

UNIT 3 : INTRODUCTION TO MACHINE LEARNING

What is Machine Learning?

- Machine Learning, often abbreviated as ML, is a subset of artificial intelligence (AI) that focuses on the development of computer algorithms that improve automatically through experience and by the use of data. In simpler terms, machine learning enables computers to learn from data and make decisions or predictions without being explicitly programmed to do so.
- At its core, machine learning is all about creating and implementing algorithms that facilitate these decisions and predictions. These algorithms are designed to improve their performance over time, becoming more accurate and effective as they process more data.
- In traditional programming, a computer follows a set of predefined instructions to perform a task. However, in machine learning, the computer is given a set of examples (data) and a task to perform, but it's up to the computer to figure out how to accomplish the task based on the examples it's given.
- For instance, if we want a computer to recognize images of cats, we don't provide it with specific instructions on what a cat looks like. Instead, we give it thousands of images of cats and let the machine learning algorithm figure out the common patterns and features that define a cat. Over time, as the algorithm processes more images, it gets better at recognizing cats, even when presented with images it has never seen before.
- This ability to learn from data and improve over time makes machine learning incredibly powerful and versatile. It's the driving force behind many of the technological advancements we see today, from voice assistants and recommendation systems to self-driving cars and predictive analytics.



The Importance of Machine Learning:

In the 21st century, data is the new oil, and machine learning is the engine that powers this data-driven world. It is a critical technology in today's digital age, and its importance cannot be overstated. This is reflected in the industry's projected growth, with the US Bureau of Labor Statistics predicting a 26% growth in jobs between 2023 and 2033.

Here are some reasons why it's so essential in the modern world:

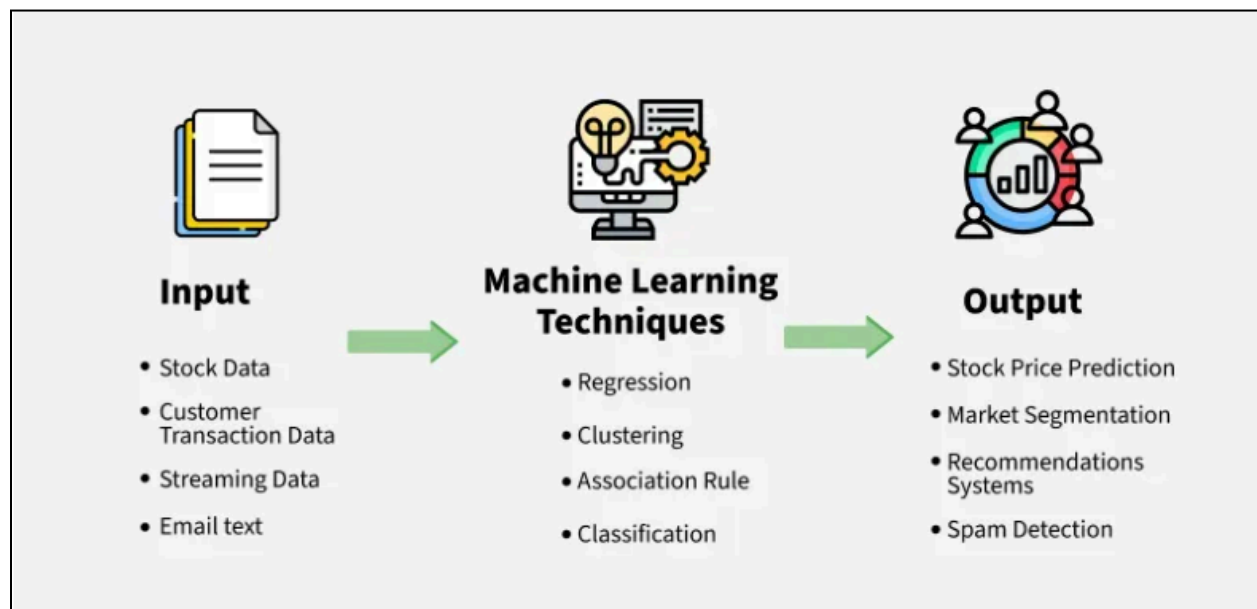
- **Data processing.** One of the primary reasons machine learning is so important is its ability to handle and make sense of large volumes of data. With the explosion of digital data from social media, sensors, and other sources, traditional data

analysis methods have become inadequate. Machine learning algorithms can process these vast amounts of data, uncover hidden patterns, and provide valuable insights that can drive decision-making.

- **Driving innovation.** Machine learning is driving innovation and efficiency across various sectors. Here are a few examples:
- **Healthcare.** Algorithms are used to predict disease outbreaks, personalize patient treatment plans, and improve medical imaging accuracy.
- **Finance.** Machine learning is used for credit scoring, algorithmic trading, and fraud detection.
- **Retail.** Recommendation systems, supply chains, and customer service can all benefit from machine learning.

