

# Practical-1

## Machine learning basics:

In this lab, we will go through the basics of machine learning. The student needs to make a soft copy note on the following topics:

### Topics:

#### 1. What is Machine learning.

Machine learning is a branch of artificial intelligence that focuses on developing algorithms and models that allow computer systems to automatically learn and improve from data, without being explicitly programmed.

#### 2. Steps in collection of data.

**1. Define the objectives:** Clearly establish the goals and objectives of the data collection effort. Determine what information you need to gather and why it is important.

**2. Determine the data requirements:** Identify the specific types of data you need to collect in order to address your objectives. This may include deciding on the variables, attributes, or metrics that are relevant to your analysis.

**3. Plan the data collection method:** Determine the most appropriate method for collecting the required data. This could involve various techniques such as surveys, interviews, observations, experiments, web scraping, or accessing existing datasets.

**4. Design data collection instruments:** If applicable, design the tools or instruments that will be used to gather the data. For surveys or interviews, this may involve creating questionnaires or interview protocols. Ensure that the instruments capture the necessary information accurately and effectively.

**5. Pilot testing:** Before conducting the full-scale data collection, it is often beneficial to conduct a pilot test. This involves running a small-scale trial of the data collection process to identify any issues, refine the instruments, and make necessary adjustments.

**6. Data collection:** Implement the planned data collection method. This may involve distributing surveys, conducting interviews, making observations, or performing experiments. Ensure that the data is collected consistently and accurately according to the established protocols.

**7. Data validation and quality control:** Review the collected data to ensure its accuracy, completeness, and consistency. Check for errors, inconsistencies, or outliers, and take necessary steps to clean or correct the data as required.

**8. Data storage and management:** Organize and store the collected data in a secure and accessible manner. Establish protocols for data backup, retention, and protection to ensure its integrity and privacy.

**9. Data documentation:** Document relevant details about the data collection process, including the methodology, instruments used, any limitations or biases, and any preprocessing steps performed. This documentation is important for ensuring transparency, reproducibility, and facilitating further analysis.

**10. Ethical considerations:** Consider and address ethical considerations related to data collection, such as obtaining informed consent from participants, protecting privacy, and ensuring compliance with applicable regulations and guidelines.

### 3. Steps in importing the data in python (Through: csv, json, and other data formats)

**1. Import the required libraries:** Begin by importing the necessary libraries or modules in Python. For CSV files, the built-in ``csv`` module is commonly used, while for JSON files, the ``json`` module is typically utilized.

#### 2. Importing CSV data:

a. Open the CSV file: Use the ``open()`` function to open the CSV file in the appropriate mode (``r`` for reading). You can specify the file path or name as a parameter.

b. Read the CSV data: Create a CSV reader object using the ``csv.reader()`` function, passing the file object as a parameter. Then, iterate over the rows to extract the data.

#### 3. Importing JSON data:

a. Open the JSON file: Use the ``open()`` function to open the JSON file in read mode.

b. Read the JSON data: Use the ``json.load()`` function to parse the JSON file and load its contents into a Python object (e.g., a dictionary or a list).

**4. Importing other data formats:** For importing data from other formats, you may need to use additional libraries or modules specific to those formats. For example, you can use the ``pandas`` library to import data from Excel spreadsheets (``xlsx`` files) or SQL databases (``sqlite``, ``mysql``, etc.).

## 4. Preprocessing

### a) Remove Outliers:

Removing outliers is a common preprocessing step to address extreme values that can significantly impact the analysis or modeling process. Here's a general approach to remove outliers:

1. Identify the outliers: Use statistical methods such as the z-score, interquartile range (IQR), or domain knowledge to detect outliers in the dataset.

2. Decide on the outlier removal strategy: Depending on the specific context and requirements, you can choose to remove outliers entirely or handle them in a different way (e.g., replacing with a more reasonable value).

3. Remove outliers: Remove the identified outliers from the dataset. This can be done by either excluding the entire data points or imputing them with more appropriate values.

