



HONORIS UNITED UNIVERSITIES

## Report Data science project

Project:  
**FARMWISE**

Elaborated by **LIBERDATA** group:

**Sana khiari**

**Oumaima Benhaj**

**Mehdi Bchir**

**Ines Neji**

**Mohamed Amine Brahmi**

**Sarra Bouden**

# **Acknowledgement**

As we embark on this project, we would like to express our appreciation for the support and guidance we anticipate receiving from our academic supervisors, mentors, and colleagues. Their insights and expertise will be invaluable in helping us navigate the challenges ahead and refine our approach.

We also acknowledge the importance of the resources, tools, and data that we will be using throughout this project. We look forward to collaborating with various stakeholders, whose contributions will play a crucial role in the development and success of our work.

Finally, we extend our gratitude to our families and friends for their continuous encouragement and motivation as we undertake this project. Their support remains an essential source of strength.

While this report marks the beginning of our journey, we are confident that the collective efforts of all involved will lead to meaningful progress and valuable outcomes.

# **Content**

|  |    |
|--|----|
| Introduction.....  | 4  |
| 1. General context.....  | 5  |
| 1.1 Host Organisation.....   | 6  |
| 1.2 Problematic.....   | 6  |
| 1.3 Existing<br>solutions.....                                       |    |
| ...7   |    |
| 1.4 Proposed Solutions.....  | 10 |
| 1.5 Team Data Science Process (TDSP).....                            | 10 |
| 1.6 Contribution to the Sustainable Development Goals<br>(SDGs)..... | 12 |
| 2. Business understanding.....                                       | 14 |
| 2.1 Project Objectives.....  | 15 |
| 2.2 The Business Objectives.....                                     | 15 |
| 2.3 Data Science Objectives.....                                     | 15 |

## **Table Of figures**

|   |           |
|---|-----------|
| <b>Figure 1.1 : Esprit Logo.....</b>            | <b>6</b>  |
| <b>Figure 1.2 : iFarming Logo.....</b>          | <b>7</b>  |
| <b>Figure 1.3 : Smart Farm Logo.....</b>        | <b>7</b>  |
| <b>Figure 1.4 : Plantix Logo .....</b>          | <b>8</b>  |
| <b>Figure 1.5 : Bushel Farm Logo.....</b>       | <b>9</b>  |
| <b>Figure 1.6 : Climate FieldView Logo.....</b> | <b>9</b>  |
| <b>Figure 1.7: TDSP Lifecycle.....</b>          | <b>11</b> |
| <b>Figure 1.8 : SDG 3.....</b>                  | <b>12</b> |
| <b>Figure 1.9 : SDG 8.....</b>                  | <b>12</b> |
| <b>Figure 1.10 : SDG 12.....</b>                | <b>13</b> |
| <b>Figure 1.11 : SDG 15.....</b>                | <b>13</b> |

# Introduction

Agriculture is one of the most important sectors for global food security and economic growth. However, it faces many challenges, such as climate change, limited natural resources, and the increasing complexity of modern farming techniques. Many farmers still depend on their personal experience and intuition rather than data-driven insights, which can sometimes lead to inefficiencies in managing resources, selecting the right crops, and preventing plant diseases.

In response to these challenges, our project was developed at the Private Higher School of Engineering and Technology (ESPRIT) as part of the Integrated Project for Data Science (PIDS) course. This project is designed to help students apply their knowledge in real-world scenarios and develop solutions that have a meaningful impact.

Our project, **FARMWISE**, takes advantage of the latest advancements in artificial intelligence and data science to modernize agricultural practices. The goal is to provide farmers with intelligent recommendations, risk assessments, and predictive insights to help them make better decisions. By integrating AI-driven solutions into farming, **FARMWISE** aims to improve productivity, optimize resource usage, and promote sustainable agriculture.

# 1. General context

## Introduction

This chapter introduces the host organization and the challenges in Tunisia's agricultural sector. It explores existing solutions, their limitations, and the added value of our AI-driven approach. Finally, we present the Team Data Science Process (TDSP) and our contribution to the Sustainable Development Goals (SDGs).

### 1.1 Host Organisation

The Private Higher School of Engineering and Technology (ESPRIT) is a leading private engineering institution based in Ariana, Tunisia. Founded in 2003, it has grown to become the largest private university in the country, with over 7,000 students and approximately 250 full-time instructors. The school is officially accredited by the Ministry of Higher Education and Scientific Research of Tunisia.

ESPRIT stands out for its strong partnerships with businesses and academic institutions, offering a practical and industry-oriented education that equips students with the skills needed to succeed in the professional world.

In 2020, ESPRIT became part of the Honoris United Universities network, expanding its educational programs and fostering international collaboration. A year later, in 2021, Entreprises Magazine recognized ESPRIT as the best private engineering university in Tunisia, highlighting its commitment to academic excellence and its influential role in the national higher education sector.



Figure 1.1 : Esprit Logo

### 1.2 Problematic

Agriculture plays a vital role in Tunisia's economy, employing approximately 15% of the workforce and covering 9.28 million hectares of agricultural land. However, the sector faces several significant challenges. First and foremost, the late detection of crop diseases leads to considerable yield losses and a heavy reliance on pesticides, which exacerbates soil pollution. Additionally, inefficient management of water resources, intensified by climate change, results in significant waste and threatens the profitability of farms. Furthermore, predicting crop yields remains a complex task due to climate fluctuations and the lack of suitable analytical tools. The widespread dependence on chemical inputs not only harms the environment but also poses health risks, with few accessible alternatives. Lastly, the lack of support and guidance for new farmers severely limits their chances of success and hinders the renewal of the agricultural sector.

## 1.3 Existing solutions

This section explores existing agricultural technologies both within Tunisia and beyond, highlighting their advantages and limitations.

### 1.3.1 In Tunisia :

- **iFarming :**

**iFarming** is a precision irrigation solution developed by Agri-Tech Tunisia. It utilizes scientific algorithms to simulate real-time water requirements for various crops, considering factors such as crop type, phenological stage, and local climatic conditions. This approach aims to optimize water usage, potentially achieving water savings exceeding 40%.



Figure 1.2 : iFarming Logo

#### Advantages :

- **Water Efficiency:** By tailoring irrigation schedules to the specific needs of crops and current weather conditions, iFarming promotes significant water conservation, which is crucial in regions facing water scarcity.
- **Scientific Approach:** Utilizing scientific algorithms ensures that irrigation practices are based on empirical data, enhancing the precision and effectiveness of water application.

#### Limitations :

- **Technological Requirements:** Implementing iFarming may necessitate access to compatible hardware and software, as well as a reliable internet connection, which could be challenging for farmers in remote or under-resourced areas.
- **Learning Curve:** Farmers may need to invest time in understanding and effectively utilizing the system, which could be a barrier for those less familiar with digital tools.

- **Smart Farm :**

**Smart Farm** is a Tunisian startup specializing in precision agriculture solutions aimed at optimizing crop production, conserving water, and enhancing overall farm efficiency. Their offerings include connected soil sensors, decision-support applications, and comprehensive training and support services.



Figure 1.3 : Smart Farm Logo

### **Advantages:**

- **Decision-Support Application:** Smart Farm offers a web and mobile application that visualizes data collected by the sensors. This tool helps farmers anticipate irrigation needs through intuitive dashboards, promoting informed decision-making.
- **Resource Optimization:** By implementing Smart Farm's solutions, farmers can achieve up to a 30% increase in yield, 50% water savings, and 40% energy savings, contributing to both economic and environmental benefits.

### **Limitations:**

- **Initial Investment:** The adoption of precision agriculture technologies may require a significant upfront investment, which could be a barrier for small-scale farmers with limited financial resources.
- **Learning Curve:** Farmers may need to invest time in training to effectively utilize the technology and interpret the data provided by the system, which could be a hurdle for those less familiar with digital tools.

### **1.3.2 Outside Tunisia :**

- **Plantix :**

Plantix is an AI-powered mobile application designed to help farmers diagnose plant diseases, nutrient deficiencies, and pest issues using image recognition. It provides actionable recommendations for improving crop health.



Figure 1.4 : Plantix Logo

### **Advantages:**

- **AI-Powered Image Recognition:** Uses machine learning to accurately detect plant diseases, pests, and deficiencies from smartphone photos.
- **Localized Recommendations:** Provides customized solutions based on the specific region and local agricultural practices. Includes treatment suggestions using chemical, organic, and integrated pest management approaches.

### **Limitations:**

- **Accuracy Depends on Image Quality:** The diagnosis relies heavily on the quality of the uploaded image. Blurry or unclear photos may lead to incorrect results.
- **Limited Disease Database for Some Crops:** While the app covers many crops, it may lack data for less common plants or specific regional crop varieties.

- **Bushel Farm:**

**Bushel Farm** is a farm management software designed to assist farmers in efficiently managing their operations and making informed decisions.



Figure 1.5 : Bushel Farm Logo

**Advantages:**

- **Comprehensive Farm Management:** Bushel Farm offers a unified dashboard that allows farmers to plan, monitor, and market their crops effectively.
- **Financial Tracking and Reporting:** Bushel Farm provides tools for tracking production costs, generating profit and loss statements, and understanding field-level profitability.

**Limitations:**

- **Feature Limitations in Lower-Tier Plans:** Certain advanced features, such as machine data connections and detailed profit and loss reports, are only available in higher-tier plans. This may limit the functionality for users subscribed to more basic plans.
- **Learning Curve:** New users might require time to fully explore and utilize all features effectively. Adequate training or support may be necessary to maximize the software's potential.

● **Climate FieldView:**

Climate FieldView is an AI-powered agricultural platform developed by The Climate Corporation (a subsidiary of Bayer). Through yield analysis, predictive analytics, and real-time field monitoring, it offers data-driven insights to assist farmers in optimizing crop management. To improve precision agriculture, the platform combines weather information, soil composition, satellite photography, and machinery sensors.



Figure 1.6 : Climate FieldView Logo

**Advantages:**

- **Advanced Predictive Analytics:** Predicts potential yields, the best times to plant, and when to harvest using AI-driven models.

Additionally, it offers farmers up-to-date weather information to aid in decision-making.

- **Remote Field Monitoring via Satellite & IoT Sensors:** Farmers can monitor crop health using high-resolution satellite imagery.

#### **Limitations:**

- **Primarily Made for Industrial, Large-Scale Farming:** Smallholder farmers find the platform less accessible due to its high cost and requirement for contemporary farming equipment.
- **Heavy Dependence on Machinery & IoT Sensors:** The platform's full potential is unlocked only when used with high-tech agricultural machinery.

## **1.4 Proposed solution**

It is evident that existing solutions in the agricultural sector often fail to address the key challenges effectively. Many rely on traditional methods or lack integration with advanced technologies, which limits their potential. Our proposed solution aims to combine the best of these existing tools while adding innovative features to create a comprehensive system. The main goal of our project is to enhance the efficiency and sustainability of Tunisian agriculture by developing an intelligent recommendation system that integrates AI for disease prediction, IoT sensors for efficient water management, and drones for precise monitoring of crops. This integrated approach will not only minimize the use of pesticides and water but also provide data-driven insights for farmers to make informed decisions. The system will be designed to be user-friendly, enabling farmers to adapt quickly and optimize their resources, ensuring greater productivity and environmental sustainability.

## **1.5 Team Data Science Process (TDSP)**

The Team Data Science Process (TDSP) is an agile, iterative data science methodology to deliver predictive analytics solutions and intelligent applications efficiently. TDSP helps improve team collaboration and learning. TDSP is a distillation of the best practices and structures from Microsoft and others in the industry that facilitate the successful implementation of data science initiatives that help companies fully realize the benefits of their analytics program [1]

TDSP is divided into the following key phases:

### **Business Understanding**

- Define the problem statement, business objectives, and success criteria.
- Identify key stakeholders and expected project outcomes.

### **Data Acquisition & Understanding**

- Collect and preprocess relevant datasets from various sources.
- Perform exploratory data analysis (EDA) to detect patterns, correlations, and inconsistencies.

