# DATA SCIENCE - ML - DL - NLP

Machine learning is a branch of Artificial Intelligence that focuses on developing models and algorithms that let computers learn from data without being explicitly programmed for every task. In simple words, ML teaches the systems to think and understand like humans by learning from the data.

—--------------------------------------------------------------------------------------------------------------------------------------

In the field of machine learning, data plays a pivotal role in training models to make accurate predictions and decisions. Two fundamental types of data are labelled and unlabeled data, each serving distinct purposes in the learning process. Understanding the difference between these two types of data is essential for leveraging them effectively in machine learning applications.

## What is Labeled Data?

Labelled data is data that has been assigned a label or category, indicating the ground truth or correct classification for each data point. This labelling is typically done by human annotators and is crucial for supervised learning tasks. In supervised learning, the model learns from labelled examples to make predictions on new, unseen data. Examples of labelled data include:

- A dataset of images with labels indicating whether each image contains a cat or a dog.
- An email dataset labelled as spam or not spam.
- A dataset of customer reviews labelled with sentiment (positive, negative, neutral).

Labelled data is valuable for training models for tasks such as classification, regression, and object detection, where the goal is to predict a specific label or value for each data point. However, obtaining labelled data can be expensive and time-consuming, as it requires human annotators to assign labels to each data point.

## What is Unlabeled Data?

Unlabeled data, on the other hand, is data that does not have any labels or categories assigned to it. The true classification or category of each data point is unknown, making unlabeled data suitable for unsupervised learning tasks. In unsupervised learning, the model must learn from the inherent structure of the data to uncover patterns or anomalies. Examples of unlabeled data include:

- A dataset of customer transactions without any labels indicating fraudulent or non-fraudulent transactions.
- A collection of text documents without any labels indicating the topic or category of each document.
- An image dataset without any labels indicating the content or objects in each image.

Unlabeled data is used in unsupervised learning tasks such as clustering, dimensionality reduction, and anomaly detection, where the goal is to find patterns or anomalies in the data without the use of labeled examples. Unlabeled data is often easier to obtain and can be generated or collected in large quantities without the need for labeling.
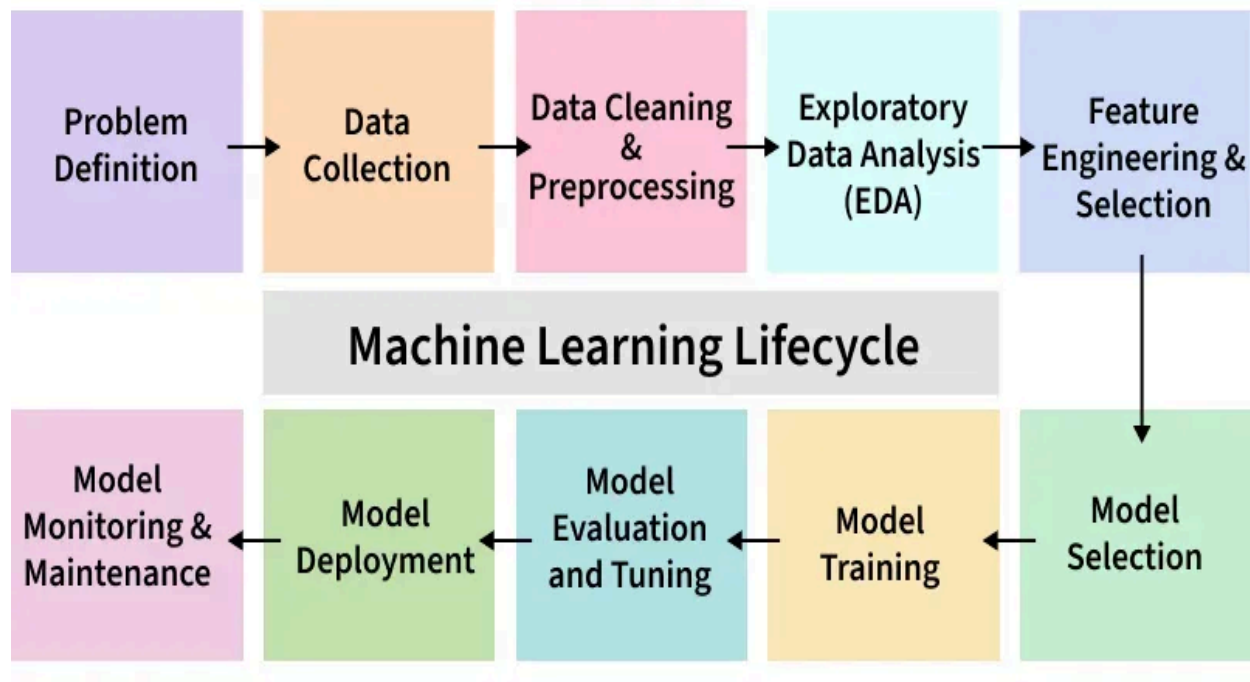
❖ Machine Learning is mainly divided into three core types: Supervised, Unsupervised and Reinforcement Learning along with two additional types, Semi-Supervised and Self-Supervised Learning.
  ● **Supervised Learning:** Trains models on labeled data to predict or classify new, unseen data.
  ● **Unsupervised Learning**: Finds patterns or groups in unlabeled data, like clustering or dimensionality reduction.
  ● **Reinforcement Learning**: Learns through trial and error to maximize rewards, ideal for decision-making tasks.