The techniques used also find applications in sectors as diverse as agriculture, education, and entertainment.

- **Enabling automation.** Machine learning is a key enabler of automation. By learning from data and improving over time, machine learning algorithms can perform previously manual tasks, freeing humans to focus on more complex and creative tasks. This not only increases efficiency but also opens up new possibilities for innovation.



3.1 HISTORY AND EVOLUTION OF ML:

The Early Days of Machine Learning

- **Philosophical Foundations:**

The roots of machine learning can be traced back to early philosophical ideas. Aristotle introduced the concept of logical reasoning, suggesting that thought processes could follow structured rules, similar to mechanical systems. René Descartes later proposed that machines might replicate aspects of human thinking, hinting at the possibility of intelligent systems. These early ideas about reasoning and logic influenced the development of both AI and ML.

- **Early Computational Devices:**

The invention of early computational devices laid the foundation for machine learning. Charles Babbage's Analytical Engine and other early machines showcased the potential of devices capable of performing complex calculations. These systems paved the way for later developments in computing and inspired researchers to explore how machines could learn from data.

- **The Turing Test (1950):**

In 1950, Alan Turing introduced the Turing Test, which evaluated a machine's ability to exhibit human-like intelligence. While the Turing Test focused on AI, it had significant implications for machine learning. It suggested that machines could learn to respond intelligently to inputs, influencing future research into learning algorithms.

- **First Neural Network (1943):**

The first mathematical model of a neural network was introduced by Warren McCulloch and Walter Pitts in 1943. Their work demonstrated that neurons could be represented mathematically and that neural processes could be simulated by machines. Although limited, this model laid the groundwork for future advancements in neural networks and shaped early research in ML.

- **Milestones in Machine Learning Development (1950-2000):**

Computer Checkers (1952)

Arthur Samuel pioneered machine learning through his work on a computer-based checkers program in 1952. Samuel's program was designed to improve its gameplay by learning from past games, marking the first practical use of ML in gaming. This

development showcased the potential of machines to learn autonomously, setting a precedent for future ML applications.

- **The Perceptron (1957):**

In 1957, Frank Rosenblatt introduced the Perceptron, a single-layer neural network model capable of recognizing patterns. Rosenblatt's Perceptron generated significant excitement as it demonstrated how machines could learn from input data. However, its inability to solve non-linear problems (like XOR) exposed the model's limitations. This sparked debates about the potential of neural networks and temporarily slowed research in the field.

- **Nearest Neighbor Algorithm (1967):**

The development of the Nearest Neighbor algorithm in 1967 marked a significant step forward in pattern recognition. This algorithm allowed computers to classify data points based on their proximity to other points in a given dataset. It became a crucial tool for tasks like handwriting recognition and clustering, illustrating the growing potential of ML in real-world applications.

- **The Backpropagation Algorithm (1974):**

The introduction of the backpropagation algorithm in 1974 was a turning point for neural networks. Backpropagation allowed multi-layer networks to learn by correcting errors through feedback loops. This breakthrough revived interest in neural networks, laying the foundation for deep learning and enabling machines to solve more complex problems effectively.

- **The Stanford Cart (1979):**

The Stanford Cart was a groundbreaking project in the field of autonomous vehicles. Developed in 1979, the cart used ML algorithms to navigate obstacles in its environment without human intervention. This project demonstrated the potential of machine learning in robotics and inspired future research in autonomous systems, including self-driving cars.

- **AI Winter:**

Despite these milestones, the AI Winter—a period of reduced funding and enthusiasm—began in the 1970s and persisted into the 1990s. Early ML models struggled with limitations such as insufficient computing power and data availability. Skepticism about the practical use of ML and AI further contributed to the decline in research

activity. However, behind the scenes, researchers continued their work, laying the groundwork for future breakthroughs.

- **The Rise of Machine Learning (2000 – Present)**

- **Machine Defeats Man in Chess (1997)**

Although technically before 2000, IBM's Deep Blue made history by defeating Garry Kasparov, the reigning world chess champion, in 1997. This event showcased the power of machine learning algorithms in decision-making and pattern recognition. It proved that machines could compete with human intelligence, sparking renewed interest in artificial intelligence (AI) and ML.

- **The Torch Software Library (2002):**

The release of Torch, an open-source software library, marked a significant shift in the development of ML. Torch allowed researchers and developers to build machine learning models efficiently, driving community-driven innovation. It paved the way for other open-source frameworks like TensorFlow and PyTorch, making ML more accessible and accelerating research.

- **Deep Learning Breakthroughs (2006):**

In 2006, Geoffrey Hinton and his team introduced breakthroughs in deep learning, enabling neural networks to process large datasets more effectively. This advancement allowed for significant improvements in fields such as speech recognition and computer vision, solidifying deep learning as a powerful subset of ML.