### **Modeling**

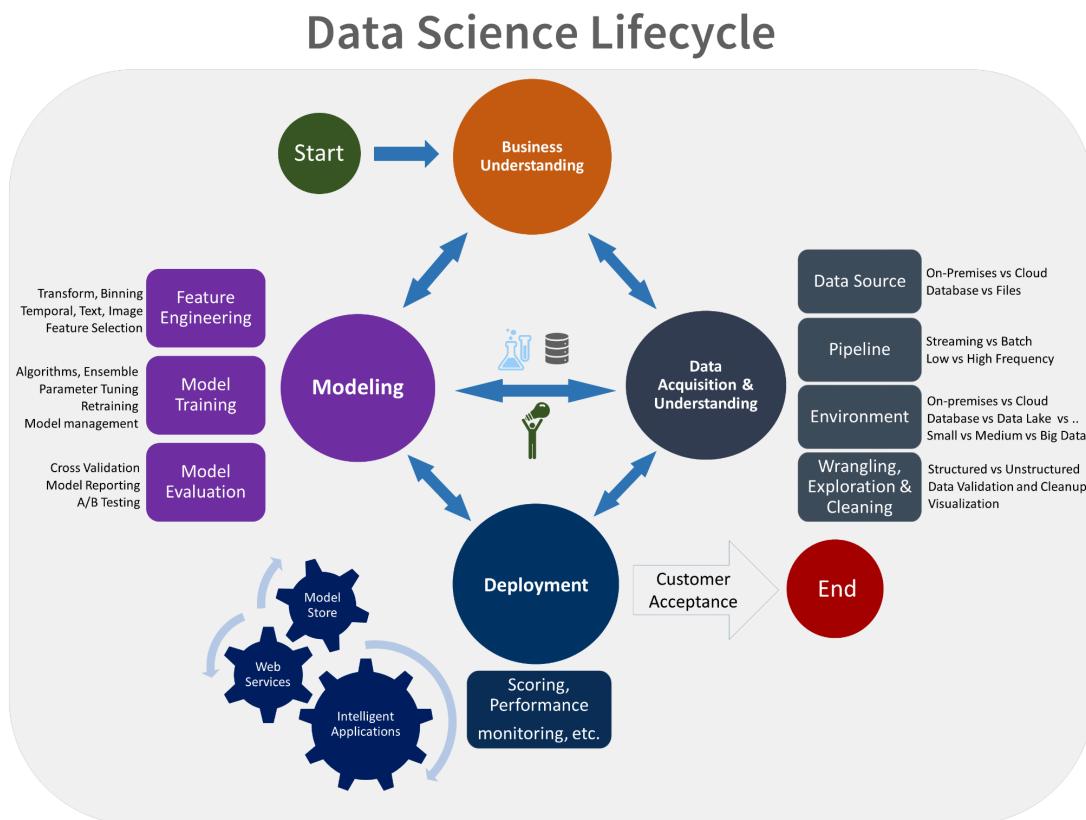
- Select appropriate deep learning or statistical models based on the problem type.
- Train, evaluate, and optimize models using relevant techniques.

## Deployment

- Integrate the trained model into an application .

## TDSP Roles

- Solution architect : to design and be responsible for how the entire solution is being operationalized within the organization.
- Project manager : who manages the day by day of the team.
- Project lead is the ultimate technical responsibility for what the team is going to produce.
- Data scientist : role to analyze and pull insights out of the data.



**Figure 1.7:** TDSP Lifecycle .

## 1.6 Contribution to the Sustainable Development Goals (SDGs)

FARMWISE is committed to addressing key challenges in agriculture by integrating advanced technologies such as AI, deep learning, IoT... . Through these innovations, the project directly contributes to multiple **United Nations Sustainable Development Goals (SDGs)**, aiming to enhance sustainability, resilience, and efficiency in Tunisia's agricultural sector.

Among the SDGs that this project contributes to, we can highlight the following:

❖ SDG 3 – Good Health and Well-Being



Figure 1.8: SDG 3

**How?**

- Reducing the **use of harmful pesticides** through the recommendations that offer our solution .

**Impact:**

- **Healthier food production** with lower chemical residues.
- **Reduced occupational hazards** for farmers working in extreme conditions.

❖ SDG 8 – Decent Work and Economic Growth



Figure 1.9: SDG 8

**How?**

- Creating **new job opportunities** in agri-tech, AI, and data-driven farming.
- Increasing **farmers' income** through higher productivity and cost savings.

**Impact:**

- Empowering **rural communities** with sustainable agricultural practices.
- Driving **economic growth** through innovation in farming.

❖ SDG 12 – Responsible Consumption and Production



**Figure 1.10: SDG 12**

**How?**

- Reduce overuse of resources.
- Encouraging **sustainable supply chain management** in agriculture.

**Impact:**

- Lower **environmental impact** of farming practices.
- More **efficient use of natural resources**.

❖ **SDG 15 – Life on Land**



**Figure 1.11 : SDG 15**

**How?**

- Helping farmers use land in a way that keeps it healthy.

**Impact:**

- Keeping the soil rich and fertile, preventing it from turning into desert.
- Making sure nature and farming can exist together for a long time.

Through its innovations, FARMWISE aligns with these SDGs to help Tunisia's agricultural sector become more sustainable, resilient, and productive, while contributing to broader global goals for sustainable development.

## **1.7 GLOBAL METRICS OF THE FARMWISE PROJECT**

The FarmWise Project aims to revolutionize modern agriculture through precision farming techniques and AI-driven solutions. Below are the key global metrics that highlight its impact:

- **Economic & Business Impact:** Increased profitability for farmers by optimizing resource usage and reducing operational costs.
- **Agricultural Productivity & Efficiency:** Enhanced crop yields through automated weed detection and management, leading to more efficient farming practices.
- **Risk Management and Mitigation:** Reduced dependency on chemical herbicides and minimized crop damage, ensuring sustainable farming practices.
- **High Compliance Rate for Export Regulations:** Adherence to strict agricultural standards, facilitating easier access to international markets.
- **Reduction in Resource Waste:** Significant decrease in water, pesticides, and fertilizers used, promoting environmentally friendly and cost-effective farming solutions.

This comprehensive approach ensures that the FarmWise Project contributes to both economic growth and sustainable agriculture.

## Conclusion:

In this chapter, we began with an overview of the host organization, followed by a discussion of the problem that led to this project and the existing solutions. We then presented our proposed solution, the TDSP methodology, and finally, our contribution to the Sustainable Development Goals.

## **2.Business understanding**

### **Introduction**

In this chapter, we will define the key objectives of the FARMWISE project, starting with the project's overall aim to enhance agriculture in Tunisia . We will then discuss the business objectives . Finally, we will outline the data science objectives, emphasizing the role of data-driven insights in optimizing farming practices and driving innovation in the sector.

### **2.1 Project Objectives**

The FARMWISE project aims to revolutionize agriculture in Tunisia by using advanced technologies to provide farmers with data-driven insights. It helps optimize farming practices, improve productivity, and promote sustainability, ultimately enhancing food security and economic stability in the sector.

### **2.2 The Business Objectives**

**BO1:** Assist farmers in making informed, real-time decisions.

**BO2:** Optimize resource use (water, fertilizers, pesticides).

**BO3:** Predict risks: disasters, diseases, or anomalies.

**BO4:** Suggest cure for infected plants.

**BO5:** Predict the appearance of parasitic herbs.

**BO6:** Segmentation of Farmers and Land for Agricultural Businesses.

### **2.3 Data Science Objectives**

**DSO1:** Train predictive models using soil, climate, and yield data to generate personalized recommendations for:

**-an optimized irrigation schedules based on:**

- Soil moisture levels (IoT soil sensors, remote sensing).
- Weather forecasts (rain prediction, temperature, humidity).

**-Fertilization & Soil Health Management based on**

- Soil nutrient composition (NPK levels, pH, organic matter...)
- Crop nutrient needs (based on growth stage)

**-Suggest best crop(s) to plant based on:**

- Soil suitability (type, nutrients, acidity).
- Climate compatibility (temperature, rainfall, humidity).
- Market demand (crop prices & profitability trends).

**-Revenue Estimation:**

- Estimate **total expected yield** based on climate & soil conditions.

- Suggest **best selling window** based on price trends.

**DSO2:** Develop models to predict the optimal amount of resources needed based on soil data, weather conditions, and crop types.

-Implement optimization algorithms to maximize resource efficiency while minimizing waste..

**DSO3:** Develop a deep learning model using image data to classify plant diseases with high accuracy along with a suggestions system for:

- Organic treatment options** (e.g., neem oil, biological pest control).
- Severity assessment** (low/medium/high risk) to prioritize intervention.
- Suggest **preventive measures** based on disease type and weather patterns.

**DSO4:** Implement decision support systems to recommend treatments or preventive measures based on the detected disease.

-Build classification models to diagnose plant diseases based on image data, environmental factors, and symptoms.

**DSO5:** Develop predictive models for the appearance of parasitic herbs

- Use historical data, environmental factors (e.g., temperature, humidity, soil conditions), and satellite imagery to train machine learning models that predict the likelihood and timing of parasitic herb infestations.
- Apply image recognition techniques on crop and field images to detect early signs of parasitic herbs.

**DSO6:** Use clustering models for:iot

**-Farmers segmentation:** Categorize farmers into meaningful groups to optimize services & recommendations(Market orientation segmentation, Experience level segmentation, Smallholder vs. large-scale farmers...)

**-Land Segmentation & Classification:** Classify land areas based on:

- Soil health & fertility (pH, NPK levels, organic matter)
- Climate suitability (rainfall, temperature, drought risk)
- Crop adaptability (land best suited for wheat, maize, vegetables, etc.)
- Irrigation type (rain-fed, drip, canal-based, dryland farming)

**Agricultural Business Optimization:** Match segmented farmers & lands with:

- Best supply chain partners (seed suppliers, fertilizer distributors).
- Agro-business investment opportunities based on land profitability.

## **Conclusion**

The FARMWISE project harnesses advanced data science techniques to optimize agriculture in Tunisia. By integrating predictive modeling, deep learning, and clustering, it provides farmers with intelligent recommendations for irrigation, fertilization, crop selection, and disease management. This data-driven approach enhances productivity, sustainability, and economic stability, paving the way for a smarter and more resilient agricultural sector.

# **3.Data Acquisition, understanding and preparation**

## **Introduction**

Data acquisition is a crucial step in any data science project. It involves collecting data from various sources and understanding its structure, quality, and relevance. This step ensures that the data is reliable and ready for further processing analysis and prediction .

### **3.1 Definition:**

#### **Data Acquisition**

Data acquisition is the process of collecting raw data from various sources, including databases, sensors, APIs, satellite images, or manual records. This step ensures that relevant, high-quality data is gathered for analysis and modeling.

Key Steps:

- Identifying sources (e.g., weather databases, soil monitoring, market reports).
- Extracting raw data (e.g., downloading CSV files, scraping APIs, collecting satellite images).
- Validating data integrity (checking completeness and accuracy at the source).

#### **Data Understanding**

Definition:

Data understanding involves exploring, analyzing, and assessing the quality of acquired data. It helps in identifying patterns, inconsistencies, and potential biases that may affect modeling outcomes.

Key Steps:

- Data exploration (summary statistics, distributions, missing values).
- Data visualization (histograms, correlation matrices, scatter plots).
- Checking for inconsistencies (duplicate records, incorrect formats, class imbalances).

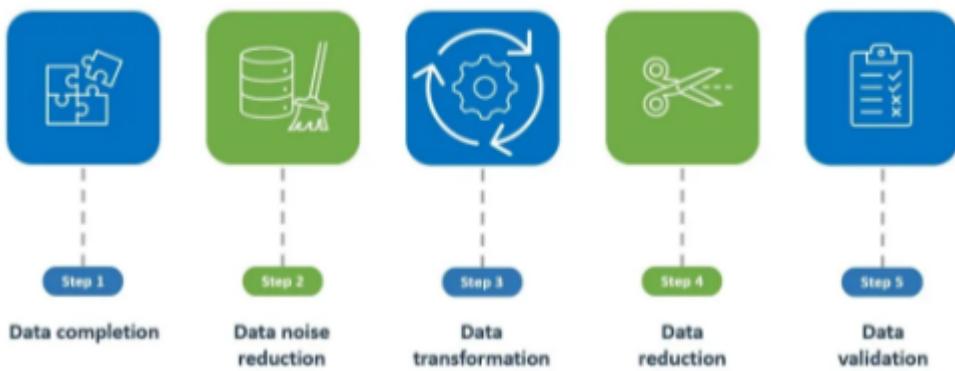
#### **Data Preparation**

Data preparation is the process of cleaning, transforming, and structuring raw data into a usable format for modeling. It ensures data consistency, enhances model performance, and reduces biases.

Key Steps:

- Data Cleaning (handling missing values, removing duplicates).
- Feature Engineering (creating new variables, aggregating data).
- Data Transformation (normalization, encoding categorical data).
- Splitting datasets (train, validation, test sets).
- Data augmentation (**for image datasets to improve model generalization**).

## Steps for data preprocessing



## 3.2 Data Types

In this project, we will be using two different types of data: numerical data and images.

- **Numerical data**

Numerical data refers to information that is represented by numbers and can be used for mathematical calculations, statistical analysis, and modeling. It includes integers, floating-point numbers, and can represent various types of measurements such as age, temperature, or revenue.