**b) Normalize Datasets, Data Encoding:**

Normalization and encoding are important preprocessing techniques to standardize and transform the data into a suitable format for machine learning algorithms. Here are the steps for normalization and data encoding:

1. Normalize the data: Normalization scales the values of different features to a standard range, usually between 0 and 1 or -1 and 1. Common normalization techniques include min-max scaling and z-score normalization.
2. Encode categorical data: Categorical variables need to be converted into numerical representations for machine learning models to process them. One-hot encoding and label encoding are commonly used methods for this purpose. One-hot encoding creates binary columns for each category, while label encoding assigns a unique integer value to each category.

**c) Handling Missing Data:**

Missing data is a common issue in datasets, and it needs to be handled appropriately before analysis or modeling. Here's an approach to handling missing data:

1. Identify missing values: Determine which variables have missing data and assess the extent of missingness. Missing data can be represented as NaN (Not a Number), null, or other special values depending on the data format.
2. Decide on the missing data strategy: Depending on the amount and nature of missing data, you can choose from various strategies such as:
  - Deleting rows or columns with missing data: If the missing data is limited, removing the corresponding rows or columns may be a viable option.
  - Imputing missing values: Replace missing values with estimated or predicted values. Common imputation techniques include mean, median, mode imputation, or more advanced methods like regression imputation or multiple imputation.
3. Implement the chosen strategy: Apply the selected missing data strategy to the dataset, either by deleting the corresponding rows/columns or filling in missing values with appropriate imputation techniques.

**5. Machine Models****a) Types of machine learning models:**

1. Supervised learning: In supervised learning, the machine learning model is trained on labeled data, where the input features and their corresponding target outputs are provided. The model learns to map the input data to the output labels and can then make predictions on new, unseen data. Examples of supervised learning algorithms include linear regression, logistic regression, decision trees, random forests, support vector machines (SVM), and neural networks.
2. Unsupervised learning: Unsupervised learning involves training a model on unlabeled data, where only the input features are available. The goal of unsupervised learning is to discover patterns, relationships, or structures within the data. Common unsupervised

learning algorithms include clustering algorithms such as k-means clustering, hierarchical clustering, and dimensionality reduction techniques like principal component analysis (PCA) and t-SNE.

3. Reinforcement learning: Reinforcement learning is a type of machine learning where an agent learns to interact with an environment and takes actions to maximize a cumulative reward signal. The agent learns through trial and error, receiving feedback in the form of rewards or penalties. Reinforcement learning algorithms are commonly used in applications like robotics, game playing, and autonomous systems.

### **b) Parameters of machine learning models:**

Machine learning models often have various parameters that can be adjusted to improve their performance and generalization. Here are some common parameters in machine learning models:

1. Learning rate: The learning rate determines the step size at which the model updates its internal parameters during training. It affects the speed and convergence of the learning process.
2. Regularization: Regularization is used to prevent overfitting, where the model becomes too complex and performs well on the training data but poorly on new, unseen data. Parameters such as L1 or L2 regularization coefficients control the penalty applied to the model's parameters, encouraging simpler and more generalized models.
3. Number of hidden layers and units: In neural networks, the architecture of the model is defined by the number of hidden layers and the number of units in each layer. These parameters affect the model's capacity to learn complex patterns and its ability to generalize.
4. Activation functions: Activation functions introduce non-linearity into the model and determine the output of a neuron or a layer. Common activation functions include sigmoid, tanh, ReLU, and softmax, each suited for different tasks and model architectures.
5. Number of clusters: In clustering algorithms, the number of clusters is a parameter that determines the desired number of distinct groups or clusters in the data.

### **6. Test-train data split: using constant ration, k-fold cross validation**

Test-train data split and k-fold cross-validation are techniques used to evaluate the performance of machine learning models. Here's an overview of both methods:

#### **1. Test-train data split:**

- In this approach, the available dataset is divided into two subsets: a training set and a test set.
- The training set is used to train the machine learning model, while the test set is used to evaluate its performance.
- The typical ratio for the split is 70-30, 80-20, or 75-25, depending on the size of the dataset and the specific problem.