- **Google Brain (2011)**

The launch of the Google Brain project applied machine learning to large-scale systems. Google used ML to improve services such as search engines and advertising platforms, demonstrating how ML could scale to handle enormous datasets. This project highlighted the role of ML in transforming industries through automation and efficiency.

- **DeepFace (2014):**

In 2014, Facebook introduced DeepFace, a facial recognition project that used deep learning to identify faces with high accuracy. This technology demonstrated the practical applications of ML in biometric security and image recognition. It showcased ML's potential to enhance authentication systems and laid the foundation for advancements in computer vision.

- **ImageNet Challenge (2017):**

The ImageNet Challenge became a benchmark for evaluating ML models in computer vision. In 2017, ML systems achieved human-level accuracy in recognizing objects, marking a major milestone in ML. The success of ImageNet highlighted the capabilities of convolutional neural networks (CNNs) and deep learning in advancing computer vision technologies.

- **Generative AI (2010s Onwards):**

The rise of generative AI models, such as GPT and DALL-E, transformed ML's role in creative fields. These models generate text, images, and even music, expanding ML's applications beyond analytics and predictions. Generative AI has opened new possibilities in fields like content creation, design, and entertainment.

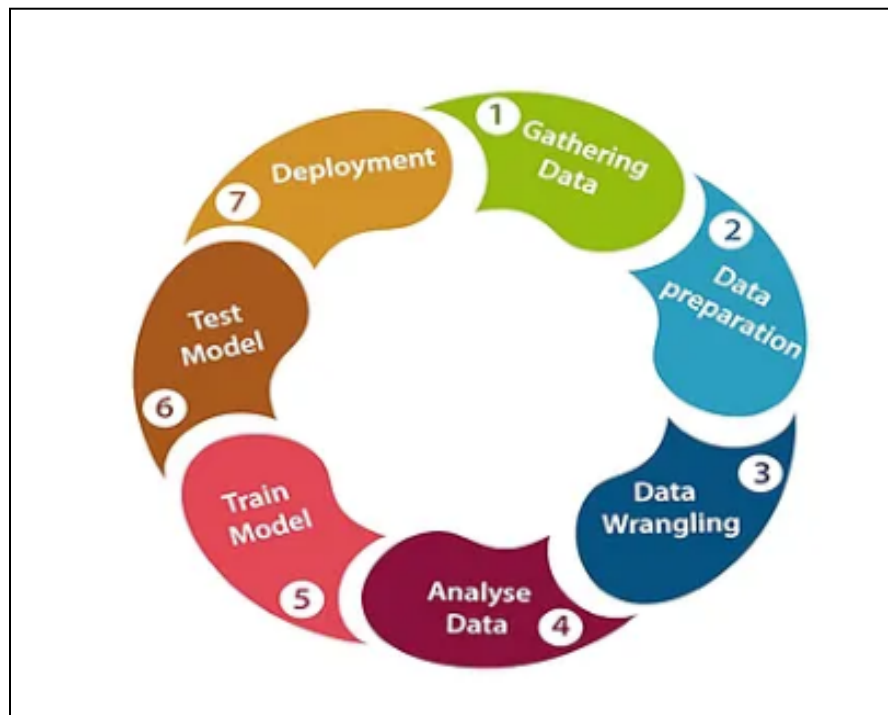
3.2 MACHINE LEARNING LIFE CYCLE:

Machine learning is about data – no lie there. There's no machine learning without a decent amount of data for the machine to learn from. The amount of available data is growing exponentially, which makes machine learning development easier than ever.

The connection between machine learning and algorithms is also on point. Indeed, there are complex mathematical methods that force machines to learn. No math – no machine learning.

Lastly, model training and data preparation is indeed the core of every ML project. Machine learning engineers spend a substantial amount of time training models and preparing datasets. That's why it's the first thing ML engineers think of.

Let's look at the steps in a flow on a very high level:



Let's dive into each step to see what it entails and the tasks involved.

1. Gathering Data: Data is the fuel that powers machine learning models. The first step in any ML project is gathering data — and lots of it!

What this step includes:

- **Identifying data sources:** Where will your data come from? This could be public datasets, company records, web scraping, or even sensor data from IoT devices.
- **Collecting data:** This can be done manually (e.g., surveys) or automatically (e.g., APIs, web scraping tools, databases).

Why it matters:

- Without high-quality data, your model will struggle to produce meaningful results. More data usually means better accuracy, so this step is critical.

2. Data Preparation: Once you've gathered your data, the next step is to prepare it for analysis. Data preparation involves cleaning, formatting, and ensuring the data is structured in a way that can be used by your machine learning algorithms.

What this step includes:

- **Removing duplicates:** Clean out any redundant records that can skew your model.
- **Handling missing data:** Decide how to deal with missing values — either by filling them in, or excluding those data points entirely.
- **Formatting data:** Ensure the data is in a consistent format, especially when combining data from multiple sources (e.g., numbers, dates).