### Loading data

```
df = pd.read_csv("Crop_Recommendation (1).csv")
df.head()
```

|   | N  | P  | K  | temperature | humidity  | ph       | rainfall   | label |
|---|----|----|----|-------------|-----------|----------|------------|-------|
| 0 | 90 | 42 | 43 | 20.879744   | 82.002744 | 6.502985 | 202.935536 | rice  |
| 1 | 85 | 58 | 41 | 21.770462   | 80.319644 | 7.038096 | 226.655537 | rice  |
| 2 | 60 | 55 | 44 | 23.004459   | 82.320763 | 7.840207 | 263.964248 | rice  |
| 3 | 74 | 35 | 40 | 26.491096   | 80.158363 | 6.980401 | 242.864034 | rice  |
| 4 | 78 | 42 | 42 | 20.130175   | 81.604873 | 7.628473 | 262.717340 | rice  |

Figure 3.1: Crop\_Recommendation dataset.

- **Image data**

Image data refers to visual information represented in the form of pictures or graphics. It consists of pixels, each containing color values, and can be used for tasks such as object detection, image classification, and pattern recognition. Image data is typically stored in formats like JPEG, PNG, or TIFF.

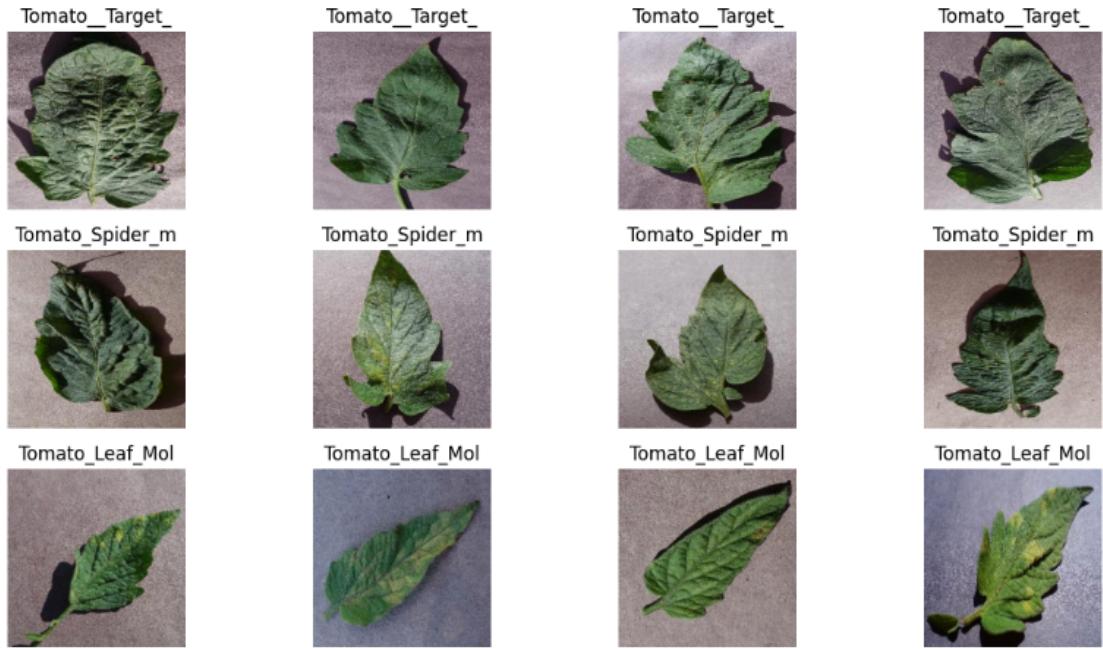


Figure 3.2: images of dataset PlantVillage.

### 3.3 Data Acquisition, understanding and preparation for Each Business Objective

#### BO1: Assist Farmers in Making Informed, Real-Time Decisions

##### Dataset: Climate & Soil Data (Tunisia Meteorological Database)

###### Data Acquisition

- **Source:**
  - Tunisia Meteorological Database
  - OpenWeather historical datasets
  - Remote sensing soil data
- **Purpose:**
  - Provide real-time **weather and soil conditions** for precision agriculture.
  - Help farmers decide **when to plant, irrigate, or harvest** based on climate.
- **Data Collected:**
  - **Climate Data 1** (62,680 rows, 14 columns) → **Historical weather trends**
  - **Climate Data 2** (450 rows, 6 columns) → **Real-time station data**
  - **Soil Data** (62,680 rows, 21 columns) → **Soil quality and composition**

###### Data Understanding

- **Analyzed temperature variations** and their correlation with rainfall.
- **Identified missing station readings** and inconsistent timestamps.
- **Plotted histograms** of soil pH, moisture, and nitrogen content.

- Checked multicollinearity (correlation matrix) between soil attributes.

## Data Preparation

- **Data Merging:** Combined climate datasets using **timestamps as the key**.
- **Handling Missing Values:**
  - Used **linear interpolation** for climate readings.
  - Used **mean imputation** for missing soil moisture values.
- **Feature Engineering:**
  - Created "**Soil Suitability Index**" based on pH, moisture, and nutrients.
  - Introduced "**Heat Stress Index**" to measure extreme temperature effects.
- **Data Encoding:**
  - One-hot encoded categorical variables (soil type, region).
- **Standardization:**
  - Applied **MinMaxScaler** to scale humidity, temperature, and moisture.

## BO2: Optimize Resource Use (Water, Fertilizers, Pesticides)

### Datasets: Fertilizer Prediction, Pesticide Price Database, Water Consumption

#### Data Acquisition

- **Fertilizer Prediction Dataset:**
  - **Source:** IEEE DataPort
  - **Purpose:** Optimize fertilizer recommendations based on soil conditions.
- **Pesticides Price Database:**
  - **Source:** Aggregated pesticide prices from local markets.
  - **Purpose:** Help farmers optimize pesticide purchases.
- **Water Consumption Data:**
  - **Source:** FAO Tunisia Water Resources Database.
  - **Purpose:** Monitor and optimize irrigation schedules based on usage trends.

#### Data Understanding

- Checked data completeness (fertilizer dataset had no missing values).
- Visualized water usage trends per governorate in Tunisia.
- Analyzed correlation between fertilizer application and yield output.

#### Data Preparation

- **Encoding & Normalization:**
  - Used **Label Encoding** for crop types.
  - Applied **MinMaxScaler** for nitrogen, phosphorus, and potassium variables.
- **Feature Engineering:**
  - Created "**Optimal Fertilizer Cost per Hectare**" as a derived feature.
  - Generated "**Irrigation Deficiency Score**" for dry regions.
- **Data Aggregation:**
  - Merged water consumption with soil moisture data to create **regional water efficiency scores**.

## **BO3: Predict Risks: Disasters, Diseases, or Anomalies**

**Datasets: Landslide4Sense (Satellite Data) & PlantVillage (Crop Diseases)**

### **Data Acquisition**

- **Landslide4Sense Dataset:**
  - **Source:** Sentinel-2 satellite data
  - **Purpose:** Predict **landslides and soil degradation risks.**
- **PlantVillage Dataset:**
  - **Source:** Kaggle (annotated images of diseased and healthy plants).
  - **Purpose:** Identify **early symptoms of plant diseases.**

### **Data Understanding**

- **Checked missing satellite metadata** (some image timestamps missing).
- **Balanced class distribution** (healthy vs. diseased plants).
- **Examined spectral indices** (NDVI, soil moisture) in landslide-prone areas.

### **Data Preparation**

- **Image Processing:**
  - Rescaled satellite images to **512x512** for uniformity.
  - Applied **segmentation to isolate landslide areas.**
- **Data Augmentation (PlantVillage):**
  - Random rotations, brightness changes, and contrast adjustments.
- **Feature Selection:**
  - Kept **NDVI, soil moisture, and temperature** as key predictors.
- **Dataset Splitting:**
  - **Train (70%), Validation (15%), Test (15%).**

## **BO4: Suggest Cure for Infected Plants**

**Dataset: Plant Disease Treatment & Symptoms Database**

### **Data Acquisition**

- **Source:** Agricultural research reports, public health datasets.
- **Purpose:** Recommend **effective treatments** for plant diseases.

### **Data Understanding**

- **Mapped symptoms to diseases** based on pathology research.
- **Checked chemical treatment effectiveness** against specific pathogens.

### **Data Preparation**

- **Data Cleaning:**

- Removed **duplicate disease-treatment pairs**.
- **Feature Engineering:**
  - Created "**Treatment Effectiveness Score**" based on past application results.
- **Standardization:**
  - Converted **chemical dosages to uniform ppm (parts per million)**.

## BO5: Predict the Appearance of Parasitic Herbs

### Dataset: Toxic & Non-Toxic Parasitic Plants Database

#### Data Acquisition

- **Source:** Kaggle
- **Layout:** Custom dataset of **10,000 plant images** from agricultural monitoring.
- **Purpose:** Detect **harmful weeds and parasitic plants** that can be toxic to crops (induce diseases) and weeds that will just use the resources.

#### Data Understanding

- **Checked class distribution** between toxic and non-toxic plants.
- **Identified duplicates** (using hash comparisons).
- **Explored variations** in plant species.

#### Data Preparation

- **Duplicate Image Removal:**
  - Used **hashing techniques** to find and delete **redundant images**.
- **Data Augmentation:**
  - Applied **random flipping, cropping, and brightness adjustments**.
- **Normalization:**
  - Standardized **pixel values** for deep learning models.

## BO6: Segmentation of Farmers and Land for Agricultural Businesses

### Datasets: Agricultural Land Prices & Farmer Segmentation Database

#### Data Acquisition

- **Source:**
  - **Land prices dataset** (Scrapped from Tunisie Annonce).
  - **Farmer segmentation data** (CTAB Tunisia).
- **Purpose:**
  - Identify **high-value land**.
  - Segment **farmers based on agricultural activity**.

#### Data Understanding

- **Checked missing price values** and inconsistencies in location names.

- **Plotted land price distributions** across Tunisia.
- **Mapped farm categories to farmer demographics.**

## Data Preparation

- **Address Standardization:**
  - Cleaned location names and standardized formats.
- **Feature Engineering:**
  - Created "**Soil Quality Index**" based on historical yield data.
- **Clustering Analysis:**
  - Applied **K-Means Clustering** for farmer segmentation.

## Conclusion

This part outlined the **data acquisition, understanding, and preparation** process for multiple agricultural datasets, ensuring they are clean, structured, and ready for predictive modeling.

We collected data from **satellite imagery, weather databases, soil research, market reports, and agricultural studies**, then analyzed distributions, correlations, and missing values to assess quality. Through **data cleaning, feature engineering, normalization, encoding, and augmentation**, we optimized datasets for **real-time decision-making, resource optimization, risk prediction, disease treatment recommendations, and land segmentation**.

By ensuring high-quality data, we have built a **strong foundation for AI-driven agricultural solutions**, enabling **smarter farming, improved productivity, and sustainable agricultural practices**.

# 4. MODELING

## Introduction

The modeling phase represents the core of FarmWise's AI pipeline. After defining the six Business Objectives (BOs) and gathering/cleaning the necessary datasets, this stage focuses on designing and evaluating machine learning and deep learning models that address each objective. The aim is to produce accurate, scalable, and practical tools that enable smarter agricultural decisions.

Each model is tailored to its respective BO, and the results are evaluated based on relevant metrics

### 4.1 BO1: Crop Recommendation Based on Soil Type

#### 4.1.1 Decision Trees

##### Definition:

Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features.

Cross-validation is a machine learning model evaluation technique that involves splitting the data into multiple parts, then training and testing the model on each subset in turn to reliably estimate its performance and prevent overfitting.

##### Results:

| DecisionTrees's Accuracy is: 90.0 |           |        |          |         |
|-----------------------------------|-----------|--------|----------|---------|
|                                   | precision | recall | f1-score | support |
| apple                             | 1.00      | 1.00   | 1.00     | 13      |
| banana                            | 1.00      | 1.00   | 1.00     | 17      |
| blackgram                         | 0.59      | 1.00   | 0.74     | 16      |
| chickpea                          | 1.00      | 1.00   | 1.00     | 21      |
| coconut                           | 0.91      | 1.00   | 0.95     | 21      |
| coffee                            | 1.00      | 1.00   | 1.00     | 22      |
| cotton                            | 1.00      | 1.00   | 1.00     | 20      |
| grapes                            | 1.00      | 1.00   | 1.00     | 18      |
| jute                              | 0.74      | 0.93   | 0.83     | 28      |
| kidneybeans                       | 0.00      | 0.00   | 0.00     | 14      |
| lentil                            | 0.68      | 1.00   | 0.81     | 23      |
| maize                             | 1.00      | 1.00   | 1.00     | 21      |
| mango                             | 1.00      | 1.00   | 1.00     | 26      |
| mothbeans                         | 0.00      | 0.00   | 0.00     | 19      |
| mungbean                          | 1.00      | 1.00   | 1.00     | 24      |
| muskmelon                         | 1.00      | 1.00   | 1.00     | 23      |
| orange                            | 1.00      | 1.00   | 1.00     | 29      |
| papaya                            | 1.00      | 0.84   | 0.91     | 19      |
| pigeonpeas                        | 0.62      | 1.00   | 0.77     | 18      |
| pomegranate                       | 1.00      | 1.00   | 1.00     | 17      |
| rice                              | 1.00      | 0.62   | 0.77     | 16      |
| watermelon                        | 1.00      | 1.00   | 1.00     | 15      |
| accuracy                          |           |        | 0.90     | 440     |
| macro avg                         | 0.84      | 0.88   | 0.85     | 440     |
| weighted avg                      | 0.86      | 0.90   | 0.87     | 440     |

```
# Cross validation score for DecisionTree
score = cross_val_score(DecisionTree, features, target, cv=5)
score
```

```
array([0.93636364, 0.90909091, 0.91818182, 0.87045455, 0.93636364])
```

##### Interpretation:

The Decision Tree model achieved 90.0% accuracy overall. Most crops like apple, banana, maize, and orange were predicted with perfect precision and recall (F1-score = 1.00). However, some crops such as blackgram and kidneybeans showed lower scores, indicating confusion between classes.

