level statistical visualizations (built on Matplotlib).
  - **ggplot2 (R)** → Layered data visualization.
  - **Tableau** → Business intelligence dashboards.
-

### 3. Machine Learning Libraries

#### Machine learning ecosystem

- Create and tune ML models

Examples:



**NumPy:** Numerical computations on data arrays

**Pandas:** Data analysis, visualization, cleaning, and preparation

**SciPy:** Computing for optimization, integration, and linear regression

**Scikit-learn:** Suite of classification, regression, clustering, and dimensionality reduction

- **NumPy** → Core math & arrays.
  - **Pandas** → Data wrangling & prep.
  - **SciPy** → Scientific computing (optimization, regression).
  - **Scikit-learn** → Classical ML models (classification, regression, clustering, dimensionality reduction).
-

## 4. 🧠 Deep Learning

### Deep learning

- Designing, training, and testing neural network-based models



Examples:

**TensorFlow:** Numerical computing and large-scale ML

**Keras:** Implementing neural networks

**Theano:** Defining, optimizing, and evaluating mathematical expressions involving arrays

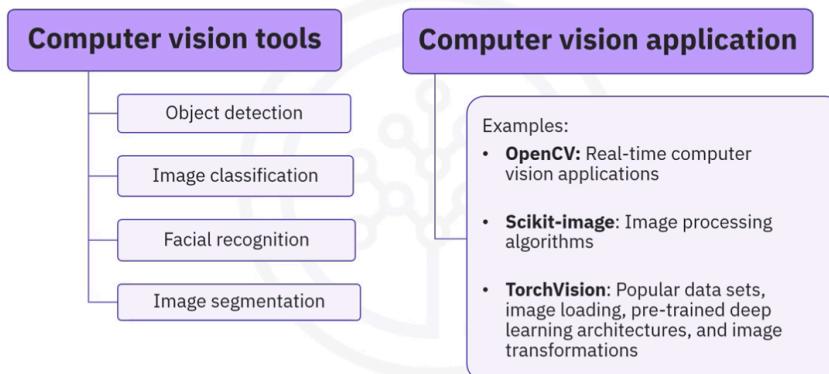
**PyTorch:** Computer vision, NLP, and experimentation

- **TensorFlow** → Large-scale ML & numerical computing.
- **Keras** → Easy-to-use neural networks.
- **Theano** → Mathematical optimization (legacy).
- **PyTorch** → Flexible, widely used for DL, CV, NLP.

---

## 5. 👁️ Computer Vision

### Types of machine learning tools



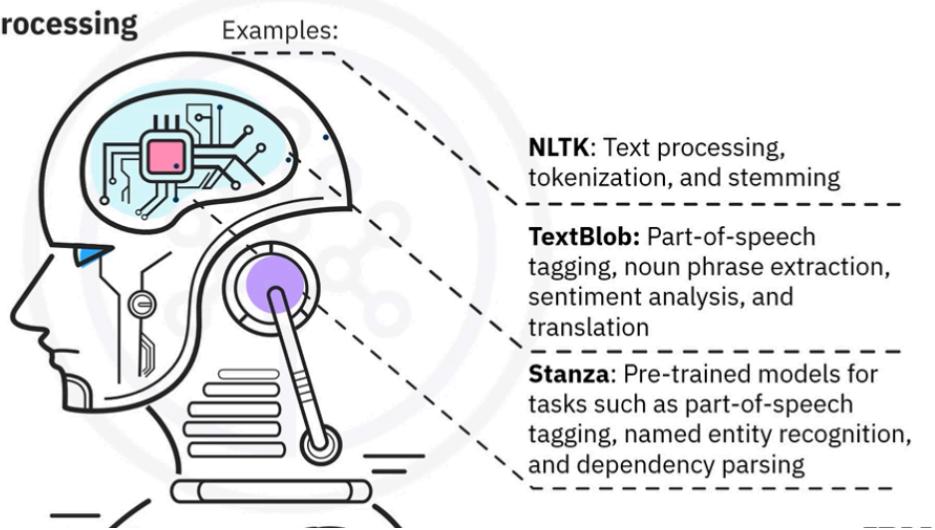
- **OpenCV** → Real-time CV apps (object detection, AR).
  - **Scikit-Image** → Image processing algorithms (filters, segmentation).
  - **TorchVision** → Pre-trained models, datasets, transformations.
- 

