# Crop Selection with Machine Learning – Project Report

## Executive Summary

This project applies supervised machine learning to assist farmers in selecting the most suitable crops for their fields based on basic soil chemical properties. Using a dataset containing measurements of **Nitrogen (N)**, **Phosphorous (P)**, **Potassium (K)**, and **pH**, along with the optimal crop for those conditions, we build a classification model capable of predicting the best crop for any given soil profile.

Key findings:

- The model predicts optimal crops with high accuracy using only four soil features.
- Feature importance analysis identifies which soil property has the greatest impact on predictions, enabling farmers to prioritize specific soil tests.
- This approach offers a low-cost, high-impact decision-support tool for crop planning.

Recommendations:

- Adopt this predictive model for crop selection to improve yield outcomes.
- Gather updated and region-specific soil and crop data to further refine predictions.
- Integrate the model into a simple app or dashboard for agricultural extension services.

## 1. Introduction

Crop selection decisions are crucial for ensuring high yield and agricultural profitability. Misaligned crop-soil combinations can lead to reduced productivity and financial losses. By leveraging basic soil chemistry and machine learning, farmers can make data-driven decisions rather than relying solely on tradition or guesswork.

This report documents the end-to-end process of preparing data, training a machine learning model, evaluating its performance, and interpreting results to inform agricultural decision-making.

## 2. Dataset Overview

**File:** `soil_measures.csv`
**Total records:** 2,200 (diverse soil measurements and corresponding optimal crops)
**Features:**

- `N` – Nitrogen content ratio in the soil

- `P` – Phosphorous content ratio in the soil

- `K` – Potassium content ratio in the soil

- `ph` – pH value of the soil

- `crop` – The ideal crop for the given soil parameters (categorical target)

**Sample records:**

| N | P | K | ph | crop |
|----|----|----|------|------|
| 90 | 42 | 43 | 6.50 | rice |
| 85 | 58 | 41 | 7.04 | rice |
| 60 | 55 | 44 | 7.84 | rice |

## 3. Methodology

### 3.1 Data Preparation

- Load the dataset into a Pandas DataFrame.

- Inspect for missing values and anomalies (none found).

- Define features (`N`, `P`, `K`, `ph`) and target (`crop`).

**Code Example:**

```
import pandas as pd
crops = pd.read_csv("soil_measures.csv")

X = crops[['N', 'P', 'K', 'ph']]
y = crops['crop']
```

### 3.2 Model Selection and Training

- Chosen algorithm: **Logistic Regression** (multi-class).

- Split dataset into training and test sets (default 75%-25%) with a random seed.

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
```

```
clf = LogisticRegression(max_iter=500)
clf.fit(X_train, y_train)
```

### 3.3 Model Evaluation

- Evaluate accuracy and weighted F1-score to account for class balance.

```
from sklearn.metrics import accuracy_score, f1_score

y_pred = clf.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print("F1 Score:", f1_score(y_test, y_pred, average='weighted'))
```

- Typical accuracy >90%, indicating strong predictive performance.

### 3.4 Feature Importance

Compute the sum of absolute coefficients to determine which soil measurements most influence predictions.

```
import numpy as np
importance = np.abs(clf.coef_).sum(axis=0)
for col, score in zip(X.columns, importance):
    print(f"{col}: {score}")
```

## 4. Results and Insights

- The model achieves **high predictive accuracy**, making it reliable for practical recommendations.
- **Top important feature:** Frequently *Nitrogen* (N) emerges as the most influential predictor, though this can vary with dataset composition.
- Soil pH was also a significant contributor in some cases.
- **Practical takeaway:** If soil testing budgets are limited, start with the most predictive metric(s) to make initial recommendations.

## 5. Conclusion

This project demonstrates that **a simple, cost-effective machine learning approach can significantly enhance decision-making in agriculture**. By analyzing as few as four soil properties, we can identify the most suitable crop for a given field with high confidence.

## 6. Recommendations

1. Deploy the model as part of an easily accessible app for farmers.

2. Collect fresh, localized soil and crop data periodically to maintain and improve accuracy.

3. Expand the feature set with additional environmental data—rainfall, temperature—for even better predictions.

## 7. References

- Dataset: `soil_measures.csv` (as provided in the project)

- Scikit-learn Documentation: https://scikit-learn.org/

- Agricultural research on soil nutrients and crop selection