- The model is trained on the training set and then tested on the independent test set to assess its generalization ability.
- The advantage of this method is its simplicity and speed, as only two subsets are created.
- However, the evaluation may be sensitive to the specific data points in the test set, and the performance estimation may not be as reliable as with other methods.

## 2. K-fold cross-validation:

- K-fold cross-validation is a more robust method that involves dividing the dataset into K subsets or folds of approximately equal size.
- The model is trained and evaluated K times, with each fold serving as the test set once and the remaining folds used for training.
- The performance metrics obtained from each fold (e.g., accuracy, precision, recall) are then averaged to give an overall performance estimate.
- Common values for K are 5 or 10, but it can vary depending on the dataset size and computational constraints.
- K-fold cross-validation provides a more reliable estimate of the model's performance by reducing the dependency on a single train-test split.
- It helps assess how well the model performs on average across different subsets of the data and can be especially useful when the dataset size is limited.

## 7. Output Inference

Output inference refers to the process of interpreting the results or predictions generated by a machine learning model. It involves understanding and extracting meaningful insights from the output to make informed decisions or take appropriate actions. Here's a brief overview of output inference:

- 1. Analyzing predictions:** Examine the output predictions generated by the model. Depending on the problem type (classification, regression, etc.), the predictions may represent class labels, probabilities, numerical values, or other relevant information.
- 2. Interpreting prediction confidence:** Consider the confidence or uncertainty associated with the model's predictions. Assess the model's certainty in its predictions by examining the prediction probabilities or confidence intervals.
- 3. Understanding feature importance:** Determine the importance of different input features in influencing the model's predictions. Feature importance analysis helps identify which features have the most significant impact on the output and can provide insights into the underlying patterns or relationships in the data.
- 4. Evaluating model performance:** Assess the overall performance of the model by comparing its predictions against ground truth or known values. Compute appropriate evaluation metrics such as accuracy, precision, recall, F1-score, mean squared error, or others, depending on the problem domain.

**5. Extracting insights:** Extract actionable insights or relevant information from the model's output. Identify patterns, trends, or relationships that can guide decision-making or help gain a deeper understanding of the underlying data.

**6. Making informed decisions:** Utilize the insights obtained from the output inference to make informed decisions or take appropriate actions based on the specific problem or application. The output can be used for various purposes, such as recommendation systems, risk assessment, fraud detection, resource allocation, or process optimization.

## 8. Validation: different metrics – Confusion Matrix, Precision, Recall, F1-score

Validation metrics play a crucial role in evaluating the performance of machine learning models. Here are explanations of some commonly used metrics for classification tasks:

**1. Confusion Matrix:** The confusion matrix is a tabular representation that summarizes the predictions of a classification model against the true labels of the data. It provides a breakdown of predicted and actual class labels, divided into four categories:

- True Positive (TP): The model correctly predicts a positive class.
- True Negative (TN): The model correctly predicts a negative class.
- False Positive (FP): The model incorrectly predicts a positive class when the true class is negative (Type I error).
- False Negative (FN): The model incorrectly predicts a negative class when the true class is positive (Type II error).

**2. Precision:** Precision measures the proportion of correctly predicted positive instances (TP) out of all instances predicted as positive (TP + FP). It indicates how well the model identifies true positives without including false positives. A higher precision value suggests a lower rate of false positives.

**3. Recall (Sensitivity or True Positive Rate):** Recall measures the proportion of correctly predicted positive instances (TP) out of all actual positive instances (TP + FN). It indicates the model's ability to identify all positives without missing any (avoiding false negatives). A higher recall value suggests a lower rate of false negatives.

**4. F1-score:** The F1-score combines precision and recall into a single metric by taking their harmonic mean. It provides a balanced measure that considers both precision and recall. The F1-score is useful when there is an imbalance between the classes or when both false positives and false negatives need to be minimized.