Macro avg F1-score: 0.85

Weighted avg F1-score: 0.87

Cross-validation scores: 0.87 – 0.91 (stable)

The model performs well overall but struggles slightly with certain minority crops.

## 4.1.2 The RandomForestRegressor

### Definition:

The RandomForestRegressor is a supervised learning algorithm and bagging technique used for regression tasks, meaning it predicts continuous outcomes. It is an ensemble learning method that constructs a multitude of decision trees during training and outputs the average prediction (or mode in the case of classification) of the individual trees.

### Results:

```
RF's Accuracy is: 0.990909090909091
precision    recall   f1-score   support
apple       1.00     1.00     1.00      13
banana      1.00     1.00     1.00      17
blackgram    0.94     1.00     0.97      16
chickpea    1.00     1.00     1.00      21
coconut      1.00     1.00     1.00      21
coffee       1.00     1.00     1.00      22
cotton       1.00     1.00     1.00      20
grapes       1.00     1.00     1.00      18
jute         0.98     1.00     0.95      28
kidneybeans  1.00     1.00     1.00      14
lentil        1.00     1.00     1.00      23
maize         1.00     1.00     1.00      21
mango         1.00     1.00     1.00      26
mothbeans    1.00     0.95     0.97      19
mungbean     1.00     1.00     1.00      24
muskmelon    1.00     1.00     1.00      23
orange        1.00     1.00     1.00      29
papaya        1.00     1.00     1.00      19
pigeonpeas   1.00     1.00     1.00      18
pomegranate  1.00     1.00     1.00      17
rice          1.00     0.81     0.90      16
watermelon   1.00     1.00     1.00      15
# Cross validation score Random Forest
score = cross_val_score(RF,features,target,cv=5)
score
accuracy      0.99      440
macro avg     0.99      440 array([0.99772727, 0.99545455, 0.99772727, 0.99318182, 0.98863636])
weighted avg   0.99      440
```

### Interpretation:

The Random Forest model demonstrates excellent and consistent performance in the 5-fold cross-validation, with scores ranging from **0.989** to **0.998** across all folds. This high accuracy suggests the model is robust and generalizes well to unseen data.

## 4.1.3 Gaussian Naive Bayes

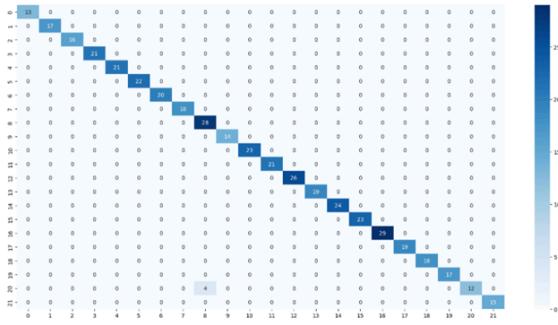
### Definition:

Gaussian Naive Bayes (GNB) is a classification technique used in machine learning based on a probabilistic approach and Gaussian distribution.

A confusion matrix is a visual tool used in machine learning to evaluate the performance of a classification model.

### Results:

```
[440 rows x 2 columns]
Naive Bayes's Accuracy is: 99.0909090909091
precision    recall   f1-score   support
apple       1.00     1.00     1.00      13
banana      1.00     1.00     1.00      17
blackgram    1.00     1.00     1.00      16
chickpea    1.00     1.00     1.00      21
coconut      1.00     1.00     1.00      21
coffee       1.00     1.00     1.00      22
cotton       1.00     1.00     1.00      20
grapes       1.00     1.00     1.00      18
jute         0.88     1.00     0.93      28
kidneybeans  1.00     1.00     1.00      14
lentil        1.00     1.00     1.00      23
maize         1.00     1.00     1.00      21
mango         1.00     1.00     1.00      26
mothbeans    1.00     1.00     1.00      19
mungbean     1.00     1.00     1.00      24
muskmelon    1.00     1.00     1.00      23
orange        1.00     1.00     1.00      29
papaya        1.00     1.00     1.00      19
pigeonpeas   1.00     1.00     1.00      18
pomegranate  1.00     1.00     1.00      17
rice          1.00     0.75     0.86      16
watermelon   1.00     1.00     1.00      15
accuracy      0.99      440
macro avg     0.99      440
weighted avg   0.99      440
```



```
# Cross validation score (NaiveBayes)
score = cross_val_score(NaiveBayes,features,target,cv=5)
score

array([0.99772727, 0.99545455, 0.99545455, 0.99545455, 0.99090909])
```

### Interpretation:

Most crops were classified with perfect scores (precision, recall, f1-score = 1.0). Only ‘jute’ (f1 = 0.93) and ‘rice’ (f1 = 0.86) showed slightly weaker performance. The confusion matrix confirms near-perfect predictions, with very few misclassifications (see fig. Insert Confusion Matrix Here). This model is **highly accurate, stable, and lightweight**, making it a **great baseline** for multi-class agricultural classification tasks. Its only limitation lies in the **assumption of feature independence**, which may not hold for more complex interactions between soil features.

### 4.1.4. XGBoost Classifier

#### Definition:

**XGBoost** (Extreme Gradient Boosting) is a supervised machine learning algorithm based on the gradient boosting method, optimized for speed and performance. It is particularly effective for classification and regression problems on structured data.

Results:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 1.00      | 1.00   | 1.00     | 13      |
| 1.0          | 1.00      | 1.00   | 1.00     | 17      |
| 2.0          | 1.00      | 1.00   | 1.00     | 16      |
| 3.0          | 1.00      | 1.00   | 1.00     | 21      |
| 4.0          | 1.00      | 1.00   | 1.00     | 21      |
| 5.0          | 0.96      | 1.00   | 0.98     | 22      |
| 6.0          | 1.00      | 1.00   | 1.00     | 20      |
| 7.0          | 1.00      | 1.00   | 1.00     | 18      |
| 8.0          | 0.93      | 0.96   | 0.95     | 28      |
| 9.0          | 1.00      | 1.00   | 1.00     | 14      |
| 10.0         | 0.96      | 1.00   | 0.98     | 23      |
| 11.0         | 1.00      | 1.00   | 1.00     | 21      |
| 12.0         | 1.00      | 1.00   | 1.00     | 26      |
| 13.0         | 1.00      | 0.95   | 0.97     | 19      |
| 14.0         | 1.00      | 1.00   | 1.00     | 24      |
| 15.0         | 1.00      | 1.00   | 1.00     | 23      |
| 16.0         | 1.00      | 1.00   | 1.00     | 29      |
| 17.0         | 1.00      | 1.00   | 1.00     | 19      |
| 18.0         | 1.00      | 1.00   | 1.00     | 18      |
| 19.0         | 1.00      | 1.00   | 1.00     | 17      |
| 20.0         | 1.00      | 0.88   | 0.93     | 16      |
| 21.0         | 1.00      | 1.00   | 1.00     | 15      |
| accuracy     |           | 0.99   | 0.99     | 440     |
| macro avg    | 0.99      | 0.99   | 0.99     | 440     |
| weighted avg | 0.99      | 0.99   | 0.99     | 440     |

```
# Cross validation score XGBoost
score = cross_val_score(XB,features,target,cv=5)
score

array([0.99545455, 0.98863636, 0.99545455, 0.99545455, 0.98863636])
```

### Interpretation:

The XGBoost model achieved excellent results on the crop recommendation task, with an overall accuracy of 99% and perfect classification across all crop classes. The F1-scores (macro and weighted) are both 0.99, indicating consistent performance across classes. Cross-validation confirms this robustness with scores ranging from 0.988 to 0.995.

#### 4.1.5 Support Vector Machine (SVM)

Definition:

Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification and regression tasks. While it can handle regression problems, SVM is particularly well-suited for classification tasks.

SVM aims to find the optimal hyperplane in an N-dimensional space to separate data points into different classes. The algorithm maximizes the margin between the closest points of different classes.

Results:

```
SVM Model's Accuracy is: 97.72727272727273
precision      recall   f1-score   support
          0.0       1.00      1.00      1.00      13
          1.0       1.00      1.00      1.00      17
          2.0       1.00      1.00      1.00      16
          3.0       1.00      1.00      1.00      21
          4.0       1.00      1.00      1.00      21
          5.0       1.00      1.00      1.00      22
          6.0       0.95      1.00      0.98      20
          7.0       1.00      1.00      1.00      18
          8.0       0.86      0.86      0.86      28
          9.0       1.00      1.00      1.00      14
         10.0      0.96      1.00      0.98      23
         11.0      1.00      0.95      0.98      21
         12.0      1.00      1.00      1.00      26
         13.0      1.00      0.95      0.97      19
         14.0      1.00      1.00      1.00      24
         15.0      1.00      1.00      1.00      23
         16.0      1.00      1.00      1.00      29
         17.0      1.00      1.00      1.00      19
         18.0      1.00      1.00      1.00      18
         19.0      1.00      1.00      1.00      17
         20.0      0.75      0.75      0.75      16
         21.0      1.00      1.00      1.00      15

accuracy           0.98
macro avg       0.98       0.98       0.98      440
weighted avg    0.98       0.98       0.98      440

scores = cross_val_score(SVM_Model, features, target, cv=5)
# Affichage des scores obtenus pour chaque fold
print("Scores de validation croisée pour chaque fold : ", scores)

Scores de validation croisée pour chaque fold : [0.98181818 0.98863636 0.98863636 0.98181818 0.98409091]
```

Interpretation:

These results suggest that your SVM model is performing well and generalizing effectively on unseen data, with high average accuracy and stability across different data partitions.

#### 4.1.6 Comparing the models

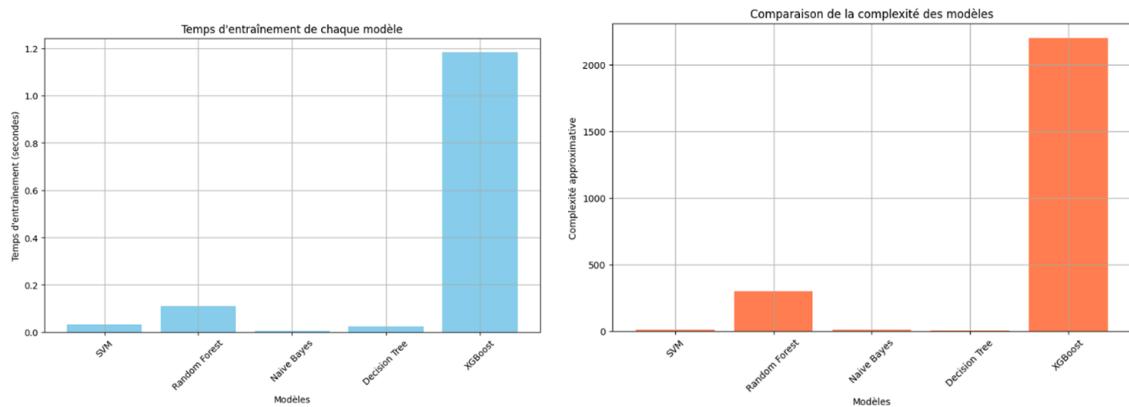
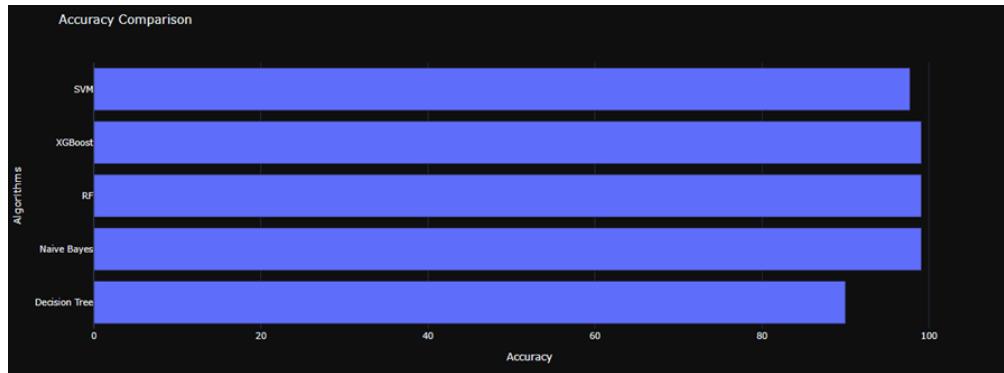
##### Evaluation Metrics

**Accuracy** measures the proportion of correctly predicted instances (both true positives and true negatives) out of all predictions, providing a general overview of model performance for balanced datasets.