Why it matters:

- Well-prepared data helps reduce noise and errors in the model, ensuring more accurate predictions later in the pipeline.

3. Data Wrangling: Now that your data is clean and well-prepared, the next step is data wrangling. This process involves transforming the data into a more suitable format for analysis. Think of it as shaping your raw data into something that can be directly used by machine learning algorithms.

What this step includes:

- **Feature selection:** Picking out the most important variables (features) that will have the greatest impact on your model.
- **Feature engineering:** Creating new features from existing data to improve the model's accuracy (e.g., calculating age from a birth date).
- **Normalization or scaling:** Adjusting your data to fit within a specific range, especially when using algorithms like neural networks or SVMs that are sensitive to large variations.

Why it matters:

- Effective data wrangling allows your model to focus on the most relevant and accurate information, leading to better results. Poor wrangling can lead to inaccurate predictions or longer training times.

4. Analyzing Data: Before you jump into model building, you need to understand your data's patterns and relationships. This is where data analysis comes into play.

What this step includes:

- **Exploratory Data Analysis (EDA):** Using visualizations (e.g., histograms, scatter plots) and statistical measures (e.g., correlation matrices) to understand the relationships and trends in your data.
- **Detecting outliers:** Identifying and possibly removing outliers that could negatively affect your model's performance.
- **Understanding distributions:** Looking at how the data is distributed (e.g., normal, skewed) and if transformations are necessary.

Why it matters:

- This step helps you discover insights and decide how to approach the problem. It's also your last chance to make changes before diving into model training.

5. Training the Model: Now comes the fun part — training your model. In this step, you'll select the machine learning algorithm that best fits your problem and train it using your cleaned and prepared data.

What this step includes:

- **Choosing an algorithm:** Based on the type of problem (regression, classification, clustering), you'll choose a suitable algorithm like Linear Regression, Decision Trees, or K-Nearest Neighbors.
- **Training the model:** Feeding the prepared data into the algorithm to teach it to make predictions.
- **Adjusting parameters:** Fine-tuning the algorithm's hyperparameters (e.g., learning rate, number of layers in a neural network) to improve performance.

Why it matters:

- A well-trained model is the core of any machine learning solution. It's where the model learns to make predictions based on your data.

6. Testing the Model: After training, it's crucial to test your model to see how well it performs on unseen data. This is where model evaluation comes in.

What this step includes:

- **Splitting data:** Dividing your dataset into training and testing sets (typically 80/20 or 70/30).
- **Evaluating performance:** Using metrics like accuracy, precision, recall, and F1-score to assess how well the model performs on the test set.
- **Cross-validation:** Running multiple rounds of testing to ensure the model generalizes well to new data.

Why it matters:

Testing ensures that your model isn't just memorizing the training data (overfitting) and can make accurate predictions on new, unseen data.

7. Deployment: Finally, after training and testing your model, it's time to deploy it. This is where your machine learning model goes live and starts making predictions in real-time environments.

What this step includes:

- **Choosing a platform:** Deciding where to deploy your model (e.g., cloud services like AWS, Google Cloud, or on-premise servers).
- **Setting up APIs:** Creating endpoints so the model can interact with other systems and provide predictions.
- **Monitoring performance:** Continuously monitoring the model's predictions and retraining it as new data becomes available.

Why it matters:

- Deployment is where your model starts delivering real value. It's essential to ensure that the model works efficiently in a production environment and is scalable.

3.3 DIFFERENT FORMS OF DATA:

3.3.1 DATA MINING:

- Data mining is the process of extracting insights from large datasets using statistical and computational techniques. It can involve structured, semi-structured or unstructured data stored in databases, data warehouses or data lakes. The goal is to uncover hidden patterns and relationships to support informed decision-making and predictions using methods like clustering, classification, regression and anomaly detection.

3.3.2 STATISTICS:

- Statistics is the discipline that concerns the collection, organization, analysis, interpretation, and presentation of data. In applying statistics to a scientific, industrial, or social problem, it is conventional to begin with a statistical population or a statistical model to be studied. Statistics deals with every aspect of data, including the planning of data collection in terms of the design of surveys and experiments.

3.3.3 DATA ANALYTICS:

- Data analytics is the practice of looking at raw data in various ways to gain information. Many of the techniques and processes of data analytics have been automated into mechanical processes and algorithms that work over raw data for human consumption. Data analytics can be used to optimize the performance of a business or help a decision-maker come to the right call based on underlying information.