## 6. Natural Language Processing (NLP)

### Types of machine learning tools

#### Natural language processing

- Help build applications that understand, interpret, and generate human language



- **NLTK** → Text processing (tokenization, stemming).
  - **TextBlob** → Easy tasks (POS tagging, sentiment analysis, translation).
  - **Stanza (Stanford)** → Pre-trained models (NER, parsing).
-

## 7. 🎨 Generative AI

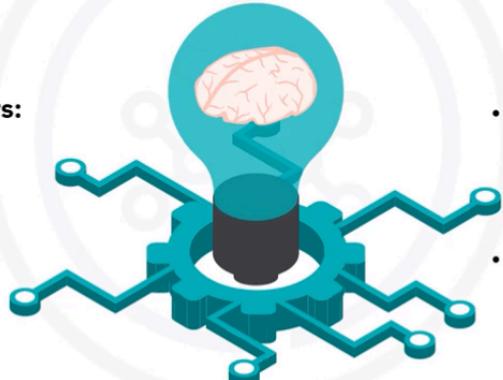
# Generative AI

## Generative AI

- Leverage AI to generate new content based on input data

Examples:

- **Hugging face transformers:** Text generation, language translation, and sentiment analysis
- **ChatGPT:** Text generation, building chatbots, etc



- **DALL-E:** Generating images from text descriptions
- **PyTorch:** Uses deep learning to create generative models such as GANs and transformers

- **Hugging Face Transformers** → NLP transformer models (text gen, translation, sentiment).
- **ChatGPT** → Conversational AI, text generation.
- **DALL·E** → Image generation from text.
- **PyTorch (GANs, Transformers)** → Generative models for text & images.

## ✓ Key Takeaways

## Recap

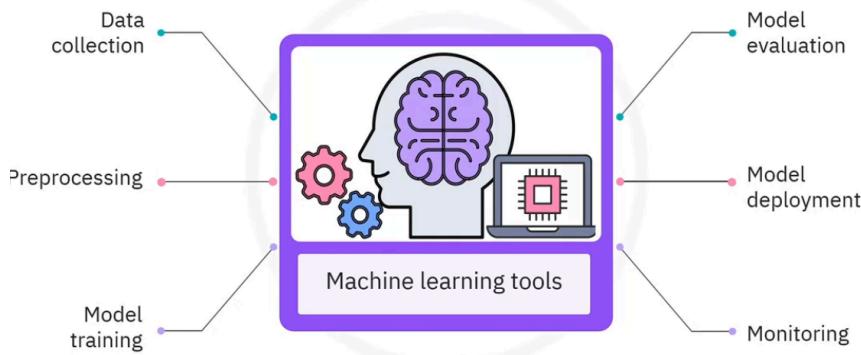
- Data is central to an ML algorithm
- ML tools simplify tasks and provide functionalities for ML pipelines
- ML programming language is used to build ML models and decode data patterns
- Common languages: Python, R, Julia, Scala, Java, and JavaScript

- Data is the **fuel** of ML.
  - Tools exist for every stage: **data processing** → **visualisation** → **ML/DL** → **CV/NLP** → **Generative AI**.  
Tools:
    - **Data processing and analytics:** PostgreSQL, Hadoop, Spark, Apache Kafka, pandas, and NumPy
    - **Data visualization:** Matplotlib, Seaborn, ggplot2, and Tableau
    - **Machine learning:** NumPy, Pandas, SciPy, and scikit-learn
    - **Deep learning:** TensorFlow, Keras, Theano, and PyTorch
    - **Computer vision:** OpenCV, scikit-image, and TorchVision
    - **NLP:** NLTK, TextBlob, and Stanza
    - **Generative AI:** Hugging face transformers, ChatGPT, DALL-E, and PyTorch
  - **Python dominates** due to ecosystem; other languages (R, Julia, Scala, Java, JS) also important.
- 

