

# AIML ASSIGNMENT

NAME-ADITYA PRATAP SINGH SISODIA

BATCH-11

SAP ID- 500121985

R2142230333

## Libraries and Modules

Pandas (pd): For data manipulation and to load the dataset from Google Drive. It makes many tasks such as reading in the CSV, handling missing values, and manipulating columns much easier.

## Scikit-learn Modules:

train\_test\_split: This function is used to split your data into training and test sets if you need them (although it isn't used here).

KFold: It provides K-fold cross-validation to help you get an idea of how a model performs on different subsets of your dataset.

cross\_val\_predict: This cross-validation estimator produces cross-validated estimates of model predictions on each fold, which facilitates calculating metrics over the whole dataset in a consistent manner.

GridSearchCV: This class is used for hyperparameter optimization. It tests several parameter combinations and finds the best one based on cross-validation scores.

RF and LR: These are the two classification algorithms selected for this model. RF is a tree-based ensemble model, while LR is a simple, linear model. It is truly interesting to see how a complex model may do compared to a simpler model.

Metrics: accuracy\_score, precision\_score, recall\_score, and many others. They show how the model has done; overall accuracy, precision, recall, F1, AUC is considered as these provide balanced estimation of prediction quality.

## Imbalanced-learn: SMOTE

SMOTE stands for Synthetic Minority Over-sampling Technique: It uses over sampling of synthetic examples in the minority class to balance out the classes. Since Imbalanced datasets-eg-fewer high foodwaste compared to the many low foodwaste countries-can make the models biased toward majority class, SMOTE helps the patterns learn both in classes.

**LabelEncoder:** It converts the categorical labels into numeric numbers, which is the requirement for the machine learning algorithms that need numeric input. It has been used here to assign country names with integers.

**StandardScaler:** This module standardizes features by subtracting the mean and scaling to unit variance. The reason for doing this is so that all features are on approximately the same scale. In many machine learning algorithms (for example, Logistic Regression), this prevents a dominant feature from dominating the entire model because of its scale.

**Loading and Preprocessing the Data:**

**Missing Value Handling:** Fill missing values in numeric columns by mean, not changing the dataset and would bring no errors during model training.

**Categorical Variable Encoding:** Label Encoding is done because machine learning algorithms work only on numerical data. Thus convert categorical text data into a number

**Target and Features Definition:**

**Threshold for Target:** When the target variable is of combined figure, for making an instance as high food waste or low food waste, a threshold is assigned to the target variable which turns this problem into a binary classification problem. This makes the problem even simpler and helps in finding the performance of a model.

**Removing Unnecessary Columns:** All the irrelevant columns like Country and Confidence in estimate are removed so that no noise comes into the model and it is ensured that only those features are considered which are relevant.

**Cross-validation and Hyperparameter Tuning:**

**K-Fold Cross Validation (KFold):** This is used in getting a more generalized measure of performance by averaging across multiple splits of the data such that there is no overfitting to a single data split. **Hyperparameter Tuning by GridSearchCV:** Optimization of the parameters for the best models, such as n\_estimators in Random Forest, C in Logistic Regression, and so on, leading to an improvement in performance of the model. **Evaluation Metrics:**

Accuracy, Precision, Recall, F1, ROC AUC: Each evaluation metric tells you something specific about how well your models are doing. For example:

Accuracy: Proportion of correct labels predicted

Precision and Recall: Balancing false positives and false negatives

F1 Score: A combination of Precision and Recall. One single performance metric

ROC AUC: Aggregate performance measure across all classification thresholds; useful if you are dealing with imbalanced datasets

Confusion Matrix: True positives, False positives, False negatives, True negatives; it diagnoses the error in predictions.