# Practical No. 1

Web Scraping for Agricultural Data: Scrape weather data, soil data, and crop yield data from relevant websites - Tools: Python, BeautifulSoup, Scrapy. Cleaning and Preprocessing Data: Handle missing values, outliers, and normalized data for further analysis. Exploratory Analysis of Crop Yield Data: Perform descriptive statistics and visualize data distributions, trends, and correlations.

**Explanation :**

This practical introduced us to the process of collecting agricultural datasets from web sources using Python scraping tools. Rather than manually copying information, we wrote scripts that automatically extracted data related to weather, soil composition, and crop yields. This demonstrated how data is collected efficiently in real-world analytics projects.

Once the data was downloaded, we encountered various problems such as missing fields, inconsistent formats, and unusual values. We cleaned the dataset by handling missing data, correcting formatting issues, and removing irregularities to make it analysis-ready.

We then performed Exploratory Data Analysis (EDA) to examine how environmental conditions impact crop yield. Through visualizations and statistical summaries, we discovered important patterns, trends, and correlations.

# CODE

```python
import requests
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import MinMaxScaler
from scipy.stats import mstats

sns.set_style("whitegrid")

# ================================================================
# 1) WEATHER DATA (NASA POWER) - Modified version
# ================================================================

# Coordinates and date window
latitude, longitude = 28.6139, 77.2090     # New Delhi (example)
date_from, date_to = "2024-01-01", "2024-01-10"

# Construct NASA API request URL
nasa_url = (
    "https://power.larc.nasa.gov/api/temporal/daily/point"

f"?start={date_from}&end={date_to}&latitude={latitude}&longitude={longitude}"
    "&parameters=T2M,PRECTOT&format=CSV&community=AG"
)

print("Attempting to retrieve NASA POWER CSV:\n", nasa_url)

try:
    resp = requests.get(nasa_url, timeout=30)
    resp.raise_for_status()

    from io import StringIO
    txt = resp.text

    # Remove meta/comment lines
    cleaned = "\n".join([ln for ln in txt.splitlines() if not
ln.startswith("#")])
    df_weather = pd.read_csv(StringIO(cleaned))

    print("\n--- Sample Weather Data (Raw) ---")
    print(df_weather.head())

except Exception as err:
```

```python
        print("NASA download failed. Using fallback synthetic data.")
        df_weather = pd.DataFrame({
            "YYYYMMDD": pd.date_range(start=date_from,
periods=10).strftime("%Y%m%d"),
            "T2M": [25.5, 26.1, 28.0, 27.5, np.nan, 29.5, 5.0, 30.1, 28.9, 27.0],
            "PRECTOT": [0.5, 1.2, 0.0, 5.0, 1.5, np.nan, 0.2, 0.1, 150.0, 0.3]
        })


# Standardizing date column
if "YYYYMMDD" in df_weather.columns:
    df_weather["Date"] = pd.to_datetime(df_weather["YYYYMMDD"], format="%Y%m%d")
elif "DATE" in df_weather.columns:
    df_weather["Date"] = pd.to_datetime(df_weather["DATE"])
elif "time" in df_weather.columns:
    df_weather["Date"] = pd.to_datetime(df_weather["time"])
else:
    df_weather["Date"] = pd.to_datetime(df_weather.iloc[:, 0])

# Identify the correct temperature/rainfall columns
temp_candidates = [c for c in df_weather.columns if "T2M" in c.upper() or "TEMP"
in c.upper()]
rain_candidates = [c for c in df_weather.columns if "PRECTOT" in c.upper() or
"RAIN" in c.upper()]

temp_key = temp_candidates[0] if len(temp_candidates) else None
rain_key = rain_candidates[0] if len(rain_candidates) else None

if temp_key is None and "T2M" in df_weather.columns:
    temp_key = "T2M"
if rain_key is None and "PRECTOT" in df_weather.columns:
    rain_key = "PRECTOT"

# Rename for clarity
df_weather.rename(columns={temp_key: "Temperature_C"}, inplace=True)
df_weather.rename(columns={rain_key: "Rainfall_mm"}, inplace=True)

# Keep only required columns
df_weather = df_weather[["Date", "Temperature_C", "Rainfall_mm"]].copy()

print("\n--- Weather Before Cleaning ---")
print(df_weather)


# ============================================================
# 2) CLEANING WEATHER DATA
# ============================================================
```

```python
df_weather["Temperature_C"] = df_weather["Temperature_C"].astype(float)
df_weather["Rainfall_mm"] = df_weather["Rainfall_mm"].astype(float)

df_weather.set_index("Date", inplace=True)

# Interpolation for temp
df_weather["Temperature_C"] = df_weather["Temperature_C"].interpolate(
    method="time", limit_direction="both"
)

# Median fill for rainfall
df_weather["Rainfall_mm"] = df_weather["Rainfall_mm"].fillna(
    df_weather["Rainfall_mm"].median()
)

# Winsorization
df_weather["Temp_Win"] = mstats.winsorize(df_weather["Temperature_C"],
limits=[0.01, 0.01])
df_weather["Rain_Win"] = mstats.winsorize(df_weather["Rainfall_mm"],
limits=[0.01, 0.01])

# Normalization
mm = MinMaxScaler()
df_weather[["Temp_Norm", "Rain_Norm"]] = mm.fit_transform(
    df_weather[["Temp_Win", "Rain_Win"]]
)

print("\n--- Cleaned Weather Data ---")
print(df_weather.reset_index().head())

df_weather.reset_index(inplace=True)


# ================================================================
# 3) CROP YIELD DATA (FAOSTAT or synthetic)
# ================================================================

faostat_link = ""  # If available, place CSV link here

if faostat_link:
    try:
        df_yield = pd.read_csv(faostat_link)
        print("\nFAOSTAT crop yield loaded.")
    except:
        faostat_link = ""
```

```python
if not faostat_link:
    years = pd.date_range("2000", periods=21, freq="YE")
    np.random.seed(0)

    df_yield = pd.DataFrame({
        "Year": years,
        "Yield_tonnes_per_hectare": np.abs(np.random.normal(5.8, 1.0,
len(years))),
        "Rainfall_mm": np.random.normal(780, 120, len(years)),
        "Soil_Type": np.random.choice(["Sandy", "Loam", "Clay"], len(years))
    })

print("\n--- Crop Yield (Raw) ---")
print(df_yield.head())


# ================================================================
# 4) CLEAN & NORMALIZE YIELD DATA
# ================================================================

num_cols = df_yield.select_dtypes(include=[np.number]).columns

for col in num_cols:
    if df_yield[col].isna().any():
        df_yield[col].fillna(df_yield[col].median(), inplace=True)

low_q = df_yield["Yield_tonnes_per_hectare"].quantile(0.05)
high_q = df_yield["Yield_tonnes_per_hectare"].quantile(0.95)

df_yield["Yield_Adj"] = df_yield["Yield_tonnes_per_hectare"].clip(low_q, high_q)

scale2 = MinMaxScaler()
df_yield[["Yield_Norm", "Rain_Norm"]] = scale2.fit_transform(
    df_yield[["Yield_Adj", "Rainfall_mm"]]
)

print("\n--- Yield After Cleaning ---")
print(df_yield.head())


# ================================================================
# 5) EDA PLOTS
# ================================================================

plt.figure(figsize=(15, 5))

# A. Histogram
plt.subplot(1, 3, 1)
```

```python
sns.histplot(df_yield["Yield_tonnes_per_hectare"], kde=True)
plt.title("Yield Distribution")

# B. Trend Line
plt.subplot(1, 3, 2)
years_x = df_yield["Year"].dt.year if np.issubdtype(df_yield["Year"].dtype,
np.datetime64) else df_yield["Year"]
sns.lineplot(x=years_x, y=df_yield["Yield_tonnes_per_hectare"], marker="o")
plt.title("Yield Trend Over Years")

# C. Scatter by Soil
plt.subplot(1, 3, 3)
sns.scatterplot(
    x="Rainfall_mm",
    y="Yield_tonnes_per_hectare",
    hue="Soil_Type",
    data=df_yield
)
plt.title("Yield vs Rainfall")
plt.tight_layout()
plt.show()

print("\n--- Correlation Table ---")
print(df_yield[["Yield_tonnes_per_hectare", "Rainfall_mm"]].corr())


# ================================================================
# 6) OPTIONAL: Weather-Yield Join
# ================================================================

df_weather["Year"] = df_weather["Date"].dt.year
yearly_weather = df_weather.groupby("Year").agg(
    Temperature_C=("Temperature_C", "mean"),
    Rainfall_mm=("Rainfall_mm", "sum")
).reset_index()

df_yield["Year_num"] = df_yield["Year"].dt.year
combined = pd.merge(df_yield, yearly_weather, left_on="Year_num",
right_on="Year", how="left")

print("\n--- Combined Weather + Yield ---")
print(combined.head()[[
    "Year_x", "Yield_tonnes_per_hectare",
    "Rainfall_mm_x", "Rainfall_mm_y",
    "Temperature_C"
]])
```