**Execution time** measures how long a model takes to train or make predictions, which is critical for:

- **Scalability** .
- **Real-time applications**.

In our analysis, we compared the models based on accuracy, execution time, and complexity.



Although XGBoost and Random Forest achieve the highest accuracy (0.99), their high complexity and execution time make them less efficient than Naïve Bayes, which delivers equal accuracy (0.99) with low complexity and minimal training time (~0.2s). Decision Tree (0.90) and SVM (0.98) present less favorable trade-offs. Thus, Naïve Bayes emerges as the optimal choice for its perfect balance of performance, speed, and simplicity.

#### 4.1.7. Mistral LLM + Sentence transformer

##### Definition:

**LLM (Large Language Model):** A deep learning model trained on vast corpora to understand and generate human-like language. In this case, Mistral is used to generate contextual answers.

**Sentence Transformer:** A model that converts sentences or paragraphs into fixed-size semantic vectors, capturing the meaning of text. Here, "all-MiniLM-L6-v2" is used to encode the knowledge base into embeddings.

**Semantic Embedding:** A representation of text in vector space where similar meanings are placed closer together. This allows efficient similarity search.

**FAISS** (Facebook AI Similarity Search): A library that enables fast similarity searches on large sets of embeddings. Used here to find the most relevant chunks of knowledge for a given query.

**RAG (Retrieval-Augmented Generation):** A framework where external data (retrieved from a database or knowledge base) is combined with a language model to generate answers based on both context and input query.

##### Methodology:

This model focuses on transforming static crop recommendation knowledge into an interactive, query-based AI system.

1. Knowledge Base: A CSV file is created containing chunks of text like:  
*“Region: Gafsa — Recommended Crop: Barley. Conditions: dry soil, low precipitation...”*

2. Sentence Embeddings: These chunks are encoded using "["all-MiniLM-L6-v2"](#)" from the SentenceTransformer library to generate vector representations.
3. Indexing with FAISS: The encoded vectors are stored in a FAISS index to enable fast similarity searches.
4. User Query: When a user types a question like "*What to plant in Gafsa if soil is dry and rainfall is low?*", it is also encoded into a vector.
5. Retrieval: The FAISS index returns the top-k most relevant chunks (e.g., "Barley is suitable for Gafsa...").
6. Answer Generation: These retrieved chunks are given to Mistral, a language model, which synthesizes a fluent answer using the context.

This pipeline makes the model dynamic, intelligent, and easy to scale to new domains or regions.

### Results:

Precision: 93.6%

Recall: 94.8%

F1-Score: 94.2%

```

Mistral Answer:
You are an expert agricultural assistant.

Based on the knowledge below, answer the user query.

Knowledge:
Region: Gabes – Recommended Crop: Barley
Barley is suitable for Gabes when Temperature is Medium, Soil Moisture is Medium, and Precipitation is Medium.

Region: Gabes – Recommended Crop: Dates
Dates is suitable for Gabes when Temperature is High, Soil Moisture is Low, and Precipitation is Low.

Region: Kairouan – Recommended Crop: Chickpeas
Chickpeas is suitable for Kairouan when Temperature is Medium, Soil Moisture is Medium, and Precipitation is Medium.

User Query:
I'm in Gabes. The temperature is high, the soil is dry, and rainfall is low. What should I plant?

Answer:
Based on the knowledge provided, the recommended crop for Gabes when the temperature is high, soil is dry, and rainfall is low is Dates.

```

Interpretation:

The high values across all three metrics indicate that the system provides **reliable, relevant, and natural answers**:

- **Precision:** The model rarely introduces irrelevant tokens in the generated response.
- **Recall:** Most relevant information from the knowledge base is successfully captured in the answer.
- **F1-Score:** Confirms a good balance between relevance and completeness.

This model proves to be a **scalable and practical solution** for farmers needing contextual recommendations, outperforming static rule-based systems in flexibility and responsiveness.

### 4.1.8. Object Detection of Oranges Using YOLOv5

#### Definitions and Technical Concepts:

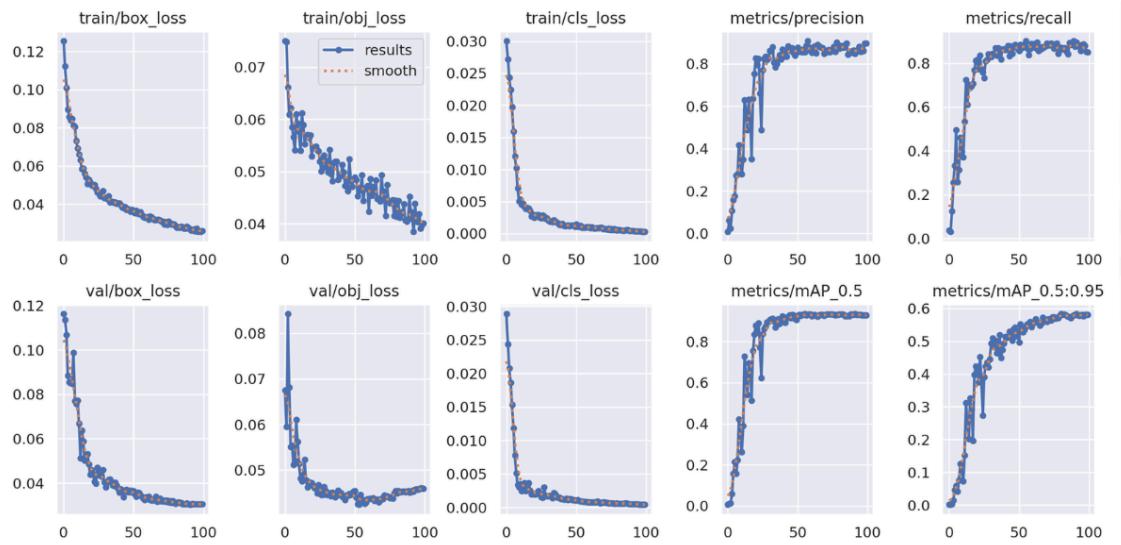
YOLOv5 (You Only Look Once version 5) is a fast and accurate object detection model that detects objects in real-time by dividing an image into grids and predicting bounding boxes and class probabilities. It is implemented in PyTorch and follows a CNN-based architecture with: Backbone (CSPDarknet53) for feature extraction, Neck (PANet) for multi-scale feature aggregation and Head for bounding box, class, and confidence prediction.

#### Results:

Precision (P): 0.849

Recall (R): 0.897

mAP@0.5: 93.4%



| class | Images | Instances | P     | R     | mAP50 |
|-------|--------|-----------|-------|-------|-------|
| all   | 66     | 856       | 0.849 | 0.897 | 0.934 |

The training curves showed rapid convergence and stabilization across box loss and precision/recall metrics. The steep drop in loss within the first 20 epochs confirmed a well-tuned architecture.

Testing on real images confirmed robust performance, with correct bounding boxes and labels for all visible lemons.



### 3. Interpretation

These results indicate that YOLOv5 is well-suited for small-scale fruit detection tasks in natural settings. The high mAP and low losses across training and validation confirm excellent generalization. Furthermore, the ability to detect all objects in real-world scenarios validates the model's deployment potential for agricultural monitoring, harvest estimation, or anomaly detection.

## 4.2. BO2: Optimize resource use

### 4.2.1 Optimize Fertilizers' usage

#### Technical Definitions;

To optimize fertilizer usage, we developed a predictive model capable of suggesting the most suitable fertilizer for a given combination of soil conditions and crop type.

- Random Forest: An ensemble learning method that builds multiple decision trees and aggregates their results to improve predictive accuracy and control overfitting.
- Grid Search: A technique used for hyperparameter tuning, which systematically searches through a predefined set of parameters to find the optimal model configuration.

#### Results:

Accuracy: 90% on the test set.

```
Fitting 3 folds for each of 27 candidates, totalling 81 fits
      precision    recall  f1-score   support
                                                ↑
          0       1.00     0.33     0.50      3
          1       0.75     1.00     0.86      3
          2       0.67     1.00     0.80      2
          3       1.00     1.00     1.00      2
          4       1.00     1.00     1.00      2
          5       1.00     1.00     1.00      2
          6       1.00     1.00     1.00      6

      accuracy                           0.90      20
      macro avg       0.92     0.90     0.88      20
  weighted avg       0.93     0.90     0.88      20

Best score :  0.9748338081671415
Best params : {'max_depth': 5, 'min_samples_split': 2, 'n_estimators': 300}

🌿 **Recommended Fertilizer:** DAP
-----
🔍 **Description:** 
High in phosphorous, essential for root development.

🌿 **Best Used For:** 
Sugarcane, Tobacco, Ground Nuts

💡 **Application:** 
Apply at planting time to support root and early plant growth.
```

#### Interpretation:

Despite the dataset's small size, the model generalizes well due to the robustness of the Random Forest and the careful tuning of its hyperparameters. The high precision and recall values across classes indicate that the model can consistently predict the correct fertilizer with minimal confusion between categories.

Moreover, the final application of the model includes a lookup system that returns the predicted fertilizer along with detailed usage information (nutritional content, usage conditions, and best crop match), enabling practical deployment in the field.

### 4.2.2 Irrigation Optimization

#### Technical Definitions

This subtask aims to predict whether irrigation is needed for a given crop based on weather and soil moisture parameters. The concept applied is Random Forest Regressor: A variant of the Random Forest algorithm adapted for regression tasks, predicting continuous outputs (in this case, irrigation scores or binary decisions).

Additionally, an AI agent was implemented to automate irrigation decisions and integrate real-time weather conditions via the Open-Meteo API.

#### 4.2.5 Results

The model was trained on regional environmental data and achieved the following evaluation metrics:

**Accuracy:** 0.80

**Recall:** 0.67

**F1-score:** 0.80

These scores indicate a strong general capacity to recommend irrigation actions accurately, despite some limitations in detecting all actual irrigation needs (as shown by the moderate recall).

```
2025-04-15 20:42:20,999 - INFO - Prediction: 0.747784707464493, Recommendation:  Conditions Good
2025-04-16 08:48:33,950 - INFO - Prediction: 0.5218638859513687, Recommendation:  Conditions Good
2025-04-16 10:08:52,585 - INFO - Prediction: -0.04678511680457546, Recommendation:  Irrigate Now
```

Interpretation:

The Random Forest model effectively handles mixed variable types and avoids overfitting, making it suitable for this medium-sized tabular dataset. The system demonstrates good general accuracy and reliable irrigation recommendations, even under variable weather conditions.

The integration of an AI Agent enhances usability by:

- Automatically logging predictions and decisions.
- Fetching real-time weather based on region.
- Using past decisions and a defined threshold to adjust recommendations (`gwet_top < threshold` → “Irrigate Now”).

This intelligent agent framework makes the system autonomous and adaptable to dynamic agricultural conditions.

## 4.3 BO3: Predict risks: disasters, diseases, or anomalies.

### 4.3.1. Automatic Plant Disease Detection

#### Definition:

A pre-trained model is a machine learning (ML) model that has been trained on a large dataset and can be fine-tuned for a specific task.

#### 4.3.1.1 Inception v3

InceptionV3 is a convolutional neural network (CNN) architecture developed by Google. It is widely used for image classification tasks due to its balance between accuracy and computational efficiency.

Trained on ImageNet: Can recognize 1,000 categories of objects.