## Machine Learning Ecosystem & Scikit-Learn

### ◆ What is the ML Ecosystem?

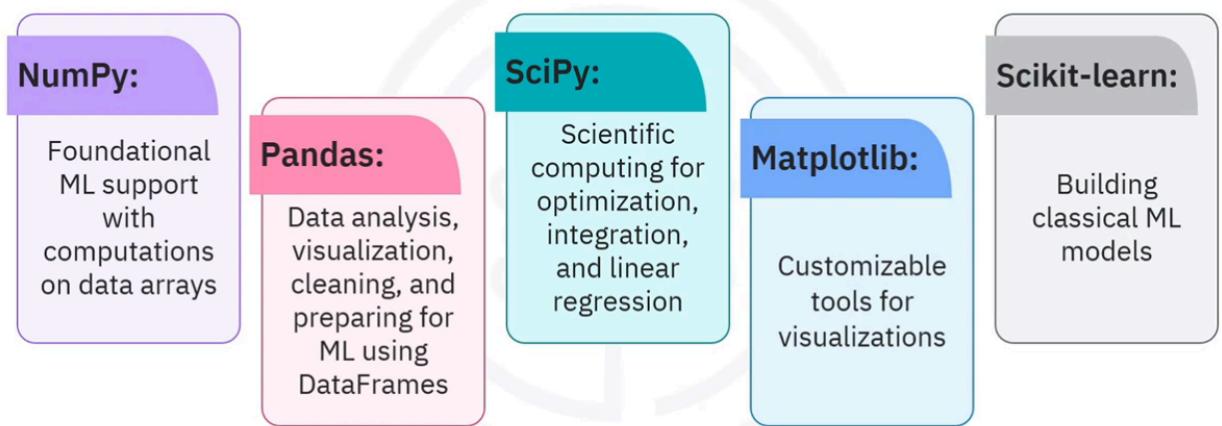
#### Machine learning ecosystem



- **Definition:** The interconnected **tools, frameworks, libraries, platforms, and processes** that support **developing, deploying, and managing ML models**.
  - **Covers:**
    - Data collection
    - Preprocessing
    - Model training
    - Model evaluation
    - Deployment & monitoring
- 

#### ◆ **Python ML Ecosystem – Core Libraries**

## Python tools and libraries



- **NumPy** → Efficient numerical computations on multidimensional arrays.
- **Pandas** → Data analysis, cleaning, wrangling (DataFrames).

- **SciPy** → Scientific computing (optimization, integration, linear regression).
  - **Matplotlib** → Visualization (highly customizable).
  - **Scikit-learn** → Classical ML models + pipeline utilities.
- 

### ◆ Scikit-learn Overview

## Scikit-learn

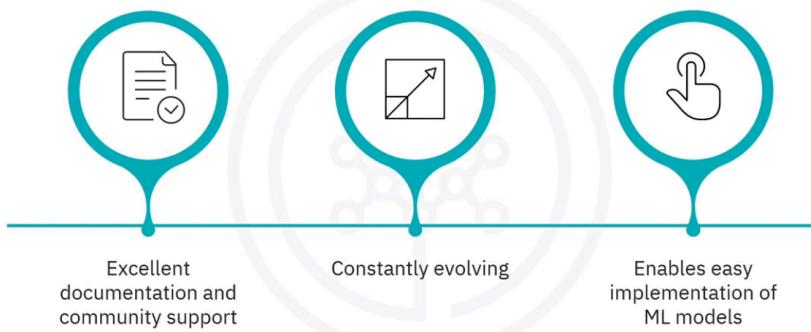
---



- **Free, open-source ML library** for Python.