3.3.4 STATISTICS DATA:

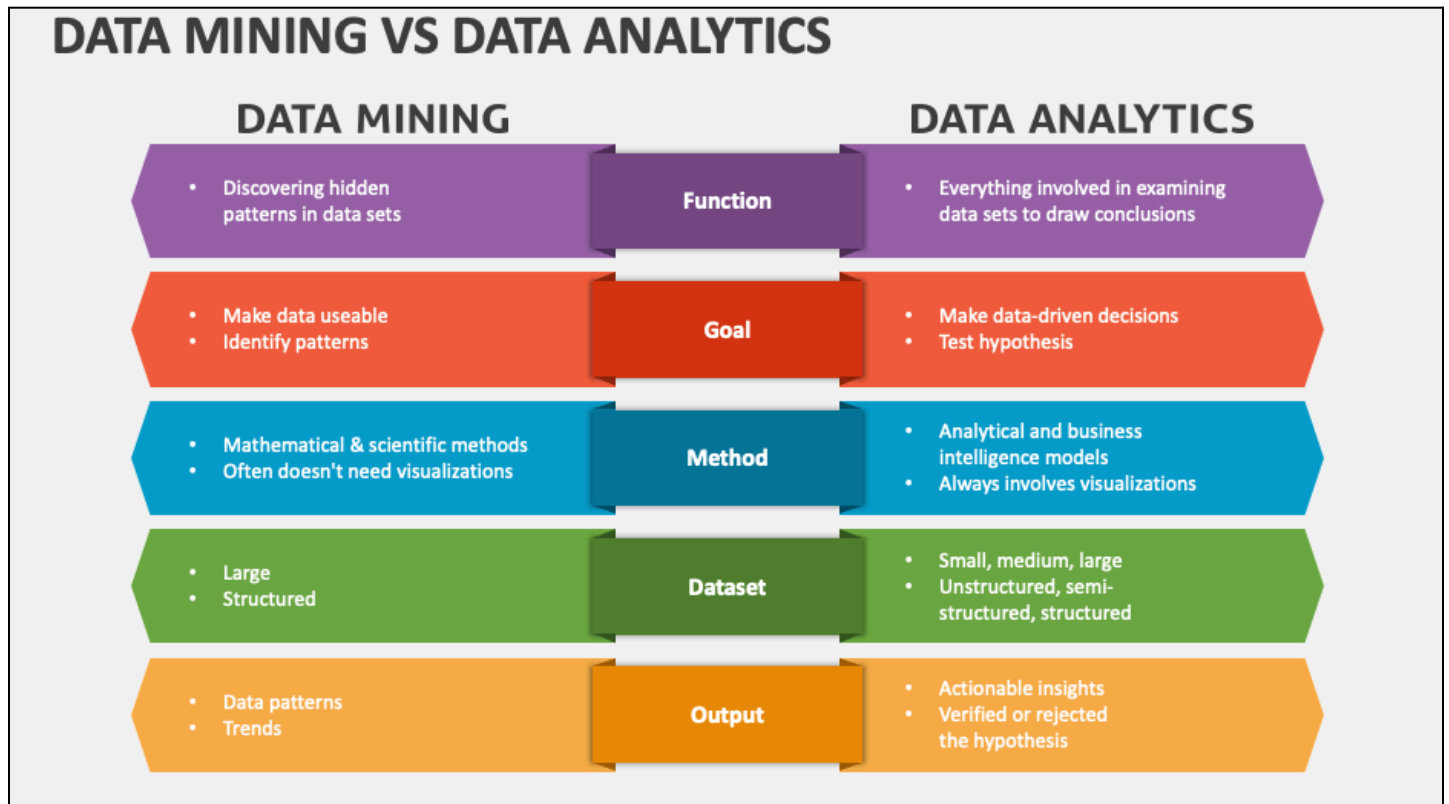
- Statistical data is data regularly collected through surveys or administrative sources as a part of the national statistical system and computed and analysed for the production of official statistics. Statistical data includes census data and surveys such as the Demographic Health Survey (DHS) and Multiple Indicator Cluster Surveys (MICS).
- Examples include height, weight, temperature, and time. Continuous data are used for measurements and observations, and they can be analyzed using mean and median, as well as continuous probability distributions like the normal distribution.

3.3.5 STATISTICS vs DATA MINING vs DATA ANALYTICS:

Difference between Data Mining And Statistics:

Aspect	Data Mining	Statistics
Definition	The process of discovering patterns, relationships, and insights from large datasets using algorithms and computational techniques.	The science of collecting, analyzing, interpreting, presenting, and organizing numerical data.
Focus	Focuses on automating the discovery of patterns and trends in large datasets.	Focuses on hypothesis testing and making inferences based on smaller samples
Techniques Used	Machine learning, clustering, classification, association rule mining, etc.	Regression analysis, hypothesis testing, probability theory, and descriptive statistics.
Nature	Applied and practical, often using tools and algorithms for large-scale analysis.	Theoretical, rooted in mathematical principles for data interpretation.
Data Size	Typically works with very large datasets (big data).	Can work with both small and large datasets but often deals with smaller samples.

Difference between Data Mining And Data Analytics:



3.4 DATASET FOR ML:

What is a Dataset?