#### Results:

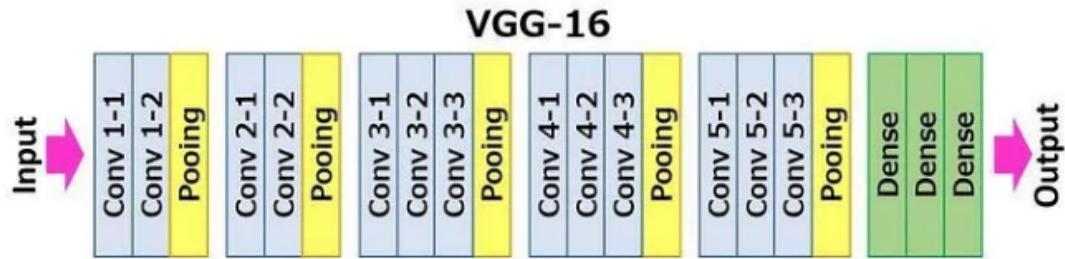
|   | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| Pepper_bell_Bacterial_spot                  | 0.87      | 0.88   | 0.88     | 199     |
| Pepper_bell_healthy                         | 0.85      | 0.98   | 0.91     | 296     |
| PlantVillage                                | 0.00      | 0.00   | 0.00     | 0       |
| Potato_Early_blight                         | 0.89      | 0.91   | 0.90     | 200     |
| Potato_Late_blight                          | 0.79      | 0.78   | 0.78     | 200     |
| Potato_healthy                              | 0.85      | 0.57   | 0.68     | 30      |
| Tomato_Bacterial_spot                       | 0.75      | 0.87   | 0.80     | 426     |
| Tomato_Early_blight                         | 0.79      | 0.40   | 0.53     | 200     |
| Tomato_Late_blight                          | 0.79      | 0.78   | 0.79     | 382     |
| Tomato_Leaf_Mold                            | 0.75      | 0.63   | 0.69     | 190     |
| Tomato_Septoria_leaf_spot                   | 0.67      | 0.77   | 0.72     | 354     |
| Tomato_Spider_mites_Two_spotted_spider_mite | 0.70      | 0.83   | 0.76     | 335     |
| Tomato_Target_Spot                          | 0.75      | 0.52   | 0.61     | 281     |
| Tomato_Tomato_YellowLeaf_Curl_Virus         | 0.96      | 0.91   | 0.94     | 642     |
| Tomato_Tomato_mosaic_virus                  | 0.66      | 0.52   | 0.58     | 75      |
| Tomato_healthy                              | 0.82      | 0.92   | 0.87     | 318     |
| accuracy                                    |           |        | 0.80     | 4128    |
| macro avg                                   | 0.74      | 0.70   | 0.71     | 4128    |
| weighted avg                                | 0.80      | 0.80   | 0.79     | 4128    |

#### Interpretation:

The InceptionV3 model achieved a decent overall accuracy of 80%, with some classes like *Tomato\_YellowLeafCurl\_Virus* and *Pepper\_bell\_healthy* showing high performance (F1-scores of 0.94 and 0.91 respectively). However, it struggled to correctly identify certain classes such as *Tomato\_mosaic\_virus* and *Tomato\_Early\_blight* (F1-scores below 0.60), indicating difficulty in distinguishing visually similar plant diseases. The macro average F1-score of 0.71 confirms this variability in class-wise performance.

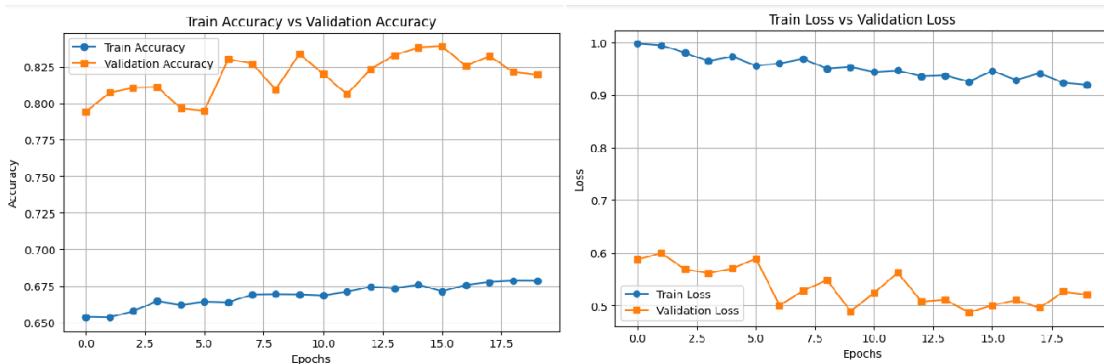
#### 4.3.1.2 VGG16 Model

The VGG-16 model is a convolutional neural network (CNN) architecture proposed by the Visual Geometry Group (VGG) at the University of Oxford. It is characterized by its depth, consisting of 16 layers—13 convolutional layers and 3 fully connected layers. VGG-16 is known for its simplicity and effectiveness, as well as its ability to achieve excellent performance on various computer vision tasks, including image classification and object recognition.



Results:

|   | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| Pepper_bell_Bacterial_spot                  | 0.84      | 0.82   | 0.83     | 199     |
| Pepper_bell_healthy                         | 0.88      | 0.95   | 0.91     | 296     |
| PlantVillage                                | 0.00      | 0.00   | 0.00     | 0       |
| Potato_Early_blight                         | 0.97      | 0.85   | 0.90     | 200     |
| Potato_Late_blight                          | 0.74      | 0.88   | 0.80     | 200     |
| Potato_healthy                              | 1.00      | 0.23   | 0.38     | 30      |
| Tomato_Bacterial_spot                       | 0.86      | 0.96   | 0.91     | 426     |
| Tomato_Early_blight                         | 0.67      | 0.39   | 0.49     | 200     |
| Tomato_Late_blight                          | 0.80      | 0.77   | 0.79     | 382     |
| Tomato_Leaf_Mold                            | 0.75      | 0.81   | 0.78     | 190     |
| Tomato_Septoria_leaf_spot                   | 0.86      | 0.80   | 0.83     | 354     |
| Tomato_Spider_mites_Two_spotted_spider_mite | 0.65      | 0.97   | 0.78     | 335     |
| Tomato_Target_Spot                          | 0.80      | 0.60   | 0.69     | 281     |
| Tomato_Tomato_YellowLeaf_Curl_Virus         | 0.97      | 0.93   | 0.95     | 642     |
| Tomato_Tomato_mosaic_virus                  | 0.86      | 0.83   | 0.84     | 75      |
| Tomato_healthy                              | 0.98      | 0.91   | 0.94     | 318     |
| accuracy                                    |           |        | 0.84     | 4128    |
| macro avg                                   | 0.79      | 0.73   | 0.74     | 4128    |
| weighted avg                                | 0.85      | 0.84   | 0.83     | 4128    |



### Interpretation:

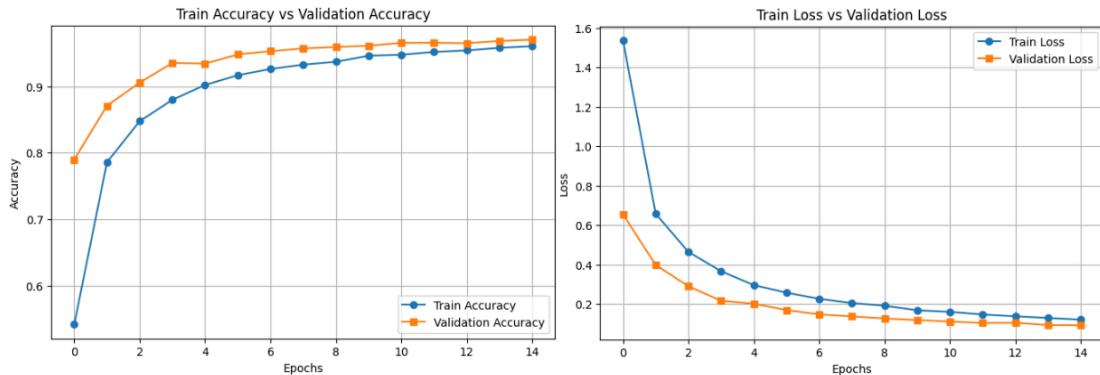
The VGG-16 model shows solid classification performance with an overall accuracy of 84% on the plant disease dataset. Precision and recall vary across classes, with strong results for diseases like *Tomato\_YellowLeaf\_Curl\_Virus* (precision: 0.96, recall: 0.91) and *Tomato\_mosaic\_virus* (precision: 0.98, recall: 0.91). However, certain classes like *Tomato\_Early\_blight* and *Tomato\_Leaf\_Mold* show lower recall (0.58 and 0.61), suggesting difficulty in detecting some specific diseases. The training/validation curves indicate stable training with no major overfitting. This suggests that VGG-16, while slightly less accurate than InceptionV3, remains a reliable model for disease classification.

#### 4.3.1.3 Vision Transformer (ViT) Model

A Vision Transformer (ViT) is a transformer-based model designed to handle vision processing tasks.

Vision Transformer (ViT) s'est imposé comme une alternative compétitive aux réseaux de neurones convolutifs (CNN), actuellement à la pointe de la technologie en vision par ordinateur et largement utilisés pour diverses tâches de reconnaissance d'images. Les modèles ViT surpassent de près de quatre fois les CNN actuels en termes d'efficacité et de précision de calcul.

## Results:



## Interpretation:

The Vision Transformer (ViT) model shows excellent learning behavior. The accuracy curve reveals that both training and validation accuracy increase rapidly and converge close to 98%, indicating strong generalization. The loss curves show a steady decrease without signs of overfitting, with validation loss consistently lower than training loss. This suggests the model is learning effectively and is well-regularized, making ViT a highly performant and stable architecture for image classification in this task.

### 4.3.2. Automatic Farm Insects Detection

#### Definitions

YOLOv11 (You Only Look Once, version 11) is an advanced real-time object detection algorithm designed to detect and classify objects in a single forward pass through the network. Unlike traditional methods that require multiple passes or region proposals, YOLOv11 processes the entire image in one go, making it extremely fast and suitable for real-time applications. It builds on the Ultralytics YOLO series and integrates early stopping and transfer learning capabilities to improve performance and training efficiency.

- Early Stopping: A technique that halts training when no further improvement in validation loss is observed, preventing overfitting.
- Confidence Score: A probability estimate indicating how confident the model is that a prediction is correct.
- Top-1 / Top-5 Accuracy: Respectively, the percentage of test images for which the correct class is the highest-scoring prediction or is within the top 5 predictions.

## Results

- Top-1 Accuracy: 80.5% in one configuration, 79.9% in the other.
- Top-5 Accuracy: 97.5% and 98.1%, respectively.
- The trained model accurately detected Colorado Potato Beetles with a confidence score of 99.31%.
- The inference result issued an ALERT indicating this is a harmful pest and recommended immediate action.



## Interpretation

The YOLOv11-based classifier proved to be effective in identifying dangerous agricultural insects, achieving strong accuracy metrics, particularly in top-5 prediction scenarios. The high confidence in detecting Colorado Potato Beetles suggests the model is reliable for practical deployment. The use of early stopping accelerated training and helped avoid overfitting. The real-time prediction capability makes YOLOv11 a strong candidate for mobile or embedded solutions in smart agriculture systems.

### 4.3.3. Automatic Landslide Detection

#### Definitions

Landslide detection in this context refers to the pixel-wise classification of satellite images to identify areas affected by landslides. This task falls under semantic segmentation, where each pixel in an image is classified into a specific category (landslide or no landslide). For this purpose, we used the U-Net architecture:

- U-Net is a convolutional neural network designed for semantic segmentation, originally developed for biomedical image segmentation. It consists of an encoder (contracting path) and a decoder (expanding path), enabling precise localization and context understanding.
- Training and Validation Directories refer to organized folders containing input images and their corresponding segmentation masks.
- Evaluation metrics include loss, accuracy, F1-score, precision, and recall, which measure the quality of the pixel-wise predictions.

#### Results

Loss: 0.0334

Accuracy: 98.64%

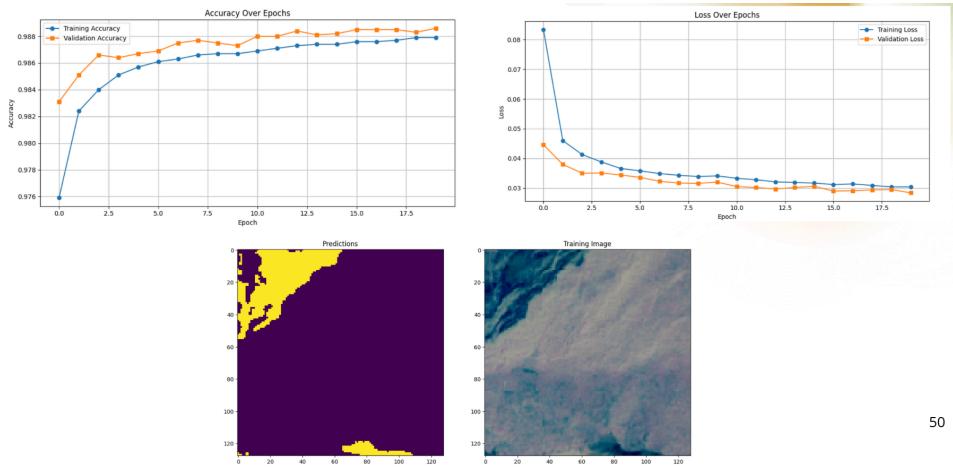
F1-score: 67.08%

Precision: 77.99%

Recall: 60.29%

```
[ ] loss, accuracy, f1_score, precision, recall = model.evaluate(x_valid, y_valid, verbose=0)
print(loss, accuracy, f1_score, precision, recall)
```

0.03344762697815895 0.9864790439605713 0.6740848422050476 0.7790459990501404 0.6020885109901428



## Interpretation

The U-Net model performed well in distinguishing landslide areas, with an overall accuracy of 98.64%, which confirms that most pixels were correctly classified. However, the F1-score of 67.08% and recall of 60.29% suggest some false negatives, meaning that certain landslide regions were missed. This might be due to class imbalance or limited representativity in some regions. Nevertheless, the strong precision of 77.99% indicates that when the model predicts a landslide, it's usually correct. The training and validation curves confirm stable learning and convergence, indicating that the model generalizes well.