*Note: The following are not part of the original three core types of ML, but they have become increasingly important in real-world applications, especially in deep learning.*

***Additional Types***:

  ● ***Self-Supervised Learning****: Self-supervised learning is often considered as a subset of unsupervised learning, but it has grown into its own field due to its success in training large-scale models. It generates its own labels from the data, without any manual labeling.*
  ● ***Semi-Supervised Learning:*** *This approach combines a small amount of labeled data with a large amount of unlabeled data. It's useful when labeling data is expensive or time-consuming.*

—-------------------------------------------------------------------------------------------------------------

------------**Machine Learning Life Cycle:** Machine Learning Lifecycle is a structured process that defines how machine learning (ML) models are developed, deployed and maintained. It consists of a series of steps that ensure the model is accurate, reliable and scalable.

Machine Learning Lifecycle

**Step 1: Problem Definition**
The first step is identifying and clearly defining the business problem. A well-framed problem provides the foundation for the entire lifecycle. Important things like project objectives, desired outcomes and the scope of the task are carefully designed during this stage.

- Collaborate with stakeholders to understand business goals
- Define project objectives, scope and success criteria
- Ensure clarity in desired outcomes

**Step 2: Data Collection**
Data Collection phase involves systematic collection of datasets that can be used as raw data to train models. The quality and variety of data directly affect the model's performance.

Here are some basic features of Data Collection:

- **Relevance:** Collect data should be relevant to the defined problem and include necessary features.
- **Quality:** Ensure data quality by considering factors like accuracy and ethical use.
- **Quantity:** Gather sufficient data volume to train a robust model.
- **Diversity:** Include diverse datasets to capture a broad range of scenarios and patterns.

**Step 3: Data Cleaning and Preprocessing**
Raw data is often messy and unstructured and if we use this data directly to train then it can lead to poor accuracy. We need to do data cleaning and preprocessing which often involves:

- **Data Cleaning:** Address issues such as missing values, outliers and inconsistencies in the data.
- **Data Preprocessing:** Standardize formats, scale values, and encode categorical variables for consistency.
- **Data Quality:** Ensure that the data is well-organized and prepared for meaningful analysis.

**Step 4: Exploratory Data Analysis (EDA)**
To find patterns and characteristics hidden in the data, Exploratory Data Analysis (EDA) is used to uncover insights and understand the dataset's structure. During EDA patterns, trends and insights are provided which may not be visible by naked eyes. This valuable insight can be used to make informed decisions.

Here are the basic features of Exploratory Data Analysis:

- **Exploration:** Use statistical and visual tools to explore patterns in data.
- **Patterns and Trends:** Identify underlying patterns, trends and potential challenges within the dataset.
- **Insights:** Gain valuable insights for informed decision making in later stages.
- **Decision Making:** Use EDA for feature engineering and model selection.

**Step 5: Feature Engineering and Selection**
Feature engineering and selection is a transformative process that involves selecting only relevant features to enhance model efficiency and prediction while reducing complexity.

Here are the basic features of Feature Engineering and Selection:

- **Feature Engineering:** Create new features or transform existing ones to capture better patterns and relationships.
- **Feature Selection:** Identify subset of features that most significantly impact the model's performance.
- **Domain Expertise:** Use domain knowledge to engineer features that contribute meaningfully for prediction.
- **Optimization:** Balance set of features for accuracy while minimizing computational complexity.

**Step 6: Model Selection**

For a good machine learning model, model selection is a very important part as we need to find a model that aligns with our defined problem, nature of the data, complexity of problem and the desired outcomes.

Here are the basic features of Model Selection:

- **Complexity:** Consider the complexity of the problem and the nature of the data when choosing a model.
- **Decision Factors:** Evaluate factors like performance, interpretability and scalability when selecting a model.
- **Experimentation:** Experiment with different models to find the best fit for the problem.

**Step 7: Model Training**
With the selected model the machine learning lifecycle moves to the model training process. This process involves exposing the model to historical data allowing it to learn patterns, relationships and dependencies within the dataset.

Here are the basic features of Model Training:

- **Iterative Process:** Train the model iteratively, adjusting parameters to minimize errors and enhance accuracy.
- **Optimization:** Fine-tune model to optimize its predictive capabilities.
- **Validation:** Rigorously train model to ensure accuracy to new unseen data.

**Step 8: Model Evaluation and Tuning**
Model evaluation involves rigorous testing against validation or test datasets to test accuracy of model on new unseen data. It provides insights into a model's strengths and weaknesses. If the model fails to achieve desired performance levels we may need to tune the model again and adjust its hyperparameters to enhance predictive accuracy.

Here are the basic features of Model Evaluation and Tuning:

- **Evaluation Metrics:** Use metrics like accuracy, precision, recall and F1 score to evaluate model performance.
- **Strengths and Weaknesses:** Identify the strengths and weaknesses of the model through rigorous testing.
- **Iterative Improvement:** Initiate model tuning to adjust hyperparameters and enhance predictive accuracy.
- **Model Robustness:** Iterative tuning to achieve desired levels of model robustness and reliability.

**Step 9: Model Deployment**

Now the model is ready for deployment for real-world application. It involves integrating the predictive model with existing systems allowing business to use this for informed decision-making.

Here are the basic features of Model Deployment:

- Integrate with existing systems
- Enable decision-making using predictions
- Ensure deployment scalability and security
- Provide APIs or pipelines for production use

**Step 10: Model Monitoring and Maintenance**

After deployment, models must be monitored to ensure they perform well over time. Regular tracking helps detect data drift, accuracy drops or changing patterns and retraining may be needed to keep the model reliable in real-world use.

Here are the basic features of Model Monitoring and Maintenance:

- Track model performance over time
- Detect data drift or concept drift
- Update and retrain the model when accuracy drops
- Maintain logs and alerts for real-time issues

Each step is essential for building a successful machine learning model that can provide valuable insights and predictions. By following the Machine learning lifecycle organizations we can solve complex problems.

—----------------------------------------------------------------------------------------------------------------------
------------