- A Dataset is a set of data grouped into a collection with which developers can work to meet their goals. In a dataset, the rows represent the number of data points and the columns represent the features of the Dataset. They are mostly used in fields like machine learning, business, and government to gain insights, make informed decisions, or train algorithms. Datasets may vary in size and complexity and they mostly require cleaning and preprocessing to ensure data quality and suitability for analysis or modeling.

S.ID	Geography	History	Arabic	Philosophy	Physics	Geometry	Algebray	Chemistry
10001	90	77	73	81	85	83	92	91
10002	89	80	78	82	90	90	80	93
10003	89	78	80	85	85	90	91	92
10004	79	88	91	79	89	86	89	79
10005	89	69	74	85	87	90	91	92
10006	91	76	72	80	84	82	91	90
10007	69	80	77	68	80	90	72	88
10008	92	79	79	83	86	87	90	92
10009	80	88	90	78	79	90	87	90
10010	91	88	79	80	79	78	87	89

Understanding Data Types: Numerical, Categorical, and Ordinal:

Data comes in various forms, each type holding distinct characteristics and implications for analysis. By categorizing data into numerical, categorical, and ordinal types, we gain valuable insights into how to interpret and process information.

1. Numerical Data:

Numerical data represents quantitative measurements and can be further classified into discrete and continuous types.

- Discrete Data:

Discrete numerical data consists of distinct, separate values often in integer form, representing counts or occurrences. In Python, consider a dataset recording the number of goals scored in soccer matches:

```
goals_scored = [2, 1, 3, 0, 2, 4]
```

- Continuous Data

Get Navneet Singh's stories in your inbox

Join Medium for free to get updates from this writer.

Continuous numerical data encompasses an infinite range of possible values. Imagine a dataset tracking daily temperatures:

```
daily_temperatures = [23.5, 25.1, 22.8, 24.6, 26.3]
```

In this case, the data points represent continuous measurements, allowing for decimals and indicating the temperature values within a range.

2. Categorical Data:

Categorical data, unlike numerical data, lacks inherent mathematical meaning and represents qualitative attributes or labels. Examples include gender, product categories, or political affiliations. In Python, a dataset capturing car colors might look like:

```
car_colors = ['red', 'blue', 'black', 'white', 'silver']
```

Each entry represents a category — color in this instance — without numerical significance beyond differentiation.

3. Ordinal Data:

Ordinal data shares characteristics of both numerical and categorical data, where categories possess a specific order or ranking.

Example: Movie Ratings

Consider a dataset with movie ratings on a 1–5 scale:

```
movie_ratings = [3, 5, 2, 4, 1]
```

While these ratings are numeric, they hold categorical meaning through their ordered nature. A rating of 1 signifies a lower evaluation compared to a rating of 5, showcasing a clear ordinal relationship.