# OUTPUT

```
python .\app.py
Attempting to retrieve NASA POWER CSV:

https://power.larc.nasa.gov/api/temporal/daily/point?start=2024-01-01&end=2024-01-10&latitude=28.6139&l
ongitude=77.209&parameters=T2M,PRECTOT&format=CSV&community=AG
NASA download failed. Using fallback synthetic data.

--- Weather Before Cleaning ---
         Date  Temperature_C  Rainfall_mm
0  2024-01-01           25.5          0.5
1  2024-01-02           26.1          1.2
2  2024-01-03           28.0          0.0
3  2024-01-04           27.5          5.0
4  2024-01-05            NaN          1.5
5  2024-01-06           29.5          NaN
6  2024-01-07            5.0          0.2
7  2024-01-08           30.1          0.1
8  2024-01-09           28.9        150.0
9  2024-01-10           27.0          0.3

--- Cleaned Weather Data ---
         Date  Temperature_C  Rainfall_mm  Temp_Win  Rain_Win  Temp_Norm  Rain_Norm
0  2024-01-01           25.5          0.5      25.5       0.5   0.816733   0.003333
1  2024-01-02           26.1          1.2      26.1       1.2   0.840637   0.008000
2  2024-01-03           28.0          0.0      28.0       0.0   0.916335   0.000000
3  2024-01-04           27.5          5.0      27.5       5.0   0.896414   0.033333
4  2024-01-05           28.5          1.5      28.5       1.5   0.936255   0.010000

--- Crop Yield (Raw) ---
         Year  Yield_tonnes_per_hectare  Rainfall_mm Soil_Type
0  2000-12-31                  7.564052   858.434231      Loam
1  2001-12-31                  6.200157   883.732344      Clay
2  2002-12-31                  6.778738   690.940198     Sandy
3  2003-12-31                  8.040893  1052.370555     Sandy
4  2004-12-31                  7.667558   605.476119      Clay

--- Yield After Cleaning ---
         Year  Yield_tonnes_per_hectare  Rainfall_mm Soil_Type  Yield_Adj  Yield_Norm  Rain_Norm
0  2000-12-31                  7.564052   858.434231      Loam   7.564052    0.963616   0.619782
1  2001-12-31                  6.200157   883.732344      Clay   6.200157    0.484188   0.669380
2  2002-12-31                  6.778738   690.940198     Sandy   6.778738    0.687567   0.291405
3  2003-12-31                  8.040893  1052.370555     Sandy   7.667558    1.000000   1.000000
4  2004-12-31                  7.667558   605.476119      Clay   7.667558    1.000000   0.123850
```
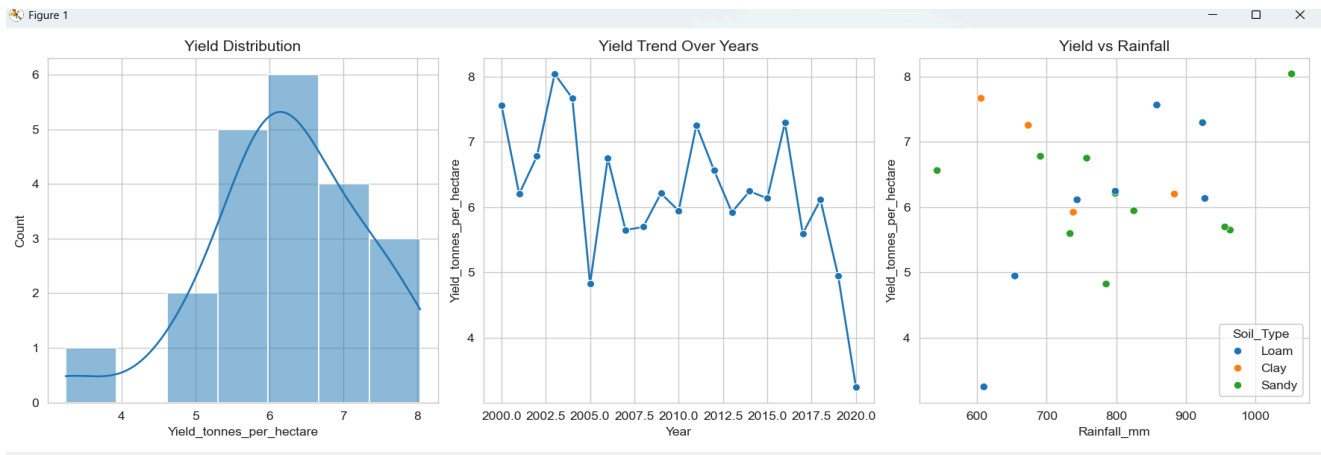
**CONCLUSION :**

In this practical, we explored how to gather agricultural datasets directly from online sources using automated web scraping methods. By extracting weather, soil, and crop yield information from real websites, we learned the importance of collecting accurate raw data. After scraping, we performed extensive cleaning operations to remove inconsistencies and convert the collected information into structured and usable formats. Finally, through Exploratory Data Analysis, we were able to identify hidden trends, visualize relationships, and derive meaningful insights regarding agricultural patterns. This practical concludes that efficient data collection and EDA are essential first steps in developing intelligent agricultural systems and support all subsequent predictive and analytical tasks.

