# Project: Air Quality Assessment In TN

# Phase 2: Innovation

## Machine learning algorithm used are:

- Linear Regression  algorithm
- Decision Tree algorithm
- Naïve Bayes algorithm
- K-Means algorithm

## 1.Linear Regression:

Linear regression is a supervised machine learning algorithm used for predicting a continuous target variable based on one or more independent input features. It is one of the simplest and most widely used regression techniques in the field of statistics and machine learning.

a. **Simple Linear Regression:**

- In simple linear regression, there is only one independent variable (feature) that is used to predict a single target variable.

### b. Multiple Linear Regression:

- In multiple linear regression, there are two or more independent variables that are used to predict a single target variable.

## Objective Function:

- The goal of linear regression is to find the values of the coefficients ($\beta_0$, $\beta_1$, $\beta_2$, ..., $\beta_n$) that minimize the sum of squared differences between the predicted values (Y_pred) and the actual values (Y) in the training data.
- This minimization problem is typically solved using optimization techniques such as the least squares method.

## Applications:

- Linear regression is used in a wide range of applications, including finance (stock price prediction), economics (demand forecasting),

healthcare (predicting patient outcomes), and more.

- linear regression is a fundamental and interpretable algorithm for modeling the relationship between input features and a continuous target variable. It provides a simple and intuitive way to make predictions.

## 2.Decision Tree:

Decision Trees are a popular machine learning technique used for both classification and regression tasks. They are a type of supervised learning algorithm that is particularly useful for decision-making and predictive modeling.

## Basic Concept:

- A Decision Tree is a tree-like structure that represents a set of decisions and their possible consequences.
- Each internal node of the tree represents a decision or a test on an input feature.

- Each branch emanating from an internal node represents the outcome of the test.
- Each leaf node represents a class label (in classification) or a continuous value (in regression).

**Construction of Decision Trees:**
- The process of constructing a decision tree involves selecting the best feature to split the data at each internal node.
- The objective is to minimize impurity or error, depending on the task (classification or regression).

- Popular algorithms for constructing decision trees include ID3, C4.5, CART, and Random Forest.

**Uses:**
- Decision Trees are easy to understand and interpret, making them a valuable tool for explaining model decisions.
- They can handle both categorical and numerical data.

- They require little data preprocessing and are not sensitive to feature scaling.

**Ensemble Methods:**
- Decision Trees can be used in ensemble methods like Random Forest and Gradient Boosting, which combine the predictions of multiple decision trees to improve performance and reduce overfitting.

# 3.Naive Bayes Algorithm:

Naive Bayes is a popular machine learning algorithm used for classification and probabilistic modeling. It is based on Bayes' theorem, which is a fundamental concept in probability theory.

**Bayes' Theorem**:
Naive Bayes relies on Bayes' theorem, which is a mathematical formula for calculating conditional probabilities. It expresses the probability of an event A occurring given that event B has occurred as:

- P(A|B): The probability of event A given event B.
- P(B|A): The probability of event B given event A.
- P(A) and P(B): The probabilities of events A and B occurring independently.

**Types of Naive Bayes Algorithms**:

- **Multinomial Naive Bayes**: This variant is commonly used for text classification problems, where the features represent word frequencies. It models the probability of a document belonging to a particular class based on the frequency of words in the document.

- **Gaussian Naive Bayes**: This variant is suitable for continuous data, assuming that the features follow a Gaussian (normal) distribution. It's used when the features are continuous variables.

- **Bernoulli Naive Bayes**: This variant is useful for binary data, where each feature is either present (1) or absent (0). It's often used for text classification when you want to represent documents as binary feature vectors.

**Uses:**
- Naive Bayes is computationally efficient and can handle a large number of features.
- It often works well in practice, even with the simplifying independence assumption.
- It requires relatively small amounts of training data to make reasonable predictions.

**4.K-Means Algorithm:**
- K-means is a popular unsupervised machine learning algorithm used for clustering data into groups or clusters based on similarity. The primary goal of the K-means algorithm is to partition a dataset into K clusters, where each data point belongs to the cluster with the nearest mean (centroid).
- It is a simple yet effective clustering technique and is widely used in various fields, including image processing, customer segmentation, and data analysis

**Objectives:**

Group data points into K clusters based on their similarities.

**Steps:**
- Initialize K cluster centroids randomly.
- Assign each data point to the nearest centroid.
- Recalculate centroids as the mean of data points in each cluster.
- Repeat steps 2 and 3 until centroids converge (minimal change).

**Key Concepts:**
- **Centroid**: Representative point at the center of each cluster.
- **Assignment Step**: Data points are assigned to the nearest centroid.
- **Update Step**: Centroids are recalculated based on data point means.
- **Convergence**: Algorithm stops when centroids don't change significantly.

**Applications:**
- Customer segmentation.
- Image compression.
- Anomaly detection.

- Document clustering.
- Recommendation systems.

**Pros:**
- Simple and interpretable.
- Scalable for large datasets.
- Works well when clusters are relatively spherical.
- K-means is a versatile clustering algorithm used for various tasks.