## Scikit-learn

---



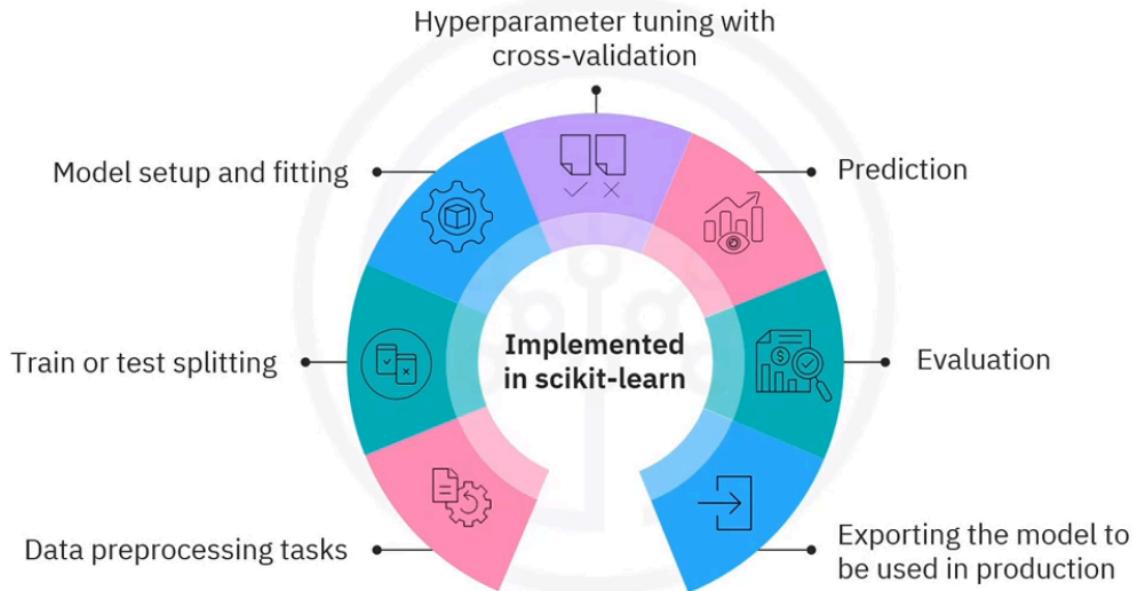
- Built on **NumPy**, **SciPy**, **Matplotlib**.
  - Supports:
    - **Classification**
    - **Regression**
    - **Clustering**
    - **Dimensionality reduction**
  - Features:
    - Easy to use (few lines of code).
    - Excellent documentation & community support.
    - Most pipeline tasks already implemented.
    - Constantly evolving (thousands of contributors).
- 

## ◆ **Scikit-learn Features in ML Pipeline**

- ✓ **Data Preprocessing** → cleaning, scaling, feature selection, feature extraction.
- ✓ **Data Splitting** → train/test split in one line.
- ✓ **Model Setup & Training** → classification, regression, etc.
- ✓ **Hyperparameter Tuning** → cross-validation.
- ✓ **Prediction & Evaluation** → accuracy, confusion matrix, etc.

✓ Model Exporting → save as pickle, reuse later.

## Machine learning pipeline tasks



### ◆ Example Workflow in Scikit-learn

#### Preprocessing

## Scikit-learn workflow

Given: Data set **X** and target variable **y** as NumPy arrays

```
from sklearn import preprocessing  
X = preprocessing.StandardScaler().fit(X).transform(X)
```

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33)
```

```
from sklearn import svm  
clf = svm.SVC(gamma=0.001, C=100.)
```

```
from sklearn.preprocessing import StandardScaler
```

```
X_scaled = StandardScaler().fit_transform(X)
```

1.

### Train-Test Split

```
from sklearn.model_selection import train_test_split  
  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33)
```

## Scikit-learn workflow

Given: **Data set X** and target **variable y** as NumPy arrays