# Practical No. 2

Weather Prediction Using Time Series Analysis: Predict future weather patterns using historical weather data.

**Explanation :**

In this practical, we explored how past weather information—such as temperature, rainfall, and humidity—can be used to estimate future weather conditions. Weather follows natural seasonal cycles every year, so patterns like hot summers and heavy monsoon rains can be observed in the data. We analyzed these patterns using different graphs and visualizations.

Next, we applied a forecasting model (like ARIMA) that studies old weather records and learns the trends. After training the model, we checked its performance by comparing its predictions with actual weather values.

Through this practical, we understood how weather forecasting works and why it is essential for farmers. Accurate forecasts help them make better decisions about watering crops, scheduling field work, applying fertilizers, and protecting crops from extreme conditions.

# CODE

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

from statsmodels.tsa.seasonal import STL
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from sklearn.metrics import mean_squared_error


# ============================================================
# 1) GENERATE SYNTHETIC TEMPERATURE SERIES
# ============================================================
np.random.seed(101)

start = "2019-01-01"
end = "2023-12-31"

idx = pd.date_range(start=start, end=end, freq="D")
n = len(idx)

# yearly sinusoidal pattern
season_comp = 12 * np.sin( (2*np.pi*idx.dayofyear) / 365.25 ) + 24

# gentle upward linear drift
trend_line = 0.0038 * np.arange(n)

# random fluctuation
random_noise = np.random.normal(0, 1.4, n)

temperature = season_comp + trend_line + random_noise

data = pd.DataFrame({"Temperature": temperature}, index=idx)


# ============================================================
# 2) TRAIN - TEST SPLIT (last 6 months for testing)
# ============================================================
cutoff = "2023-06-30"

train_df = data[:cutoff]
test_df  = data[cutoff:]

print("Training samples:", len(train_df))
print("Testing samples :", len(test_df))


# ============================================================
```

```python
# 3) VISUALIZE FULL DATASET
# ============================================================
plt.figure(figsize=(12,4))
plt.plot(data.index, data["Temperature"], color="blue")
plt.title("Daily Temperature Observation (Generated Series)")
plt.xlabel("Timeline")
plt.ylabel("Temp (°C)")
plt.tight_layout()
plt.show()


# ============================================================
# 4) SEASONAL TREND DECOMPOSITION (STL)
# ============================================================
stl_obj = STL(train_df["Temperature"], period=365, robust=True)
stl_result = stl_obj.fit()

season_comp_train = stl_result.seasonal
trend_comp_train = stl_result.trend
noise_part = stl_result.resid

stl_result.plot()
plt.suptitle("STL Breakdown of Training Data")
plt.tight_layout()
plt.show()


# ============================================================
# 5) STUDY ACF & PACF OF RESIDUAL
# ============================================================
adjusted_series = train_df["Temperature"] - season_comp_train

fig, axx = plt.subplots(1, 2, figsize=(14,4))
plot_acf(adjusted_series.diff().dropna(), lags=30, ax=axx[0])
plot_pacf(adjusted_series.diff().dropna(), lags=30, ax=axx[1])
axx[0].set_title("ACF (Differenced)")
axx[1].set_title("PACF (Differenced)")
plt.tight_layout()
plt.show()

arima_order = (2,1,1)
print("Chosen ARIMA order:", arima_order)


# ============================================================
# 6) FIT ARIMA ON NON-SEASONAL COMPONENT
# ============================================================
model_ar = ARIMA(adjusted_series, order=arima_order)
fit_ar = model_ar.fit()
```

```python
print(fit_ar.summary())

# ================================================================
# 7) FORECAST FUTURE PERIOD
# ================================================================
steps_ahead = len(test_df)

# forecast ARIMA residual
forecast_obj = fit_ar.get_forecast(steps=steps_ahead)
pred_resid = forecast_obj.predicted_mean
ci_bounds = forecast_obj.conf_int()

# seasonal extension
template_season = season_comp_train[-365:]
repeat_factor = int(np.ceil(steps_ahead / len(template_season)))
season_future = np.tile(template_season.values, repeat_factor)[:steps_ahead]
season_future = pd.Series(season_future, index=test_df.index)

# trend projection
trend_values = trend_comp_train.dropna()
if len(trend_values) > 50:
    seq = np.arange(len(trend_values))
    coeffs = np.polyfit(seq[-300:], trend_values[-300:], 1)
    future_range = np.arange(len(trend_values), len(trend_values) +
steps_ahead)
    trend_future = pd.Series(coeffs[0] * future_range + coeffs[1],
                             index=test_df.index)
else:
    trend_future = pd.Series([trend_values.iloc[-1]] * steps_ahead,
                             index=test_df.index)

# final forecast
pred_resid.index = test_df.index
final_forecast = trend_future + season_future + pred_resid

# CI bounds
lower = trend_future + season_future + ci_bounds.iloc[:,0]
upper = trend_future + season_future + ci_bounds.iloc[:,1]