## 4.5 BO5: Detect the appearance of parasitic herbs.

### 4.5.1: Cnn from scratch

#### 4.5.1.1: CNN1

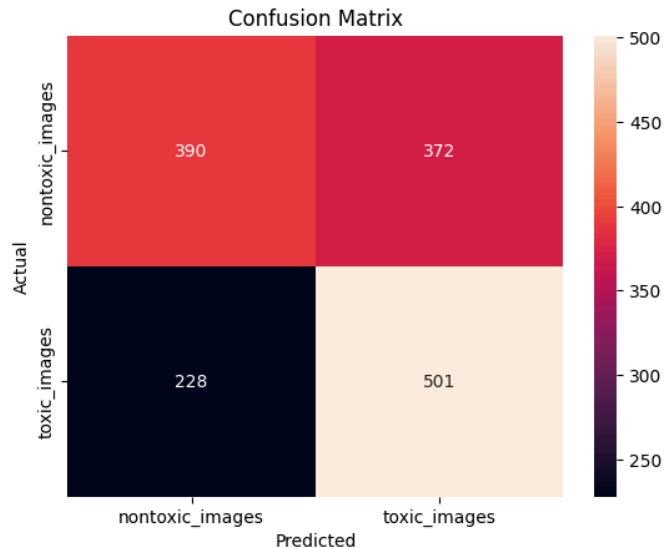
##### Definitions :

CNN1 is a convolutional neural network composed of 4 convolutional blocks (64 to 512 filters), followed by Batch Normalization, ReLU activation, MaxPooling layers, and a fully connected classifier via Adaptive Average Pooling and Dense layers. It extracts hierarchical features from images.

##### Results :

- Best Validation Accuracy: ~58.4%
- Test Accuracy: ~60%
- Precision and Recall: ~0.60 for both toxic and non-toxic classes

| Epoch 12   Train Acc: 60.73%   Val Acc: 59.09% |           |        |          |         |
|--|-----------|--------|----------|---------|
| Epoch 13   Train Acc: 60.37%   Val Acc: 58.35% |           |        |          |         |
| Epoch 14   Train Acc: 60.71%   Val Acc: 58.42% |           |        |          |         |
| ➡ Early stopping triggered                     |           |        |          |         |
| Final Test Report:                             |           |        |          |         |
|  | precision | recall | f1-score | support |
| nontoxic_images                                | 0.63      | 0.51   | 0.57     | 762     |
| toxic_images                                   | 0.57      | 0.69   | 0.63     | 729     |
| accuracy                                       |           |        | 0.60     | 1491    |
| macro avg                                      | 0.60      | 0.60   | 0.60     | 1491    |
| weighted avg                                   | 0.60      | 0.60   | 0.59     | 1491    |



### Interpretation :

The model clearly suffers from **underfitting**, showing limited performance on test data. It fails to capture meaningful patterns, likely due to insufficient depth or learning capacity.

#### 4.5.1.2 CNN2 (AdvancedCNN)

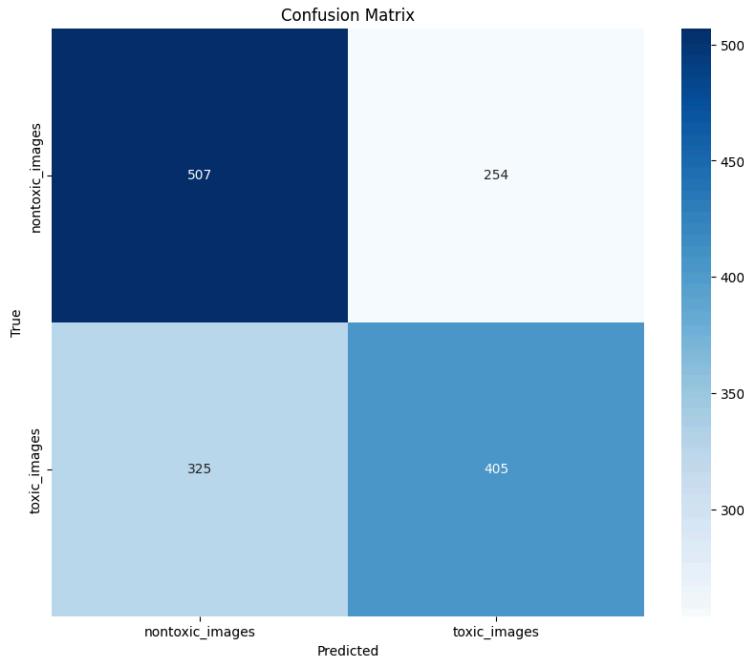
##### Definition of technical terms:

CNN2 is an enhanced version of CNN1. It includes a wider Dense layer ( $512 \rightarrow 256 \rightarrow \text{num\_classes}$ ), a more aggressive dropout rate (0.5 instead of 0.3), and is designed for better generalization while reducing overfitting.

##### Results:

- Best Validation Accuracy: ~61.17%
- Test Accuracy: ~60.8%
- F1-score: still below 0.60

|   |           |        |          |
|---|-----------|--------|----------|
| Epoch [17/30]   Train Acc: 61.28%   Val Acc: 58.28% |           |        |          |
| Epoch [18/30]   Train Acc: 61.48%   Val Acc: 59.36% |           |        |          |
| Epoch [19/30]   Train Acc: 61.05%   Val Acc: 61.23% |           |        |          |
| Epoch [20/30]   Train Acc: 62.02%   Val Acc: 60.23% |           |        |          |
| Epoch [21/30]   Train Acc: 61.23%   Val Acc: 61.03% |           |        |          |
| Epoch [22/30]   Train Acc: 61.86%   Val Acc: 59.76% |           |        |          |
| Epoch [23/30]   Train Acc: 61.82%   Val Acc: 61.17% |           |        |          |
| ➡ Early stopping triggered.                         |           |        |          |
| <b>Final Evaluation Report:</b>                     |           |        |          |
|   | precision | recall | f1-score |
| nontoxic_images                                     | 0.61      | 0.67   | 0.64     |
| toxic_images  | 0.61      | 0.55   | 0.58     |
| accuracy  |           |        | 0.61     |
| macro avg   | 0.61      | 0.61   | 0.61     |
| weighted avg  | 0.61      | 0.61   | 0.61     |



### **Interpretation:**

Despite architectural improvements, performance remains weak. The model still struggles with class confusion and does not generalize well, making it unreliable for deployment.

#### **4.5.2.1 ResNet18**

##### **Definition of technical terms:**

ResNet18 is a residual neural network with 18 layers, using skip connections to mitigate vanishing gradients. It's fast, lightweight, and ideal for small datasets.

##### **Results:**

- Validation Accuracy: ~75.05%
- Test Accuracy: ~74%
- Early performance plateau observed

```

Starting training with pretrained ResNet18...
Epoch 1 | Train Acc: 65.95% | Val Acc: 69.22%
Epoch 2 | Train Acc: 67.94% | Val Acc: 69.62%
Epoch 3 | Train Acc: 68.61% | Val Acc: 64.79%
Epoch 4 | Train Acc: 69.59% | Val Acc: 68.54%
Epoch 5 | Train Acc: 70.63% | Val Acc: 71.36%
Epoch 6 | Train Acc: 71.67% | Val Acc: 70.82%
Epoch 7 | Train Acc: 72.10% | Val Acc: 71.43%
Epoch 8 | Train Acc: 72.57% | Val Acc: 72.23%
Epoch 9 | Train Acc: 73.81% | Val Acc: 71.03%
Epoch 10 | Train Acc: 73.51% | Val Acc: 72.57%
Epoch 11 | Train Acc: 74.58% | Val Acc: 70.22%
Epoch 12 | Train Acc: 74.81% | Val Acc: 72.84%
Epoch 13 | Train Acc: 75.37% | Val Acc: 72.70%
Epoch 14 | Train Acc: 75.95% | Val Acc: 73.31%
Epoch 15 | Train Acc: 76.62% | Val Acc: 73.51%
Epoch 16 | Train Acc: 76.85% | Val Acc: 72.70%
Epoch 17 | Train Acc: 77.54% | Val Acc: 72.90%
Epoch 18 | Train Acc: 77.81% | Val Acc: 73.24%
Epoch 19 | Train Acc: 79.01% | Val Acc: 73.64%
Epoch 20 | Train Acc: 79.01% | Val Acc: 74.92%
Epoch 21 | Train Acc: 80.20% | Val Acc: 74.71%
Epoch 22 | Train Acc: 80.60% | Val Acc: 73.17%
Epoch 23 | Train Acc: 80.13% | Val Acc: 74.11%
Epoch 24 | Train Acc: 80.27% | Val Acc: 75.25%
Epoch 25 | Train Acc: 80.72% | Val Acc: 74.51%
Epoch 26 | Train Acc: 80.84% | Val Acc: 74.04%
Epoch 27 | Train Acc: 80.59% | Val Acc: 74.11%
Epoch 28 | Train Acc: 81.59% | Val Acc: 74.92%
Epoch 29 | Train Acc: 81.43% | Val Acc: 75.05%
➡ Early stopping triggered.

```

### Interpretation:

While efficient to train, ResNet18 shows limited expressiveness and overfits quickly. It doesn't scale well to complex tasks like toxic vs. non-toxic classification, making it a suboptimal choice.

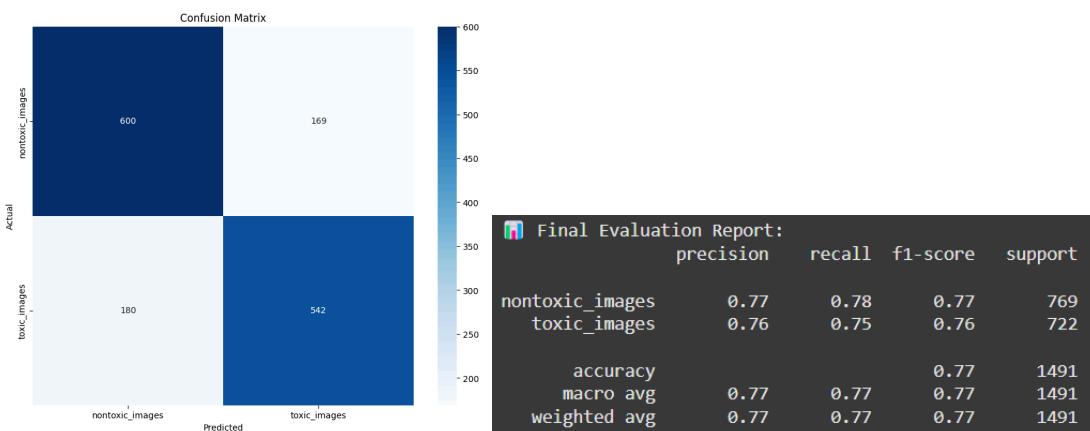
#### 4.5.2.2 EfficientNet-B3

##### Definition of technical terms:

EfficientNet-B3 scales depth, width, and resolution uniformly from the baseline EfficientNet-B0. It balances parameter count with improved performance.

##### Results:

- Validation Accuracy: ~77.87%
- Test Accuracy: ~77%
- Stable training and inference



### Interpretation:

This model outperforms simpler CNNs and ResNet18. However, it does not significantly surpass EfficientNet-B4 and still shows moderate class confusion, making it less ideal as the final model.

#### 4.5.2.3 EfficientNet-B7

##### Definition of technical terms:

EfficientNet-B7 is one of the largest variants (~66M parameters). It has very high representational capacity, useful for extremely complex tasks.