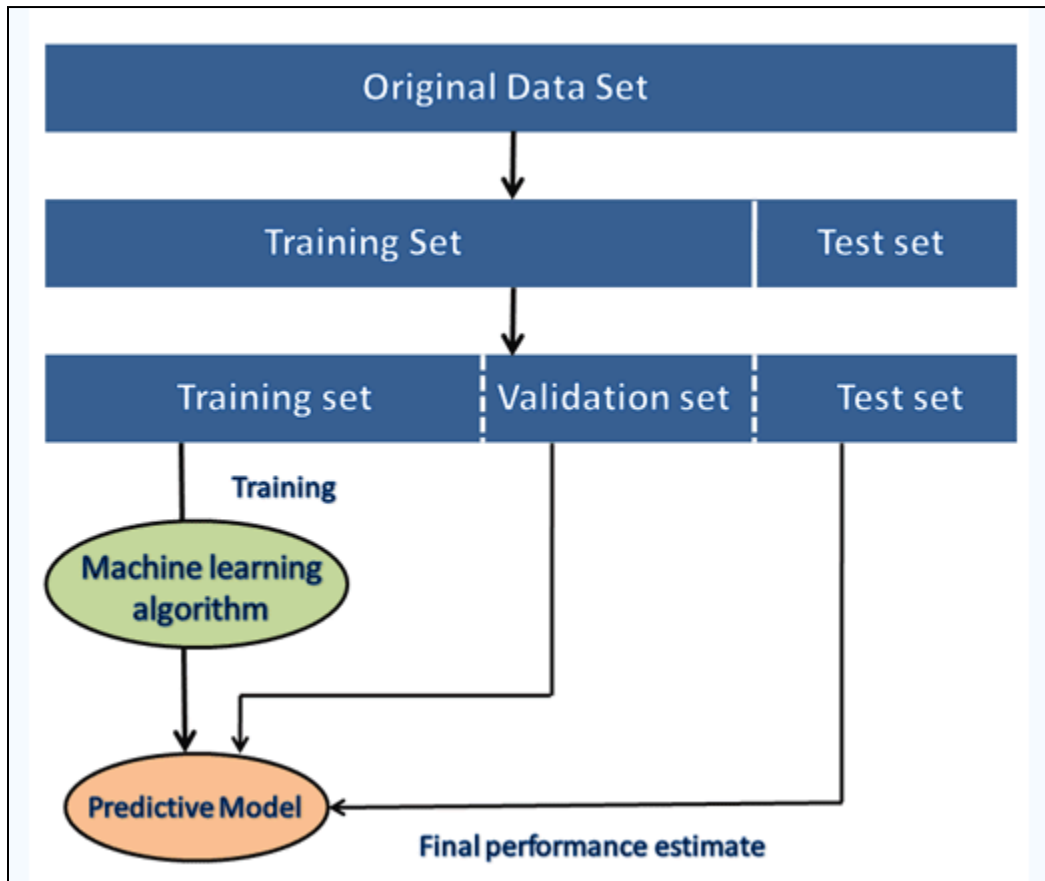
- **Need of Dataset:**

To work with machine learning projects, we need a huge amount of data, because, without the data, one cannot train ML/AI models. Collecting and preparing the dataset is one of the most crucial parts while creating an ML/AI project.

The technology applied behind any ML projects cannot work properly if the dataset is not well prepared and pre-processed.

During the development of the ML project, the developers completely rely on the datasets. In building ML applications, datasets are divided into two parts:

1. **Training dataset**
2. **Test Dataset**



3.4.1 TRAINING DATA SET:

There are two key types of data used for machine learning training and testing data. They each have a specific function to perform when building and evaluating machine learning models. Machine learning algorithms are used to learn from data in datasets. They discover patterns and gain knowledge. make choices, and examine those decisions.

What is Training data?

Testing data is used to determine the performance of the trained model, whereas training data is used to train the machine learning model. Training data is the power that supplies the model in machine learning, it is larger than testing data. Because more data helps to more effective predictive models. When a machine learning algorithm receives data from our records, it recognizes patterns and creates a decision-making model.

Algorithms allow a company's past experience to be used to make decisions. It analyzes all previous cases and their results and, using this data, creates models to score and predict the outcome of current cases. The more data ML models have access to, the more reliable their predictions get over time.

- Training data preparation:

Data is all around us. The global population generates immense amounts of data every second of the day. But raw data is typically not useful for model training. Quality assurance is critical. First, data must be pre-processed through a multi-step data pipeline. This can be an involved process for data scientists, comprising a large portion of the scope of a machine learning project, requiring sophisticated data science tools and infrastructure. Poor quality data can introduce noise and bias, which prevents machine learning models from making accurate predictions, but high-quality training data allows models to produce more reliable results across innumerable use cases, from automation to translation to data-driven decision-making:

1. Data collection:

First data must be collected. For AI systems like autonomous vehicles or smart homes, data collection might happen using sensors or IoT devices. Government agencies, research institutions and businesses often provide public datasets. Advertisers use clickstreams, form submissions and behavioral data from users.

2. Data cleaning and transformation:

Raw data often contains missing values, duplicates and other errors. Once data is collected, it must be cleaned to correct these errors. This can be as straightforward as standardizing formats, like ensuring that dates appear as MM/DD/YYYY. After cleaning, data often needs to be transformed into a format that is easier for algorithms to process. Feature engineering preprocesses raw data into a machine-readable format. It optimizes ML model performance by transforming and selecting relevant features.

3. Splitting the dataset:

To evaluate how well a model generalizes to new data, the dataset is typically divided into three sets. The first is a training set which is used to adjust a model's parameters to find the best match between its predictions and the data, a training process called "fitting." The second is a validation data set which is used to fine-tune hyperparameters and prevent overfitting. Finally a testing data set is used for final evaluation of model performance.

4. Data labelling:

Sometimes called "human annotation," data labelling is the process of adding meaningful labels to raw data so that a model can learn from it. Labels can describe any property of data. For example, a social media post saying "This product is terrible," could be labeled as a "negative sentiment" in a process known as sentiment analysis. A human annotator could label a photo of a dog as "dog." A bank transaction could be labeled as "fraudulent."

Further steps may include data structuring, augmentation, and versioning. Some workflows include a feedback loop wherein analysis reveals where more or better data is needed, or where unuseful data can be filtered out.

3.4.2 TRAINING DATA SET:

Testing data is used to assess a model's performance after the training portion, and specifically its accuracy, reliability and robustness in real-world conditions. Of course, this data must be different from the data used in the training of the model.

- To illustrate the point, say you're building a machine learning model to classify emails as either "spam" or "not spam."
- To train the model, you might collect a dataset of labeled emails where each one is marked as "spam" or "not spam." For example, emails with phrases like "You've been hacked" may be labeled as "spam," while emails with personalized content like "It was great meeting you yesterday" are labeled "not spam."
- The model analyzes the patterns in these emails and learns the differences between spam and non-spam emails. Thus far, this is all based on the training data.
- Once the model is trained, it needs to be tested to ensure it works well on real-world, unseen data.