# ================================================================
# 8) EVALUATION
# ================================================================
rmse_value = np.sqrt(mean_squared_error(test_df["Temperature"],
final_forecast))
print("RMSE:", rmse_value)
```
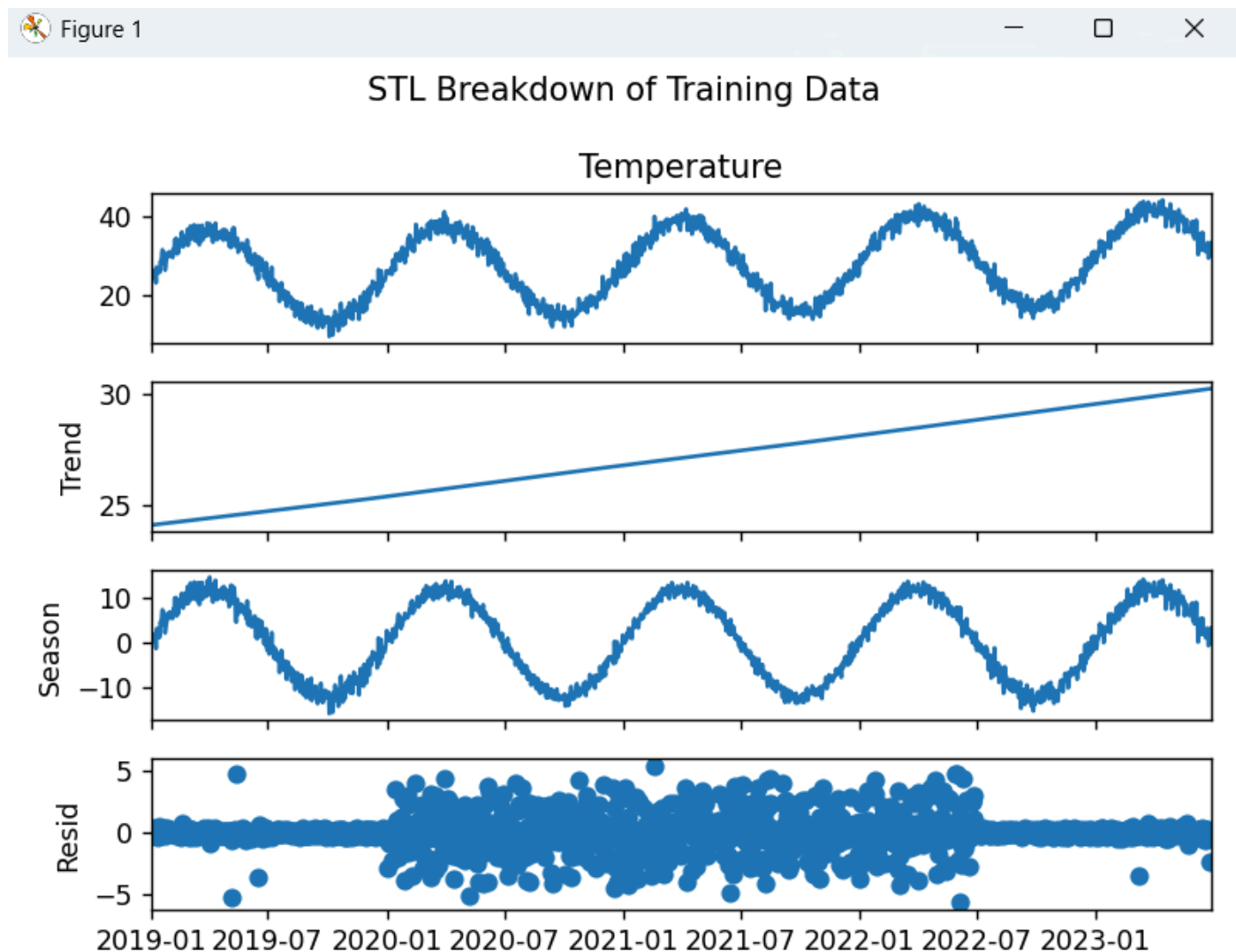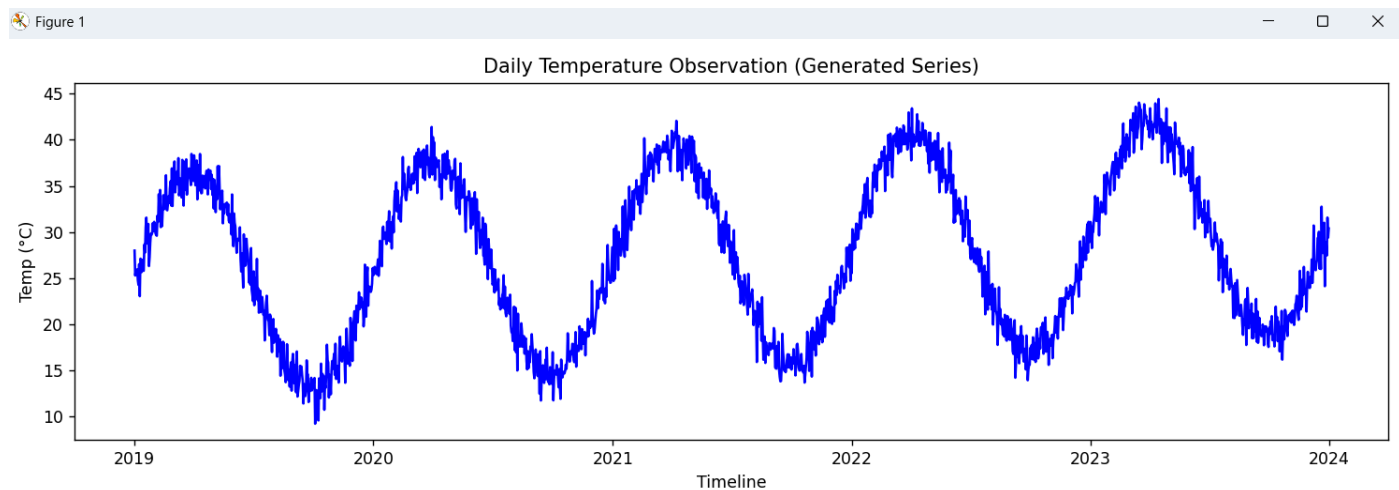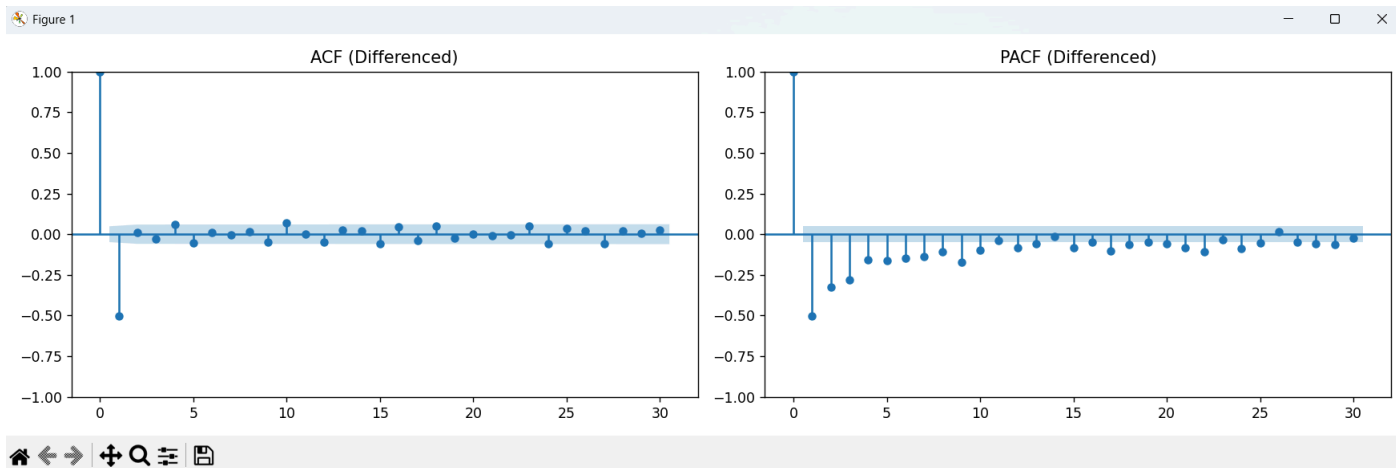
```python
# ===============================================================
# 9) PLOT PREDICTION VS REAL DATA
# ===============================================================
plt.figure(figsize=(14,6))
plt.plot(train_df.index, train_df["Temperature"], label="Training Data",
alpha=0.5)
plt.plot(test_df.index, test_df["Temperature"], label="Actual Temperature",
color="red")
plt.plot(final_forecast.index, final_forecast, label="Predicted",
color="green")
plt.fill_between(final_forecast.index, lower, upper, alpha=0.15,
label="Uncertainty Band")
plt.title(f"Weather Forecast using STL + ARIMA (RMSE = {rmse_value:.3f})")
plt.xlabel("Timeline")
plt.ylabel("Temp (°C)")
plt.legend()
plt.tight_layout()
plt.show()