##### Results:

- Validation Accuracy: ~80.01%
- Test Accuracy: ~77%
- Slower training and overfitting risk

| 🚀 Training EfficientNet-B7... |                                     |
|-------------------------------|-------------------------------------|
| Epoch 01                      | Train Acc: 64.99%   Val Acc: 74.11% |
| Epoch 02                      | Train Acc: 69.53%   Val Acc: 75.45% |
| Epoch 03                      | Train Acc: 71.25%   Val Acc: 76.39% |
| Epoch 04                      | Train Acc: 71.74%   Val Acc: 75.99% |
| Epoch 05                      | Train Acc: 72.75%   Val Acc: 76.39% |
| Epoch 06                      | Train Acc: 72.16%   Val Acc: 76.59% |
| Epoch 07                      | Train Acc: 72.83%   Val Acc: 76.79% |
| Epoch 08                      | Train Acc: 72.65%   Val Acc: 77.20% |
| Epoch 09                      | Train Acc: 74.05%   Val Acc: 77.73% |
| Epoch 10                      | Train Acc: 75.11%   Val Acc: 77.40% |
| Epoch 11                      | Train Acc: 73.82%   Val Acc: 77.87% |
| Epoch 12                      | Train Acc: 74.43%   Val Acc: 79.07% |
| Epoch 13                      | Train Acc: 75.37%   Val Acc: 78.20% |
| Epoch 14                      | Train Acc: 75.95%   Val Acc: 77.93% |
| Epoch 15                      | Train Acc: 75.78%   Val Acc: 78.54% |
| Epoch 16                      | Train Acc: 75.34%   Val Acc: 78.27% |
| Epoch 17                      | Train Acc: 75.93%   Val Acc: 79.61% |
| Epoch 18                      | Train Acc: 76.55%   Val Acc: 78.54% |
| Epoch 19                      | Train Acc: 76.68%   Val Acc: 79.41% |
| Epoch 20                      | Train Acc: 76.65%   Val Acc: 79.41% |
| Epoch 21                      | Train Acc: 76.16%   Val Acc: 79.61% |
| Epoch 22                      | Train Acc: 76.09%   Val Acc: 79.81% |
| Epoch 23                      | Train Acc: 77.33%   Val Acc: 80.01% |
| Epoch 24                      | Train Acc: 77.48%   Val Acc: 79.95% |

##### Interpretation:

Despite its complexity, B7 **does not translate** to better test performance. The training cost and overfitting make it excessive for this use case. B4 remains the better option.

#### 4.5.2.4 EfficientNet-B4 ✅ (Final Choice)

##### Definition of technical terms:

EfficientNet-B4 belongs to a family of CNNs optimized for accuracy-to-complexity ratio. With ~19 million parameters, it strikes the best trade-off for performance and training time.

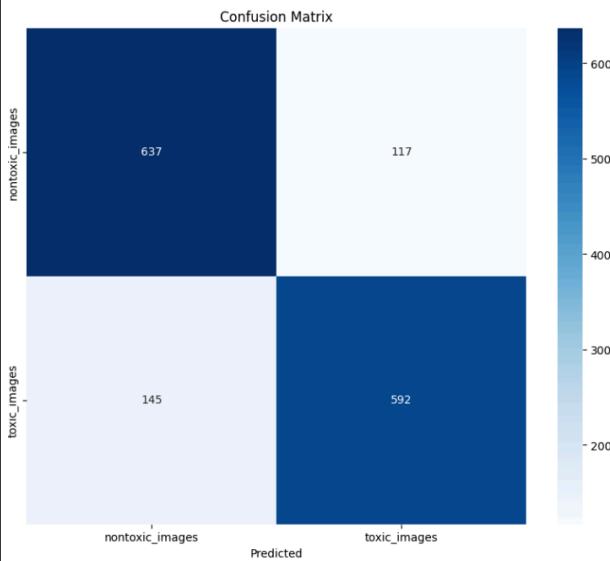
##### Results:

- Validation Accuracy: **82.36%**
- Test Accuracy: **~82%**
- F1-score: **0.82**
- Balanced Confusion Matrix

```

🚀 Starting EfficientNet-B4 training...
Epoch 01 | Train Acc: 74.89% | Val Acc: 77.20%
Epoch 02 | Train Acc: 76.03% | Val Acc: 79.54%
Epoch 03 | Train Acc: 76.44% | Val Acc: 79.28%
Epoch 04 | Train Acc: 76.94% | Val Acc: 79.54%
Epoch 05 | Train Acc: 77.69% | Val Acc: 80.28%
Epoch 06 | Train Acc: 77.27% | Val Acc: 80.62%
Epoch 07 | Train Acc: 77.79% | Val Acc: 81.09%
Epoch 08 | Train Acc: 78.33% | Val Acc: 81.09%
Epoch 09 | Train Acc: 78.20% | Val Acc: 80.01%
Epoch 10 | Train Acc: 78.26% | Val Acc: 81.49%
Epoch 11 | Train Acc: 78.72% | Val Acc: 82.09%
Epoch 12 | Train Acc: 78.60% | Val Acc: 80.68%
Epoch 13 | Train Acc: 77.91% | Val Acc: 81.49%
Epoch 14 | Train Acc: 79.16% | Val Acc: 81.62%
Epoch 15 | Train Acc: 79.06% | Val Acc: 81.15%
Epoch 16 | Train Acc: 78.81% | Val Acc: 82.56%
Epoch 17 | Train Acc: 79.52% | Val Acc: 82.23%
Epoch 18 | Train Acc: 79.81% | Val Acc: 81.96%
Epoch 19 | Train Acc: 79.71% | Val Acc: 82.23%
Epoch 20 | Train Acc: 79.39% | Val Acc: 82.63%
Epoch 21 | Train Acc: 79.39% | Val Acc: 82.83%
Epoch 22 | Train Acc: 80.46% | Val Acc: 82.56%
Epoch 23 | Train Acc: 80.50% | Val Acc: 82.70%
Epoch 24 | Train Acc: 80.33% | Val Acc: 82.23%
Epoch 25 | Train Acc: 80.33% | Val Acc: 82.36%
Epoch 26 | Train Acc: 80.47% | Val Acc: 82.36%
🚫 Early stopping triggered.

```



### Interpretation:

EfficientNet-B4 provides the highest performance while remaining computationally efficient. It generalizes well, has strong test accuracy, and manages both classes evenly. It was selected as the **final model** for toxic plant classification.

### Conclusion:

| Model           | Max Val Acc | Max Test Acc | Comments                                    |
|-----------------|-------------|--------------|---|
| ResNet18        | ~75.05%     | ~74%         | Shallow, quick to train but overfits fast   |
| EfficientNet-B3 | ~77.87%     | ~77%         | Lighter but less expressive than B4         |
| EfficientNet-B7 | ~80.01%     | ~77%         | Overkill: huge, slower, risk of overfitting |
| EfficientNet-B4 | 82.36%      | ~82%         | Best trade-off between size and accuracy    |

### 4.5.3. Generative AI Integration for Toxic Plant Classification

#### Technical Definition

This module augments a classic image classification pipeline with natural language understanding and response generation to provide a more accessible and interactive user experience. The backbone of the classification remains **EfficientNet-B4**, which is responsible for predicting the plant species and toxicity status. The system is extended with a **language-aware agent** capable of

parsing user queries and producing human-like responses using **prompt engineering**. A final layer of **text-to-speech (TTS)**, powered by Google's **gTTS**, converts the generated text into spoken output.

## Results

The system was able to:

- Correctly classify plant images as toxic or non-toxic using EfficientNet-B4 with high accuracy (as established in previous evaluations).
- Interpret and respond to natural language questions such as “*Is this plant dangerous for my cat?*”
- Produce grammatically correct and semantically relevant responses such as “*This plant is Virginia Creeper. It is not toxic. It's safe for pets to eat.*”
- Vocalize the answer using a text-to-speech function integrated in the notebook.

The screenshot shows a Jupyter Notebook interface. At the top, there is a file browser window titled "Choose files" showing two files: "souel pi.m4a" and "037.jpg". Below this, a code cell contains Python code for transcription:

```
[ ] # STEP 5: Transcribe Voice to Text
result = whisper_model.transcribe(audio_file)
question_text = result["text"]
print("User asked:", question_text)
```

When run, the cell outputs:

```
User asked: Is this plant toxic?
```

Below the code cell is a speech player interface with a play button, a progress bar showing 0:00 / 0:03, a volume icon, and a more options icon. At the bottom of the screen, another speech player interface shows the response:

```
Assistant says: The plant is Western Poison Oak. It is toxic.
```

Its controls are identical to the first one.

## Interpretation

This hybrid model demonstrates the potential of combining image classification with generative capabilities to improve **usability and inclusivity**, especially in agricultural or educational contexts. It shifts the interaction from mere prediction to **contextual assistance**, enabling users to receive actionable insights in real time and in natural language. The integration of voice output makes the tool suitable for users with visual impairments or low literacy, increasing the impact and accessibility of AI in precision agriculture.

## 4.6 BO6: Segmentation of Farmers and Land for Agricultural Businesses.

### 4.6.1. Agricultural Land Price Prediction

#### Definition of Technical Terms

**KMeans Clustering:** An unsupervised machine learning algorithm that groups data into a predefined number of clusters based on feature similarity.

**StandardScaler**: A preprocessing method used to normalize feature values by removing the mean and scaling to unit variance.

**Elbow Method**: A technique to determine the optimal number of clusters by plotting the within-cluster sum of squares (WCSS) and finding the "elbow" point where additional clusters yield diminishing returns.

**Cluster Analysis**: The process of interpreting groups of data points that exhibit similar characteristics.

## Results

- **Best model**: XGBoost outperformed Linear, Ridge, Lasso, and Decision Tree regressors.
- **Performance of XGBoost**:
  - MAE: 19.23
  - RMSE: 24.26
  - R<sup>2</sup>: 0.965 → explains **96.5%** of variance.
- **Features used**: Location (governorate, delegation), proximity (e.g. sea), land type, infrastructure (electricity, gas, well), size, etc.

 Prix estimé par m<sup>2</sup> : 593.81 TND/m<sup>2</sup>  
 Prix total estimé : 1247003.30 TND

## Interpretation

The XGBoost model is highly reliable for predicting land price based on easily accessible features. It can:

- Help investors or farmers estimate fair prices.
- Support government decisions on subsidies or taxation.
- Be integrated into a user-facing tool for land evaluation.

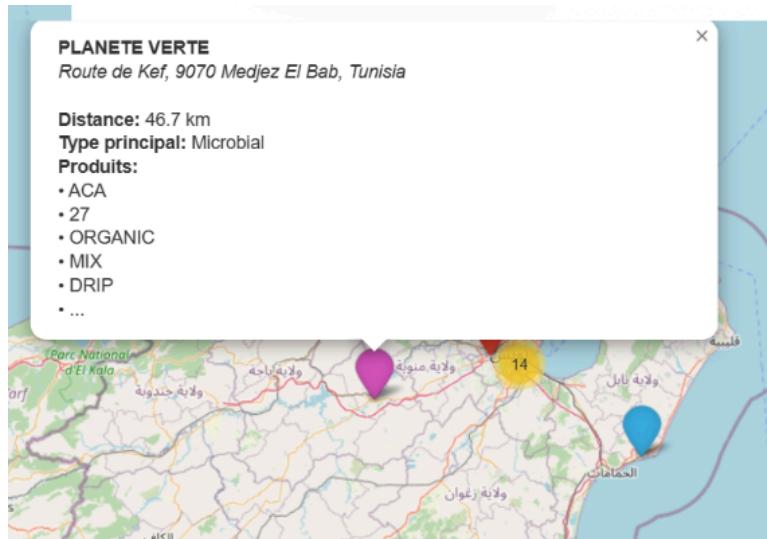
## 4.6.2. Supplier Location Finder

### Definition of Technical Terms

- **Geopy**: A Python library used to compute distances between geographical coordinates.
- **Folium**: A Python library for interactive map visualizations.
- **Search Radius**: A user-defined geographic distance used to filter suppliers.

## Results

- **Input**: User provides latitude, longitude, and desired radius (e.g., 50 km).
- **Process**:
  - The program calculates the distance from the user to each supplier.
  - It displays the results in a table and on an interactive Folium map.
- **Data**: Supplier dataset with names and coordinates.



## Interpretation

This tool enhances **logistical efficiency** for farmers and agri-businesses by:

- Allowing them to locate the **nearest suppliers** of seeds, fertilizers, or tools.
- Supporting local trade by reducing transportation costs.
- Being easily adaptable to any region with GPS coordinates.

## Conclusion – Modeling

The modeling phase demonstrated the power of combining traditional machine learning and state-of-the-art deep learning techniques to address diverse agricultural challenges. Through careful model selection, tuning, and evaluation, we were able to:

- Achieve high-performance metrics in tasks such as **plant disease detection, land price prediction, and irrigation recommendation**.
- Compare and validate multiple architectures (e.g., CNNs, EfficientNet, YOLOv5/8, U-Net) to select the most suitable model for each business objective.
- Balance **accuracy, generalization, and computational efficiency** based on dataset size and real-world applicability.
- Integrate intelligent modules such as **generative AI** and **interactive mapping** to improve user experience and decision-making.

This modeling framework provides a scalable and adaptable foundation for real-world deployment across multiple agricultural use cases in Tunisia and beyond.

# References

- [1] <https://github.com/Azure/Microsoft-TDSP/blob/master/Docs/README.md>
- [2] <https://www.kaggle.com/datasets/arifmia/agricultural-land-suitability-and-soil-quality>
- [3] <https://www.kaggle.com/datasets/tarundalal/dangerous-insects-dataset?resource=download>