```
clf.fit(X_train, y_train)  
  
clf.predict(X_test)  
  
yhat = clf.predict(X_test)  
  
from sklearn.metrics import confusion_matrix  
print(confusion_matrix(y_test, yhat, labels=[1,0]))  
  
import pickle  
s = pickle.dumps(clf)
```

2.

### Model Setup

```
from sklearn.svm import SVC  
  
clf = SVC(gamma=0.001, C=100.)
```

3.

### Training

```
clf.fit(X_train, y_train)
```

4.

## Prediction

```
y_pred = clf.predict(X_test)
```

5.

## Evaluation

```
from sklearn.metrics import confusion_matrix  
  
print(confusion_matrix(y_test, y_pred))
```

6.

## Save Model

```
import pickle  
  
pickle.dump(clf, open("model.pkl", "wb"))
```

7.

---

## ✓ Key Takeaways

## Recap

---

- ML ecosystem: Interconnected tools, frameworks, libraries, platforms, and processes that support ML models
- Python tools and libraries used: NumPy, pandas, SciPy, Matplotlib, and scikit-learn
- Scikit-learn:
  - Free ML library for Python
  - ML pipeline tasks are already implemented
- ML workflow using the scikit-learn library
- ML ecosystem = full pipeline (tools, frameworks, processes).

- Core Python ML stack: **NumPy, Pandas, SciPy, Matplotlib, Scikit-learn**.
  - **Scikit-learn** = easy, powerful, community-supported library with most ML tasks built-in.
  - ML workflow = **Preprocessing → Splitting → Training → Prediction → Evaluation → Deployment**.
- 

## Module 1 Summary and Highlights

Congratulations! You have completed this lesson. At this point in the course, you know that:

- Artificial intelligence (AI) simulates human cognition, while machine learning (ML) uses algorithms and requires feature engineering to learn from data.
- Machine learning includes different types of models: supervised learning, which uses labeled data to make predictions; unsupervised learning, which finds patterns in unlabeled data; and semi-supervised learning, which trains on a small subset of labeled data.
- Key factors for choosing a machine learning technique include the type of problem to be solved, the available data, available resources, and the desired outcome.
- Machine learning techniques include anomaly detection for identifying unusual cases like fraud, classification for categorizing new data, regression for predicting continuous values, and clustering for grouping similar data points without labels.
- Machine learning tools support pipelines with modules for data preprocessing, model building, evaluation, optimization, and deployment.
- R is commonly used in machine learning for statistical analysis and data exploration, while Python offers a vast array of libraries for different machine learning tasks. Other programming languages used in ML include Julia, Scala, Java, and JavaScript, each suited to specific applications like high-performance computing and web-based ML models.

- Data visualization tools such as Matplotlib and Seaborn create customizable plots, ggplot2 enables building graphics in layers, and Tableau provides interactive data dashboards.
- Python libraries commonly used in machine learning include NumPy for numerical computations, Pandas for data analysis and preparation, SciPy for scientific computing, and Scikit-learn for building traditional machine learning models.
- Deep learning frameworks such as TensorFlow, Keras, Theano, and PyTorch support the design, training, and testing of neural networks used in areas like computer vision and natural language processing.
- Computer vision tools enable applications like object detection, image classification, and facial recognition, while natural language processing (NLP) tools like NLTK, TextBlob, and Stanza facilitate text processing, sentiment analysis, and language parsing.
- Generative AI tools use artificial intelligence to create new content, including text, images, music, and other media, based on input data or prompts.
- Scikit-learn provides a range of functions, including classification, regression, clustering, data preprocessing, model evaluation, and exporting models for production use.
- The machine learning ecosystem includes a network of tools, frameworks, libraries, platforms, and processes that collectively support the development and management of machine learning models.

**Completed**