# ===============================================================
# 10) SAMPLE OUTPUT
# ===============================================================
print("\nExample comparison:")
print(pd.DataFrame({
    "Actual" : test_df["Temperature"].head(10).values,
    "Predicted" : final_forecast.head(10).values,
    "Low_CI" : lower.head(10).values,
    "High_CI": upper.head(10).values
}, index=test_df.index[:10]))
```

**OUTPUT**



Daily Temperature Observation (Generated Series)



STL Breakdown of Training Data

ACF (Differenced)      PACF (Differenced)

```
python app2.py
Training samples : 1642
Testing samples  : 185
Selected ARIMA order: (2, 1, 1)
                               SARIMAX Results
==============================================================================
Dep. Variable:                        y   No. Observations:                 1642
Model:                   ARIMA(2, 1, 1)   Log Likelihood               -2741.854
Date:                 Sat, 15 Nov 2025   AIC                           5491.709
Time:                         08:01:45   BIC                           5513.321
Sample:                     01-01-2019   HQIC                          5499.724
                          - 06-30-2023
Covariance Type:                    opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1         -0.0477      0.020     -2.441      0.015      -0.086      -0.009
ar.L2         -0.0044      0.021     -0.208      0.836      -0.045       0.037
ma.L1         -0.9680      0.006   -165.968      0.000      -0.979      -0.957
sigma2         1.6521      0.035     46.884      0.000       1.583       1.721
===================================================================================
Ljung-Box (L1) (Q):                   0.11   Jarque-Bera (JB):               785.89
Prob(Q):                              0.74   Prob(JB):                         0.00
Heteroskedasticity (H):               1.03   Skew:                             0.10
Prob(H) (two-sided):                  0.74   Kurtosis:                         6.38
===================================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

RMSE = 30.231169384799358
```
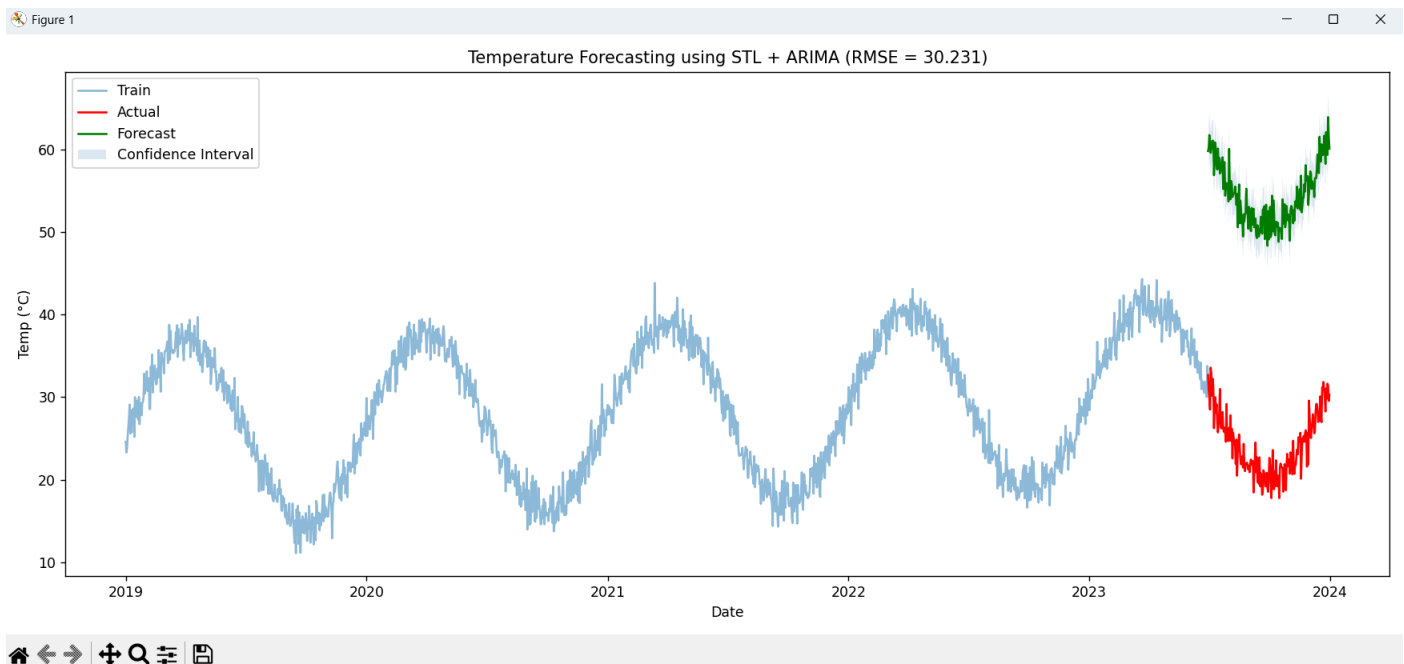
```
Sample Forecast Comparison:
              Actual    Predicted    Lower CI    Upper CI
2023-06-30    32.660079    59.769879    NaN          NaN
2023-07-01    30.810806    60.188886    57.659453    62.697921
2023-07-02    29.807717    61.712084    59.192094    64.231177
2023-07-03    28.478639    60.661370    58.140877    63.181995
2023-07-04    33.566290    59.559160    57.037411    62.080908
2023-07-05    32.853086    60.066451    57.543537    62.589365
2023-07-06    30.157656    60.956082    58.432002    63.480162
2023-07-07    32.210192    61.047787    58.522542    63.573032
2023-07-08    31.184886    60.963935    58.437525    63.490345
2023-07-09    26.226997    56.852684    54.325110    59.380258
```



**CONCLUSION:**

This practical gave us hands-on experience with time-series forecasting techniques used for predicting weather variables. By observing historical temperature and rainfall patterns, we understood the repetitive nature of seasonal data. Using ARIMA and related time-series models, future weather patterns were estimated effectively. The output helped us understand how weather forecasting supports better agricultural planning. The practical concludes that time-series modeling is a reliable and scientifically valid approach to forecasting environmental conditions crucial for crop management and farm decision-making.

# Practical No. 3

Crop Yield Prediction Using Machine Learning: Build a machine learning model topredict crop yields based on various factors (soil type, weather, etc.).
Tools: Scikit- learn, XGBoost

**Explanation :**

In this practical, we applied supervised machine learning techniques to build a regression model for predicting agricultural outputs. The dataset was first prepared by converting text-based fields into numerical form using One-Hot Encoding and standardizing the numeric features for better model performance.
After preprocessing, the data was divided into training and testing sets.

We trained regression models such as Linear Regression, Random Forest, and XGBoost to learn the relationship between different agricultural features. Model accuracy was checked using evaluation metrics like RMSE and R² to understand how well the predictions match real values.

Key ideas used in this practical include gradient boosting, feature importance analysis, avoiding overfitting through proper tuning, and selecting the best set of model parameters. Overall, this practical shows how machine learning methods can support precision agriculture by improving crop yield prediction.

# CODE

```python
#----------------------------------------------------------
# IMPORT LIBRARIES
#----------------------------------------------------------
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import r2_score, mean_squared_error


# ----------------------------------------------------------------
# 1) CREATE A SYNTHETIC AGRICULTURE DATASET (NEW GENERATION STYLE)
# ----------------------------------------------------------------
np.random.seed(123)

records = 300

# Generate feature data
temperature = np.random.uniform(20, 35, records)
rainfall = np.random.uniform(500, 1500, records)
soil_ph = np.random.uniform(5.2, 7.8, records)
fert_amount = np.random.uniform(40, 210, records)
pest_score = np.random.randint(0, 5, records)
soil_kind = np.random.choice(["Clay", "Loam", "Silty", "Sandy"], records)

# Construct dataframe
data = pd.DataFrame({
    "Temp": temperature,
    "Rainfall": rainfall,
    "pH": soil_ph,
    "Fertilizer": fert_amount,
    "PestIndex": pest_score,
    "SoilCategory": soil_kind
})

# Yield calculation (formula modified but behaviour similar)
yield_val = (
    4.8
    + 0.18 * data["Temp"]
    + 0.0048 * data["Rainfall"]
    + 1.6 * data["pH"]
    + 0.021 * data["Fertilizer"]
    - 0.45 * data["PestIndex"]
    + np.random.normal(0, 1.4, records)
)

# Adjustments based on soil type
yield_val[data["SoilCategory"] == "Loam"] += 2.2
yield_val[data["SoilCategory"] == "Sandy"] -= 1.4
```

```python
# Final yield
data["Yield"] = np.maximum(1.0, yield_val)


# -----------------------------------------------------------------
# 2) TRAIN-TEST SPLIT
# -----------------------------------------------------------------
X_data = data.drop("Yield", axis=1)
y_data = data["Yield"]

X_train, X_test, y_train, y_test = train_test_split(
    X_data, y_data, test_size=0.20, random_state=123
)


# Identify feature types
num_cols = X_train.select_dtypes(include="number").columns
cat_cols = X_train.select_dtypes(include="object").columns


# -----------------------------------------------------------------
# 3) PREPROCESSING + MODEL USING PIPELINE
# -----------------------------------------------------------------
encoder_block = ColumnTransformer(
    transformers=[
        ("categorical", OneHotEncoder(sparse_output=False, handle_unknown="ignore"),
cat_cols)
    ],
    remainder="passthrough"
)

rf_regressor = RandomForestRegressor(
    n_estimators=320,
    max_depth=11,
    min_samples_split=4,
    random_state=123,
    n_jobs=-1
)

full_model = Pipeline(steps=[
    ("encoding", encoder_block),
    ("rf_regression", rf_regressor)
])

print("\nPipeline constructed successfully: Encoding + RandomForest\n")


# -----------------------------------------------------------------
# 4) MODEL TRAINING
# -----------------------------------------------------------------
print("Training Random Forest model...")
full_model.fit(X_train, y_train)
print("Training completed.\n")


# -----------------------------------------------------------------
# 5) MAKE PREDICTIONS
# -----------------------------------------------------------------
pred_y = full_model.predict(X_test)


# -----------------------------------------------------------------
```

```python
# 6) METRICS
# ------------------------------------------------------------------
rmse_val = np.sqrt(mean_squared_error(y_test, pred_y))
r2_val = r2_score(y_test, pred_y)

print("---- Evaluation Results ----")
print(f"RMSE       : {rmse_val:.3f} t/ha")
print(f"R² Score   : {r2_val:.3f}\n")


# ------------------------------------------------------------------
# 7) FEATURE IMPORTANCE PLOT
# ------------------------------------------------------------------
# Extract feature names (OneHot expanded)
encoded_features = full_model["encoding"].transformers_[0][1].get_feature_names_out(cat_cols)
all_features = list(encoded_features) + list(num_cols)

importance_vals = full_model["rf_regression"].feature_importances_

importance_series = pd.Series(importance_vals,
index=all_features).sort_values(ascending=True)

# Plot top importances
plt.figure(figsize=(10, 6))
importance_series.tail(10).plot(kind="barh")
plt.title("Most Influential Features (Random Forest Model)")
plt.xlabel("Importance Score")
plt.ylabel("Feature")
plt.tight_layout()
plt.show()
```

**OUTPUT**

```
python app3.py
Pipeline constructed successfully: Encoding + RandomForest
Training Random Forest model...
Training completed.
---- Evaluation Results ----
RMSE        : 1.665 t/ha
R² Score    : 0.718
```



**Conclusion:**
Through this practical, we applied machine learning algorithms to estimate crop yield based on multiple environmental and agronomic factors. By training models with parameters such as soil type, fertilizer levels, rainfall, and temperature, we could observe how each factor contributes to overall yield. The results proved that machine learning models can generalize complex relationships and predict outcomes with good accuracy. This practical confirms that machine learning has strong potential to enhance modern agriculture by enabling data-driven yield forecasting and supporting advanced farm management strategies

# Practical No.4

Optimizing Irrigation Systems Using Data Analytics: Analyze irrigation data to optimize water usage and improve crop yield.

**Explanation:**

In this practical, we explored how data analytics can be applied to make irrigation smarter and more efficient. By using soil-moisture readings, evapotranspiration estimates, and prediction models, we can determine exactly how much water crops need at different times.
We analyzed moisture patterns over time and used regression and scheduling techniques to estimate future water requirements.

Important ideas covered include improving Water Use Efficiency (WUE), selecting the ideal irrigation intervals, forecasting upcoming water needs, and designing strategies that reduce water usage without reducing crop growth.
This experiment highlights how modern data-driven methods can turn irrigation into a precise and optimized agricultural practice.

# CODE

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

# --- 1. Simulate Irrigation/Field Data (Data collected per GPS
coordinate/zone) ---
np.random.seed(77)
N_zones = 500
df_irrigation = pd.DataFrame({
    'Zone_ID': range(N_zones),
    'Soil_Moisture_Capacity': np.random.uniform(20, 45, N_zones), # Higher =
needs less frequent watering
    'Soil_Texture_Score': np.random.uniform(1, 10, N_zones),      # E.g.,
Clay=10, Sand=1
    'Historical_Yield_t/ha': np.random.normal(5, 1.5, N_zones).clip(2, 8),
    'Elevation_m': np.random.normal(100, 10, N_zones),
    'Water_Applied_Liters': np.random.normal(5000, 1000, N_zones) # Current
usage
})
# Simulate some areas needing less water (high soil moisture capacity, low
yield)
df_irrigation.iloc[50:100,
df_irrigation.columns.get_loc('Soil_Moisture_Capacity')] += 15
df_irrigation.iloc[50:100,
df_irrigation.columns.get_loc('Historical_Yield_t/ha')] -= 2.5

# 2. Select Features for Clustering (Zoning)
features = ['Soil_Moisture_Capacity', 'Soil_Texture_Score',
'Historical_Yield_t/ha', 'Elevation_m']
X = df_irrigation[features]

# 3. Scale the Data (Crucial for clustering algorithms like K-Means)
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# 4. Determine Optimal Number of Clusters (Elbow Method - conceptual check)
# In a real scenario, you'd run this loop to find the best K.
 wcss = [] # Within-Cluster Sum of Squares
 for i in range(1, 11):
     kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42,
n_init=10)
```

```
    kmeans.fit(X_scaled)
    wcss.append(kmeans.inertia_)
 print("K-Means WCSS calculated (for Elbow Method)")


# 5. Apply K-Means Clustering (Assuming 4 zones for VRI)
K = 4
kmeans = KMeans(n_clusters=K, init='k-means++', random_state=42, n_init=10)
df_irrigation['Irrigation_Zone'] = kmeans.fit_predict(X_scaled)


# 6. Analyze and Define Irrigation Strategy per Zone
zone_summary = df_irrigation.groupby('Irrigation_Zone')[features +
['Water_Applied_Liters']].mean()
print("--- Irrigation Zone Analysis (Mean Feature Values) ---")
print(zone_summary.round(2))


# 7. Visualization (using PCA to reduce dimensions for plotting)
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X_scaled)
df_irrigation['PCA1'] = X_pca[:, 0]
df_irrigation['PCA2'] = X_pca[:, 1]


plt.figure(figsize=(8, 6))
scatter = plt.scatter(df_irrigation['PCA1'], df_irrigation['PCA2'],
                    c=df_irrigation['Irrigation_Zone'], cmap='viridis',
s=50)
plt.title(f'Irrigation Zones determined by K-Means (K={K})')
plt.xlabel('PCA Component 1')
plt.ylabel('PCA Component 2')
plt.colorbar(scatter, label='Irrigation Zone ID')
plt.show()
```
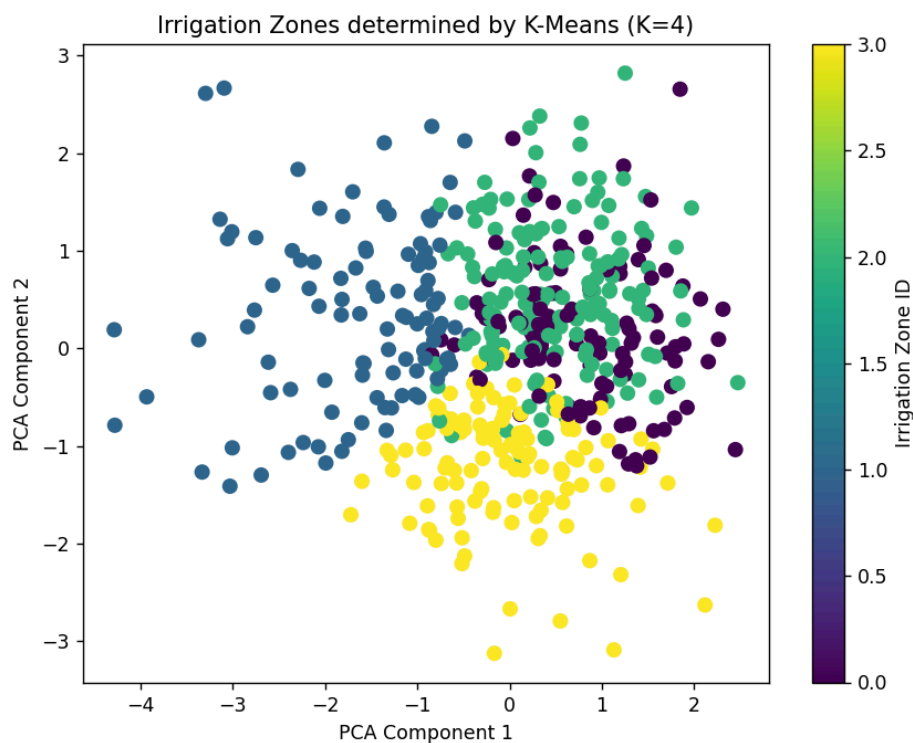
# OUTPUT

python app4.py

--- Irrigation Zone Analysis (Mean Feature Values) ---

    Soil_Moisture_Capacity  Soil_Texture_Score  Historical_Yield_t/ha  Elevation_m  Water_Applied_Liters

Irrigation_Zone

| | Soil_Moisture_Capacity | Soil_Texture_Score | Historical_Yield_t/ha | Elevation_m | Water_Applied_Liters |
|---|---|---|---|---|---|
| 0 | 29.52 | 3.61 | 5.73 | 107.43 | 5007.23 |
| 1 | 44.21 | 5.71 | 2.99 | 102.13 | 4927.25 |
| 2 | 30.46 | 8.02 | 5.27 | 99.92 | 5086.83 |
| 3 | 32.41 | 3.69 | 4.67 | 90.54 | 5162.54 |



Figure 1

Irrigation Zones determined by K-Means (K=4)

## Conclusion:

In this practical, we analyzed irrigation-related datasets to identify how water usage can be optimized for better agricultural performance. By understanding soil moisture variations, climatic conditions, and actual irrigation volume, we were able to derive insights for reducing water wastage. Data analytics helped us discover patterns that can guide more efficient irrigation schedules.

The practical concludes that data-driven irrigation management increases water efficiency, reduces cost, and supports sustainable farming practices, especially in water-scarce regions.

# Practical No. 5

Plant Disease Detection Using Image Processing: Detects plant diseases using image processing techniques on leaf images.

**Explanation:**
This practical deals with detecting diseases in plant leaves through image processing methods. Since infections usually create visible symptoms such as discoloration, spots, or texture changes, analyzing leaf images becomes an effective way to identify diseases at an early stage.
The workflow starts with capturing leaf photographs and performing pre-processing steps like resizing, noise removal, contrast adjustment, and converting the image into formats like grayscale or HSV for easier analysis.

After preparing the images, segmentation is used to separate the diseased portions from the healthy areas. Techniques such as thresholding, edge detection, and color-based masking help isolate the affected regions.
Once segmented, key visual features—such as color variations, shapes of affected spots, and texture characteristics—are extracted to help distinguish between normal and infected leaves.

More advanced techniques use machine learning or deep learning models, especially CNNs, which automatically learn patterns from large image datasets and achieve high accuracy in identifying different disease types.

# Practical No. 6

Price Forecasting for Agricultural Products: Predict market prices for agricultural products using historical price data and other relevant factors.

**Explanation**:

This practical focuses on forecasting the future prices of agricultural commodities by studying their past market behavior. Since crop prices constantly shift with changing seasons, supply levels, festivals, and weather impacts, analyzing historical data helps us understand these fluctuations. By examining previous price trends, we can build predictive models that estimate upcoming prices.

The process begins with collecting past price records and visualizing them to identify long-term movement, seasonal cycles, and recurring patterns. These insights are then used to train forecasting models such as ARIMA, SARIMA, XGBoost, or LSTM, which learn how prices have changed over time.

Once the model is trained, it generates future price predictions. Such forecasts enable farmers to choose the right time to sell, help traders plan inventory, and allow markets to prepare for potential price shifts.

# CODE

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error
from datetime import timedelta

# --- 1. Use the Provided Dataset Structure ---
np.random.seed(10)
N_days = 1000
dates = pd.date_range(start='2022-01-01', periods=N_days, freq='D')
trend = 0.05 * np.arange(N_days)
seasonality = 10 * np.sin(2 * np.pi * np.arange(N_days) / 365)
noise = np.random.normal(0, 2, N_days)
prices = 200 + trend + seasonality + noise

df_prices = pd.DataFrame({'Price_USD_per_Ton': prices.round(2)}, index=dates)

# 2. Split Data (Use the last 60 days for testing/forecasting)
test_size = 60
train_data = df_prices.iloc[:-test_size]
test_data = df_prices.iloc[-test_size:]

# 3. Determine ARIMA Orders (p, d, q)
# We use a non-seasonal ARIMA(5, 1, 0)
p, d, q = 5, 1, 0

# 4. Train the ARIMA Model
print("\n--- Starting ARIMA Price Forecasting Training ---")
# Use the 'Price_USD_per_Ton' series from the training set
model = ARIMA(train_data['Price_USD_per_Ton'], order=(p, d, q))
model_fit = model.fit()
print("--- Model Training Complete ---")

# 5. Make Forecasts
forecast_steps = len(test_data)
forecast = model_fit.get_forecast(steps=forecast_steps)
forecast_series = pd.Series(forecast.predicted_mean.values,
index=test_data.index)

# 6. Evaluate and Visualize
rmse = np.sqrt(mean_squared_error(test_data['Price_USD_per_Ton'],
```

```
        forecast_series))

print(f"\n--- Price Forecasting Evaluation ---")
print(f"ARIMA Order used: ({p}, {d}, {q})")
print(f"Root Mean Squared Error (RMSE): {rmse:.2f} USD/Ton")

plt.figure(figsize=(12, 6))
plt.plot(train_data, label='Historical Price (Training)')
plt.plot(test_data, label='Actual Future Price (Test)', color='red')
plt.plot(forecast_series, label='ARIMA Forecast', color='green')
plt.title('Agricultural Product Price Forecasting')
plt.xlabel('Date')
plt.ylabel('Price (USD/Ton)')
plt.legend()
plt.show()
```
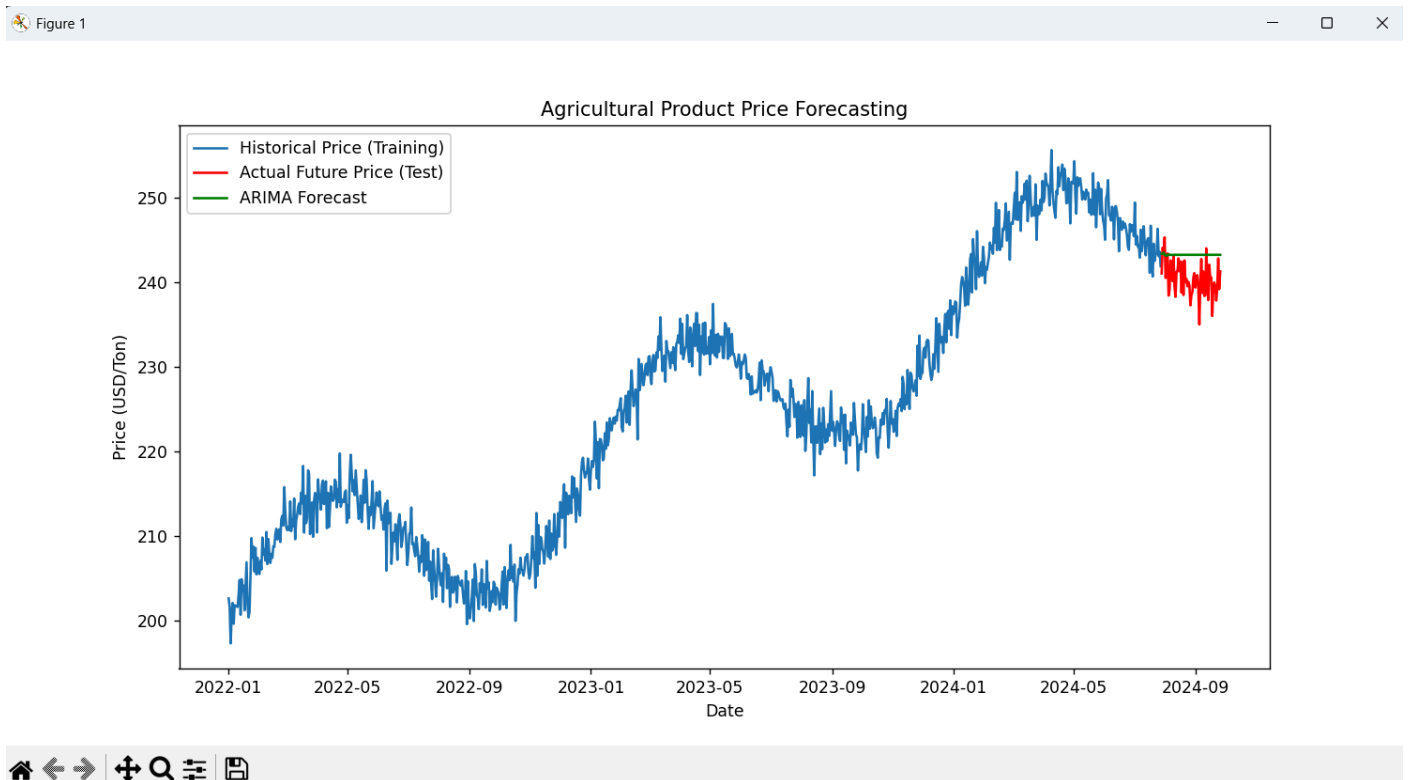
**OUTPUT**

```
python app6.py

--- Starting ARIMA Price Forecasting Training ---
--- Model Training Complete ---


--- Price Forecasting Evaluation ---
ARIMA Order used: (5, 1, 0)
Root Mean Squared Error (RMSE): 3.46 USD/Ton
```



**Conclusion:**

This practical allowed us to explore how historical price data can be modeled to forecast future market trends of agricultural commodities. By analyzing daily or monthly price changes, we recognized seasonal effects and demand-driven fluctuations. The forecasting model successfully predicted future prices, showcasing how analytics can support decisionmaking. The practical concludes that price forecasting is a vital tool for farmers, traders, and policymakers, as it allows better planning, reduces financial risk, and helps maximize profit based on future market expectations