The testing data might be from the same dataset (though not the same actual data) and will be similarly labeled. However the model does not know the labels.

If the model is fed an email from the testing data, "Your account details are exposed," it will likely mark it as spam based on its training. This can then be compared to the actual label to assess performance.

- You will need unknown information to test your machine learning model after it was created (using your training data). This data is known as testing data, and it may be used to assess the progress and efficiency of your algorithms' training as well as to modify or optimize them for better results.
 - Showing the original set of data.
 - Be large enough to produce reliable projections

3.4.3 TRAINING VS TESTING:

The Following is the difference between training and testing dataset:

Features	Training Data	Testing Data
Purpose	The machine-learning model is trained using training data. The more training data a model has, the more accurate predictions it can make.	Testing data is used to evaluate the model's performance.
Exposure	By using the training data, the model can gain knowledge and become more accurate in its predictions.	Until evaluation, the testing data is not exposed to the model. This guarantees that the model cannot learn the testing data by heart and produce flawless forecasts.
Distribution	This training data distribution should be similar to the distribution of actual data that the model will use.	The distribution of the testing data and the data from the real world differs greatly.
Use	To stop overfitting, training data is utilized.	By making predictions on the testing data and comparing them to the actual labels, the performance of the model is assessed.

Why do we need Training data and Testing data:

Training data teaches a machine learning model how to behave, whereas testing data assesses how well the model has learned.

- **Training Data:** The machine learning model is taught how to generate predictions or perform a specific task using training data. Since it is usually identified, every data point's output from the model is known. In order to provide predictions, the model must first learn to recognize patterns in the data. Training data can be compared to a student's textbook when learning a new subject. The learner learns by reading the text and completing the tasks, and the book offers all the knowledge they require.
- **Testing Data:** The performance of the machine learning model is measured using testing data. Usually, it is labeled and distinct from the training set. This indicates that for every data point, the model's result is unknown. On the testing data, the model's accuracy in predicting outcomes is assessed. Testing data is comparable to the exam a student takes to determine how well-versed in a subject they are. The test asks questions that the student must respond to, and the test results are used to gauge the student's comprehension.

Why is it important to use separate training and testing data?

To avoid overfitting, it is essential to use separate training and testing data. When a machine learning model learns the training data too well, it becomes hard to generalize to new data. This may happen if the training data is insufficient or not representative of the real-world data on which the model will be used.

We can confirm that the model is learning the fundamental patterns and relationships in the data and not simply memorizing the training data by using separate training and testing sets. This will assist the model in making more accurate predictions based on new data.

3.4.3 DATA CLEANING: MISSING DATA , OUTLIERS :

What is Data Cleaning?

- In data science and machine learning, the quality of input data is paramount. It's a well-established fact that data quality heavily influences the performance of machine learning models. This makes data cleaning, detecting, and correcting (or removing) corrupt or inaccurate records from a dataset a critical step in the data science pipeline.
- Data cleaning is not just about erasing data or filling in missing values. It's a comprehensive process involving various techniques to transform raw data into a format suitable for analysis. These techniques include handling missing values, removing duplicates, data type conversion, and more. Each technique has its specific use case and is applied based on the data's nature and the analysis's requirements.

Common Data Cleaning Techniques:

1. **Handling Missing Values:** Missing data can occur for various reasons, such as errors in data collection or transfer. There are several ways to handle missing data, depending on the nature and extent of the missing values.
 - **Imputation:** Here, you replace missing values with substituted values. The substituted value could be a central tendency measure like mean, median, or mode for numerical data or the most frequent category for categorical data. More sophisticated imputation methods include regression imputation and multiple imputation.
 - **Deletion:** You remove the instances with missing values from the dataset. While this method is straightforward, it can lead to loss of information, especially if the missing data is not random.
2. **Removing Duplicates:** Duplicate entries can occur for various reasons, such as data entry errors or data merging. These duplicates can skew the data and lead to biased results. Techniques for removing duplicates involve identifying these redundant entries based on key attributes and eliminating them from the dataset.

3. Data Type Conversion: Sometimes, the data may be in an inappropriate format for a particular analysis or model. For instance, a numerical attribute may be recorded as a string. In such cases, data type conversion, also known as datacasting, is used to change the data type of a particular attribute or set of attributes. This process involves converting the data into a suitable format that machine learning algorithms can easily process.

4. Outlier Detection: Outliers are data points that significantly deviate from other observations. They can be caused by variability in the data or errors. Outlier detection techniques are used to identify these anomalies. These techniques include statistical methods, such as the Z-score or IQR method, and machine learning methods, such as clustering or anomaly detection algorithms.