# CROP PREDICTION  MODEL REPORT

## Introduction:

The goal of this project was to develop a predictive model that recommends the best crops based on given environmental conditions. The dataset contains information about crops, including nutrient levels (N, P, K), environmental factors (temperature, humidity, pH, rainfall), and corresponding labels. The predictive model was built using various machine learning algorithms and evaluated based on its accuracy in predicting crop labels.

## Data Preprocessing

**Feature Renaming:** Renamed columns for better readability and understanding. This step involved renaming columns such as 'N', 'P', 'K', and 'ph' to 'Nitrogen', 'Phosphorus', 'Potassium', and 'pH value of the soil', respectively.

**Label Encoding:** The dataset was already Label encoded and given as a new column but for the model training it's better to use numerical encoded data so I used the encoded column and dropped the Non encoded Object data type column. Also I cache the label encoding for the upcoming purposes.

**Outlier Detection:** Employing the Interquartile Range (IQR) method, we identified outliers in the numerical columns of the dataset. Outliers have the potential to substantially impact the performance of our models; thus, their detection and appropriate handling are paramount for ensuring the robustness of our models and the reliability of our predictions.

**Data Distribution and Visualization:** Visualize the Data Distribution to observe some obvious facts. Trends , patterns like Label count, type of distributions. This helps to understand the dataset and give the idea for the next step.

**Feature Scaling:** Used StandardScaler to scale numerical features. Scaling ensures that all features have the same scale, preventing certain features from dominating others during model training.

**Correlation Analysis:**  we noticed the presence of a derived metric, 'temperature_humidity', which combines temperature and humidity values. Considering that this derived metric might introduce redundancy into the model, we evaluated its relevance alongside the individual features 'humidity' and 'temperature'. Given that 'temperature_humidity' encapsulates information already present in 'humidity' and 'temperature', we decided to drop the 'temperature_humidity' column to avoid multicollinearity and streamline the model. By eliminating redundant features, we aim to enhance the model's interpretability and reduce complexity, thus facilitating more efficient training and inference processes.

**PDA:** PDA was an attempt to utilize synthetic principal components derived from PCA for predicting the target variable. However, it was observed that the accuracy achieved with these synthetic components was lower compared to using the original features. As a result, the decision was made to drop the principal components and rely on the original features for predicting the target variable. This decision was based on the understanding that the original features contain more relevant information for the predictive task at hand, leading to better model performance.

Now our EDA and data cleaning is finished we can move to the model part

# Recommendations for Improvement:

- Advanced Feature Engineering: Explore advanced feature engineering techniques such as polynomial features, interaction terms, or domain-specific transformations to extract more meaningful information from the dataset. This can help capture complex relationships and improve the model's predictive performance.
- Algorithm Selection and Hyperparameter Tuning: Experiment with a wider range of machine learning algorithms beyond the ones used in the initial analysis. Consider ensemble methods, gradient boosting algorithms, or deep learning architectures to discover models that better capture the underlying patterns in the data. Perform thorough hyperparameter tuning to optimize the performance of selected algorithms.
- Incorporate Additional Data Sources: Seek out and incorporate additional data sources relevant to crop cultivation, such as soil composition data, weather forecasts, satellite imagery, or agricultural reports. Integrating diverse data sources can provide a more comprehensive understanding of the factors influencing crop growth and enable more accurate predictions.
- Error Analysis and Model Refinement: Conduct a detailed analysis of model errors, misclassifications, and areas of uncertainty. Identify patterns in mispredictions and explore potential reasons behind them, such as data quality issues, feature importance, or model biases. Use insights from error analysis to refine the model architecture, feature selection, or data preprocessing steps.
- Data Augmentation and Collection: Continue efforts to collect and augment the dataset with additional samples, features, or data points. Collaborate with agricultural experts, research institutions, or farming communities to gather field data, conduct experiments, or participate in data-sharing initiatives. Increasing the diversity and size of the dataset can lead to more robust and generalizable models.
- Longitudinal Studies and Seasonal Variations: Consider conducting longitudinal studies or capturing data across multiple growing seasons to account for seasonal variations, crop phenology, and environmental changes over time. Analyzing data over different time periods can reveal insights into long-term trends, climate impacts, and crop adaptation strategies.
- Model Interpretability and Transparency: Prioritize model interpretability and transparency by using techniques such as feature importance analysis, partial dependence plots, or model-agnostic methods like SHAP (SHapley Additive exPlanations). Understanding how the model makes predictions can help build trust with stakeholders and provide actionable insights for decision-making.

- Validation and External Validation: Validate the model's performance on external datasets or real-world scenarios to assess its generalization capabilities and robustness. Collaborate with domain experts or conduct field trials to validate model predictions under different environmental conditions, geographical regions, or farming practices.

By implementing these recommendations, you can enhance the quality, reliability, and applicability of the predictive model for crop recommendation and cultivation, ultimately contributing to sustainable agriculture practices and food security.
ion

# Instructions for Running the Code:

- Install Dependencies: Open your terminal or command prompt and navigate to the directory containing the 'requirements.txt' file. Then, run the following command to install all the required dependencies:

    pip install -r requirements.txt

- Install XGBoost (if using conda environment): If you are using a conda environment, you can install XGBoost separately by running the following command:

    conda install -c conda-forge xgboost

- Create Models Folder: Create a folder named 'models' in the directory where you'll run the Jupyter Notebook. You can do this by running the following command in your terminal or command prompt:

    mkdir models

- Place Training Dataset: Ensure that the training dataset file is placed in the same directory as the Jupyter Notebook file ('Crop_Prediction.ipynb').

- Run Jupyter Notebook: Open the Jupyter Notebook file ('Task.ipynb') and execute all the cells. This notebook contains the code for data preprocessing, model training, evaluation, and prediction.

- Assign Values for Prediction: In the last block of the notebook, assign the values you want to predict. These values should correspond to the environmental conditions for which you want to recommend crops.

By following these instructions, you'll be able to set up the environment, run the Jupyter Notebook, and make predictions based on the provided environmental conditions.

# References:

Scikit-learn documentation: https://scikit-learn.org/stable/documentation.html
Pandas documentation: https://pandas.pydata.org/docs/
NumPy documentation: https://numpy.org/doc/
Matplotlib documentation: https://matplotlib.org/stable/contents.html
Seaborn documentation: https://seaborn.pydata.org/documentation.html

# Acknowledgments: