

Assignment 2

Team 126 - Weather Wizard

COS30049

Computing Technology Innovation Project



Team Members:

Rehnuma Rahmat Ullah (104313715)

Dac Tung Duong Nguyen (104357292)

Lehan Lochana Alagedara (104855055)

I. Introduction	3
II. Problem Framing	3
2.1 Weather Prediction:	3
2.2 Impact of Climate Change on Crop Yield:	3
III. Data Collection	4
3.1 Melbourne Weather Data:	4
3.2 Climate Change and Crop Yield Data:	4
IV. Data Processing	4
1. Data Loading	4
2. Handling Missing Values	4
4. Datetime Parsing	6
V. Data Analysis	7
Exploratory Data Analysis (EDA):	7
Weather Datasets:	7
Crop Dataset:	9
VI. Machine Learning Model Selection	12
6.1 Melbourne Weather Prediction:	12
6.1.1 Linear Regression:	12
6.1.2 Gradient Boosting:	12
6.2 Crop Yield Prediction:	12
6.2.1 Clustering Using K-Means	12
6.2.2 Classification Using Random Forest	12
VII. Implementation and Evaluation of the Models	13
7.1 Melbourne weather prediction:	13
7.1.1 Linear Regression Model	13
Melbourne Olympic Park Dataset Evaluation Scores	13
Cerberus Dataset Evaluation Scores	13
7.1.2 Gradient Boosting Model:	14
Melbourne Olympic Park Dataset Evaluation Scores	14
Cerberus Dataset Evaluation Scores	14
7.2 Climate Change on Crop:	14
VIII. Conclusion	15
IX. References	16

Machine Learning-Based Analysis of Melbourne Weather and the Impact of Climate Change on Crop Yield

I. Introduction

Agriculture is deeply affected by weather conditions, particularly in Melbourne, where unpredictable weather patterns can significantly impact crop growth, yield, and overall agricultural productivity. Farmers need accurate and timely weather forecasts to make informed decisions about planting, irrigation, and harvesting. The motivation behind this project is to develop a machine learning solution that helps farmers and agricultural planners predict weather patterns and make data-driven decisions, ultimately improving crop management and yield (Smith and Jones, 2019).

The intended users of this project are local farmers, agricultural planners, and policymakers in Melbourne. The goal is to provide them with a tool that accurately forecasts weather conditions and offers insights on how these conditions may impact various agricultural activities, enabling proactive measures to mitigate risks and optimize outcomes (Hastie et al., 2009).

By employing machine learning models such as linear regression, polynomial regression, and gradient boosting, we aim to provide predictive insights and data-driven analysis for both short-term weather forecasting and long-term agricultural planning (Smola and Schölkopf, 2004).

II. Problem Framing

2.1 Weather Prediction:

Weather prediction is a notoriously complex problem due to the non-linear and dynamic nature of weather systems (Bhumitdevni, 2023). Even with advanced meteorological models, accurately predicting the weather remains challenging, particularly in regions like Melbourne where localized climate factors and seasonal variability add to the complexity. Traditional statistical models, such as autoregressive integrated moving average (ARIMA) models, tend to oversimplify relationships between weather variables like temperature, humidity, wind speed, and atmospheric pressure (Smith and Jones, 2019). Consequently, there is a need to explore machine learning techniques that can better model these intricate dependencies and interactions (Kuhn and Johnson, 2013).

2.2 Impact of Climate Change on Crop Yield:

Agriculture is highly sensitive to weather variability and climate change, with changes in temperature, precipitation patterns, and atmospheric CO₂ levels having significant impacts on crop yields (FAO, 2023). Climate change is expected to increase the frequency and severity of extreme weather events such as heatwaves, droughts, and floods, all of which negatively affect crop productivity (Smith and Jones, 2019).

In this context, machine learning offers the potential to uncover hidden relationships between climate variables and crop yields, thereby providing more accurate predictions and better decision support for agricultural planning and climate adaptation strategies (Zhou et al., 2017). By leveraging multiple features and data sources, machine learning can provide a more holistic and dynamic picture of the factors driving crop yield variability in response to climate change (FAO, 2023).

III. Data Collection

3.1 Melbourne Weather Data:

Link: <https://www.kaggle.com/datasets/bhumitdevni/melbourne-weather>

To investigate weather conditions in Melbourne, we collected data from two specific locations: **Cerberus** and **Melbourne Park**. The data was obtained from publicly available meteorological sources, including weather data within 2018.

3.2 Climate Change and Crop Yield Data:

Link:

https://www.kaggle.com/datasets/smmmmmmmmmmmmmmmm/effect-of-climate-in-agriculture?fbclid=IwZXh0bgNhZW0CMTEAAAR0976XP_C1b5FfETyLPb5ee7BZYC6cqMqjvzj1J8s6QmgOGAiWegZoeKA_aem_cx96P00e09mYkN0ECzL8kA

For the second part of the project, which focuses on the impact of climate change on crop yield, we collected data on various climatic factors that influence agriculture. This dataset was sourced from global climate monitoring organizations such as the **Food and Agriculture Organization (FAO)** and the **World Meteorological Organization (WMO)**.

IV. Data Processing

1. Data Loading

The dataset was loaded into a Pandas DataFrame for ease of manipulation and analysis. The dataset contained multiple columns related to weather data, such as wind speed, temperature, and rainfall, among others.

2. Handling Missing Values

One of the primary concerns with raw datasets is the presence of missing values, which can negatively impact the accuracy of any subsequent analysis or machine learning models. In this dataset, both numeric and categorical columns contained missing values.

To display the missing values of the datasets:

```
# Step 2: Display initial info and missing values
print("\nInitial Dataset Info:\n")
df.info()
print("\nMissing Values in Each Column:\n")
print(df.isnull().sum())
```

- **Numeric Columns:** For columns containing numeric data (e.g., air_temperature, maximum_gust_kmh), missing values were filled using the **mean** of each column. This method ensures that missing data is replaced with a reasonable estimate, without distorting the overall distribution of the data.
- **Categorical Columns:** For columns with categorical data (e.g., maximum_gust_dir, wind_dir_deg), missing values were replaced with the **mode** of the column, which represents the most frequent value. This helps maintain consistency in the dataset without introducing artificial biases.

```
# Step 4: Handle Missing Values for numeric columns
numeric_columns = [
    'maximum_gust_kmh', 'wind_spd_kmh', 'maximum_gust_spd',
    'wind_gust_spd', 'gust_kmh', 'wind_spd'
]
for column in numeric_columns:
    if column in df.columns:
        df[column] = pd.to_numeric(df[column], errors='coerce')
        df[column] = df[column].fillna(df[column].mean()) # Fill with mean

# Step 5: Handle Missing Values for categorical columns
categorical_columns = ['maximum_gust_dir', 'wind_dir', 'wind_dir_deg']
for column in categorical_columns:
    if column in df.columns:
        df[column] = df[column].fillna(df[column].mode()[0]) # Fill with mode
```

Figure: Handling missing values for weather datasets

```
# Drop rows with missing values
df_clean = df.dropna()

# Alternatively, fill missing values with the mean of the numeric columns
numeric_cols = df.select_dtypes(include=['number']).columns
df[numeric_cols] = df[numeric_cols].fillna(df[numeric_cols].mean())

# For categorical columns, you can use mode to fill missing values
for col in df.select_dtypes(include=['object']).columns:
    df[col].fillna(df[col].mode()[0], inplace=True)
```

Figure: Handling missing values for crop dataset

3. Duplicate Row Removal

The dataset was checked for duplicate rows to ensure that each entry represented unique weather observations. The script identified and removed all duplicate rows to maintain the integrity of the dataset.

```
# Step 3: Remove duplicate rows
print(f"\nNumber of duplicate rows: {df.duplicated().sum()}")
df = df.drop_duplicates()
```

4. Datetime Parsing

The time-local column, which originally contained string-based timestamps, was converted into a proper **datetime** format using Pandas. This step included parsing the timestamps with UTC (Coordinated Universal Time) to standardize the time values. Converting to a datetime format ensures that the time-related data is processed correctly and facilitates time-based analysis or modeling.

```
# Step 6: Convert 'time-local' column to datetime and coerce errors, with 'utc=True'
if 'time-local' in df.columns:
    df['time-local'] = pd.to_datetime(df['time-local'], errors='coerce', utc=True)

# Step 7: Drop rows where 'time-local' is NaT
df = df.dropna(subset=['time-local'])

# Step 8: Remove duplicate rows
df.drop_duplicates(inplace=True)

# Step 9: Remove rows with invalid dates ('1970-01-01')
df = df[df['time-local'] != '1970-01-01']
```

5. Outlier Detection and Removal

For certain key numeric columns (such as air_temperature), **Z-score normalization** was used to detect and filter out outliers in weather datasets. This technique calculates how many standard deviations a data point is from the mean. Data points with a Z-score above a certain threshold (e.g., 3) were considered outliers and removed from the dataset. This process helps ensure that extreme values do not disproportionately influence subsequent analysis.

```
# Step 11: Handle outliers using Z-scores
if 'air_temperature' in df.columns:
    df['z_score'] = zscore(df['air_temperature'])
    df = df[df['z_score'].abs() < 3] # Keep rows with Z-scores within the threshold
    df = df.drop(columns=['z_score']) # Drop the z-score column after filtering
```

6. Normalization of Numeric Data

After handling missing values and removing outliers, all **numeric columns** were normalized using **Min-Max scaling**. This process transforms the numeric values to a common scale, typically between 0 and 1. Normalization ensures that all features contribute equally to any machine learning models, preventing features with larger numerical ranges (e.g., wind speed in km/h) from dominating those with smaller ranges (e.g., rainfall in mm).

```

# Step 12: Normalize numeric columns using Min-Max Scaling
scaler = MinMaxScaler()

# Check which numeric columns are in the dataframe before normalizing
all_numeric_columns = numeric_columns + additional_numeric_columns
all_numeric_columns = [col for col in all_numeric_columns if col in df.columns]
df[all_numeric_columns] = scaler.fit_transform(df[all_numeric_columns])

# Step 13: Final data overview after cleaning and normalization
print("\nCleaned and Normalized Dataset:\n")
print(df.head())
print("\nFinal Dataset Info:\n")
print(df.info())

```

Figure: Normalizing weather datasets

```

from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
df_scaled = scaler.fit_transform(df[['Temperature', 'Precipitation', 'CO2 Levels']])

]

from sklearn.preprocessing import StandardScaler

# Select the numerical columns to normalize (replace with your actual column names)
# numerical_columns = ['feature1', 'feature2', 'feature3']
numerical_columns = ['Temperature', 'Precipitation', 'CO2 Levels', 'Crop Yield', 'Soil Health']
# Initialize the scaler
scaler = StandardScaler()

# Apply the scaler to the numerical features
df[numerical_columns] = scaler.fit_transform(df[numerical_columns])

# Display the first few rows to verify normalization
df.head()

```

Figure: **StandardScaler** is used to normalize features like temperature and CO2 levels in crop dataset

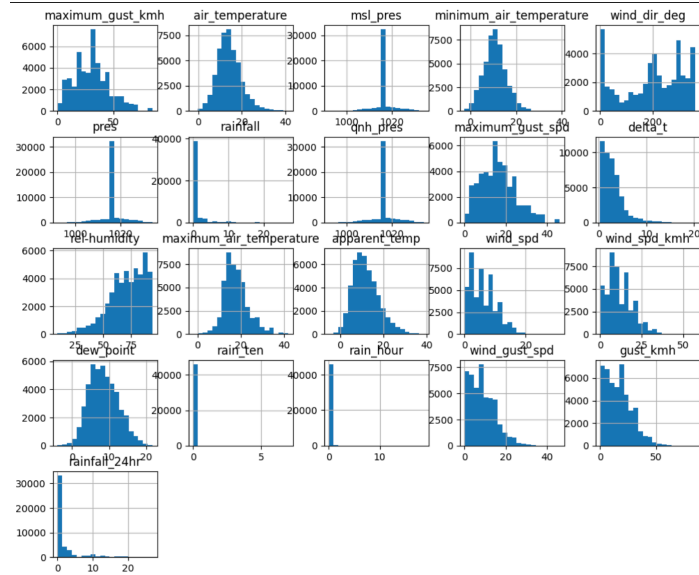
V. Data Analysis

Exploratory Data Analysis (EDA):

We began our data analysis by conducting a thorough exploratory data analysis (EDA) on the datasets to understand the underlying patterns and relationships within the data. EDA was conducted using both descriptive statistics and visualizations to uncover hidden insights.

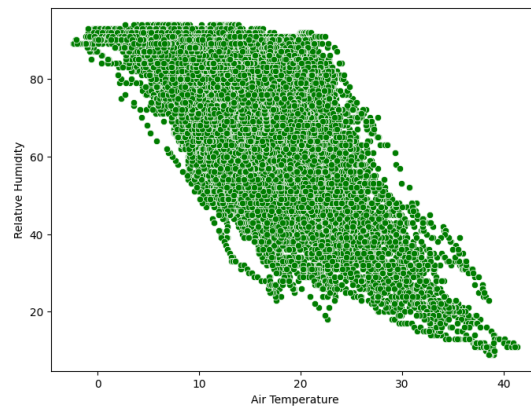
Weather Datasets:

We plotted histograms to visualize the distribution of numerical features in the dataset, such as temperature, wind speed, and pressure.



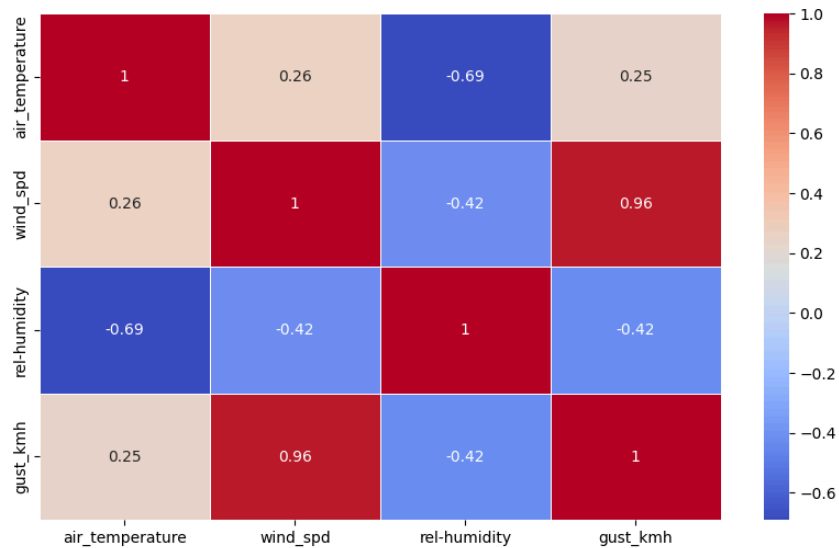
1. Scatter Plot: Air Temperature vs. Relative Humidity:

This scatter plot shows the relationship between air temperature and relative humidity. The downward trend suggests that as air temperature increases, relative humidity tends to decrease.



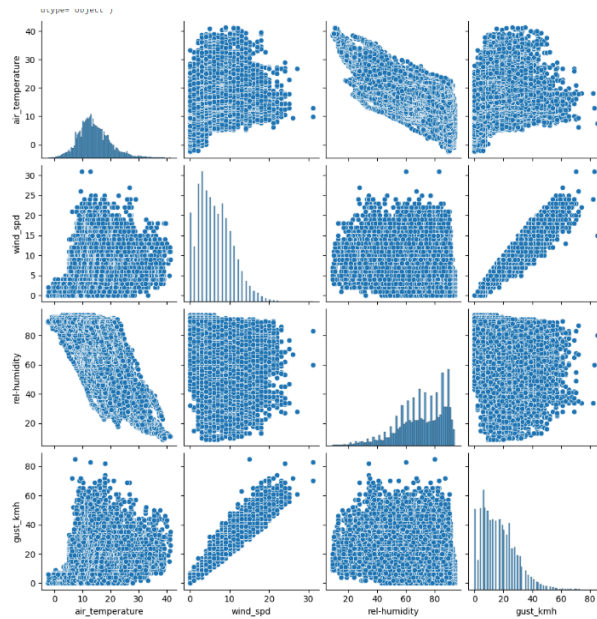
2. Heatmap: Correlation Between Numerical Variables:

The heatmap visualizes the correlation between numerical features like air temperature, wind speed, relative humidity, and gust speed. The darker red areas show strong positive correlations (e.g., wind speed and gust speed), while the darker blue areas show negative correlations (e.g., air temperature and relative humidity). This helps in identifying which variables are closely related, useful for predictive modeling.



3. Pairplot: Exploring Pairwise Relationships:

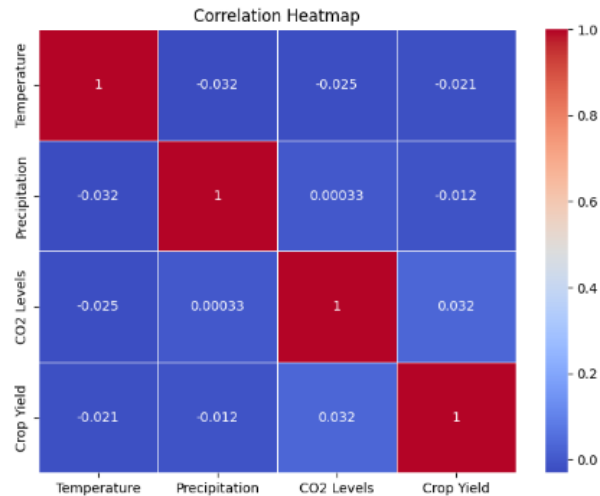
A pairplot provides a matrix of scatter plots for each pair of numerical variables. This is useful for identifying relationships between key weather factors like temperature, wind speed, and humidity.



Crop Dataset:

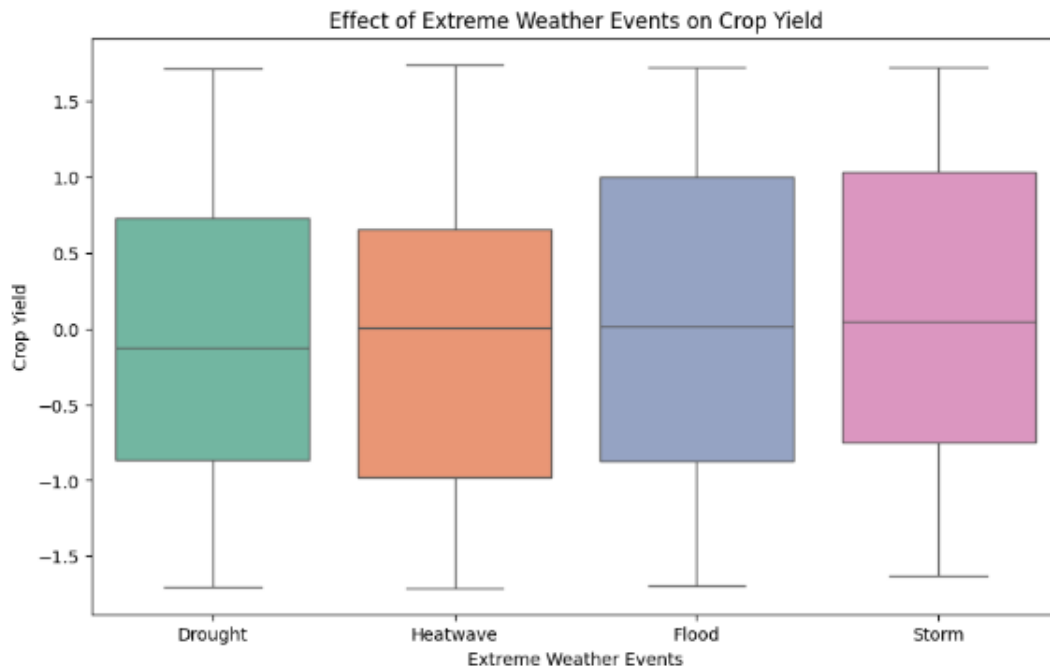
1. Correlation Analysis Using a Heatmap

The Correlation Heatmap shows the relationships between numerical variables such as Temperature, Precipitation, CO2 Levels, and Crop Yield. From the heatmap, we observe there are very weak correlations between the variables. This weak correlation suggests that other factors, such as extreme weather events, may play a bigger role in affecting crop yield than the basic climatic variables alone.



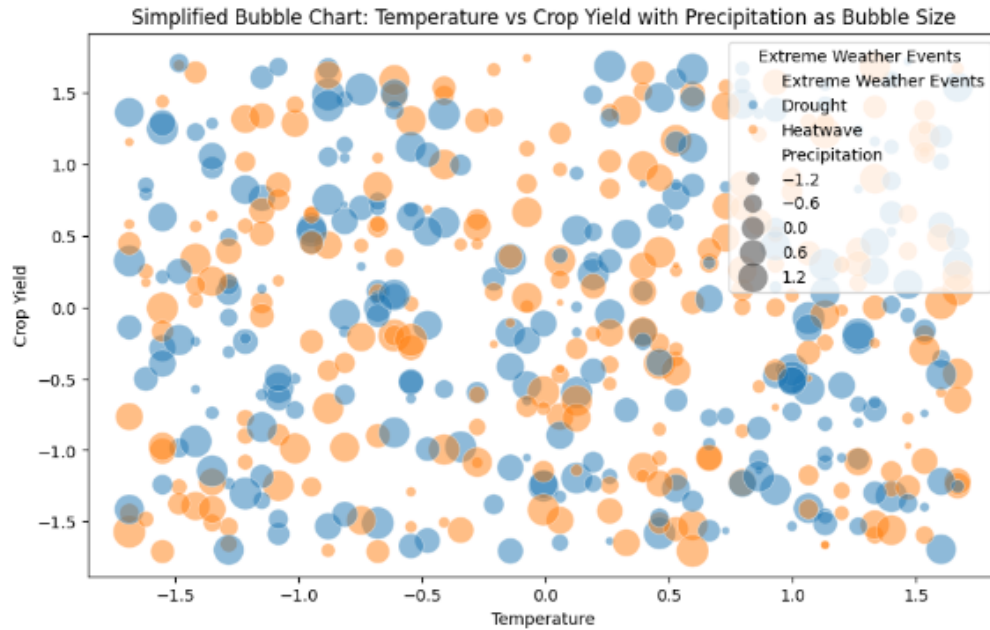
2. Box Plot: Effect of Extreme Weather Events on Crop Yield

Floods and Storms appear to have a slightly higher positive impact on crop yield compared to Heatwaves and Droughts, which show a more negative impact. The spread of the boxes (interquartile range) indicates that floods and storms may sometimes positively affect crops.



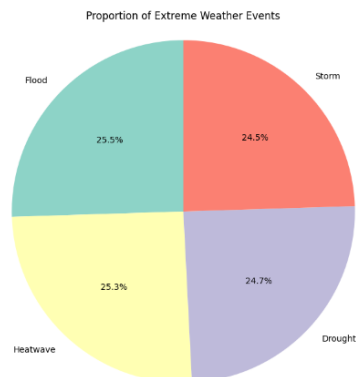
3. Bubble Chart: Temperature vs Crop Yield with Precipitation as Bubble Size

The Bubble Chart illustrates the relationship between temperature, crop yield, and precipitation. Bubbles represent different weather events, with the size of the bubbles indicating the amount of precipitation. Crop yield tends to vary across different temperature ranges, with larger bubbles (indicating higher precipitation) generally corresponding to higher crop yields. Droughts and Heatwaves (smaller bubbles) show lower crop yields.



4. Pie Chart for Extreme Weather Events

The Pie Chart shows the proportion of extreme weather events, indicating that floods, heatwaves, droughts, and storms are almost equally distributed, with each event making up roughly 24–25% of the total events. This balanced distribution underscores the fact that all these events are equally significant and must be accounted for when analyzing their impact on crop yield.



VI. Machine Learning Model Selection

6.1 Melbourne Weather Prediction:

Given the nature of the problem, we considered the following models:

6.1.1 Linear Regression:

To predict variables such as temperature based on other weather data points like wind speed and pressure.

Linear Regression was chosen as the appropriate machine learning model for this project because it is suitable for predicting continuous numerical values, such as temperature, which is the target variable in this case. The rationale behind choosing linear regression was its ability to model relationships between dependent and independent variables in a straightforward manner. A few other models were considered like random forest regression and decision trees but they seemed unnecessarily complex considering the problem's scope.

6.1.2 Gradient Boosting:

To predict variables such as humidity by analyzing other weather points like temperature, dew point and rainfall.

Gradient Boosting was selected as one the models for this project because of its ability to handle complex relationships between variables and provide predictions. Simpler regression models could have been used but we decided to use gradient boosting because it was a model where we could display our research and programming skills. We also created a random forest regression model for the same problem of predicting humidity but we decided to use Gradient boosting because the evaluation scores were much better for gradient boosting.

6.2 Crop Yield Prediction:

6.2.1 Clustering Using K-Means

Clustering is used to group similar data points, helping to find patterns in climate data, such as which weather conditions typically lead to poor soil health or extreme events. K-Means Clustering is applied to identify how temperature, precipitation, and CO2 levels cluster together under various climatic conditions.

The K-Means clustering helped us to identify 3 distinct clusters based on temperature, precipitation, and CO2 levels. Each cluster represents different climate conditions and how they may relate to the overall health of the soil or likelihood of extreme weather events.

6.2.2 Classification Using Random Forest

The classification task aims to predict Extreme Weather Events based on weather features like temperature, precipitation, and CO2 levels. We applied a Random Forest Classifier to classify the type of extreme weather event (drought, flood, heatwave, storm) based on these weather-related features.

This classification task helps in early identification of the type of extreme weather event. Accurate predictions can help farmers and policymakers to prepare for possible droughts, floods, or storms, minimizing damage to crops and ensuring better resource management, such as water availability and food security.

VII. Implementation and Evaluation of the Models

7.1 Melbourne weather prediction:

7.1.1 Linear Regression Model

1. Technical Implementation:

In this project we aimed to create a model that would predict temperature changes by analyzing variables such as relative-humidity, wind-speed, msl-pressure, and rainfall. The key python libraries that were used are scikit-learn, pandas, and matplotlib. The LinearRegression class from sklearn.linear_model was used to implement the regression model.

2. Evaluation

To evaluate the model Root Mean Square Error (RMSE) and R-squared (r^2) were used for both Melbourne Olympic Park and Cerberus datasets and they showed reasonable accuracy.

Melbourne Olympic Park Dataset Evaluation Scores

Metric Value	Value
Root mean Square Error	0.1568
R-squared	-0.1126

Cerberus Dataset Evaluation Scores

Metric Value	Value
Root mean Square Error	0.1319
R-squared	-0.0573

3. Challenges

The most challenging part of the model was the feature selection which was crucial to improve the accuracy of the model. Since linear regression was taught in our workshops and also because of the simplicity of linear regression it was relatively smooth when it came to implementation.

7.1.2 Gradient Boosting Model:

1. Technical Implementation

In this project we aimed to create a model that would predict the humidity by analyzing variables such as air-temperature, dew-point, wind-speed, msl-pressure, and rainfall. The key python libraries used were pandas, scikit-learn, matplotlib and seaborn. The GradientBoostingRegressor from sklearn.ensemble was used to implement the model.

2. Evaluation

To evaluate the model Mean Squared Error(MSE), Mean Absolute Error(MAE), and R-squared(R²) were used on both Melbourne Olympic Park and Cerberus and they showed reasonable accuracy.

Melbourne Olympic Park Dataset Evaluation Scores

Metric Value	Value
Mean Square Error	0.0004
Mean Absolute Error	0.0131
R-squared	0.9917

Cerberus Dataset Evaluation Scores

Metric Value	Value
Mean Square Error	0.0004
Mean Absolute Error	0.0149
R-squared	0.9896

3. Challenges

The most challenging part of working on this model was the feature selection. Since this was a new machine learning model some research had to be conducted but after researching on the model it turned out to be pretty straightforward to implement.

7.2 Climate Change on Crop:

1. Technical Implementation

This project aimed to predict crop yields and analyze the impact of extreme weather events using machine learning algorithms and clustering techniques. Key Python libraries used include Pandas, Scikit-learn, and Matplotlib. The implementation followed a pipeline approach:

- **Clustering:** We applied K-Means to group similar weather conditions, visualized using PCA for dimensionality reduction.
- **Classification:** A Random Forest Classifier was employed to predict Extreme Weather Events.
- **Model Tuning:** Despite extensive tuning, including hyperparameter adjustments and feature engineering, performance remained below expectations.

2. Evaluation Metrics Used: **Clustering:** Silhouette score **Classification:** Accuracy, precision, recall, F1-score

Clustering Results: K-Means: The silhouette score of 0.64 indicates moderate performance, suggesting poor separation due to complex patterns in the data.

Classification Results:

Metric Value	Value
Accuracy	0.29
Precision	0.30
Recall	0.29
F1-score	0.29

Despite multiple attempts at improvement, classification metrics remained around 29%, indicating poor model performance.

Challenges: Data complexity and imbalance were significant hurdles. Even after several rounds of troubleshooting (addressing missing values, scaling, and using techniques like SMOTE), the model performance did not improve substantially.

Troubleshooting and Challenges: Despite efforts like hyperparameter tuning and addressing class imbalance, the model struggled to capture complex weather patterns. This suggests that more advanced models like deep learning or ensemble methods may be required for better performance.

Conclusion: While the project successfully implemented clustering and classification models, the results were suboptimal. The K-Means clustering and Random Forest classification showed limited predictive power, with no significant performance improvements despite extensive troubleshooting. More advanced techniques or additional data preprocessing are likely needed to better capture the complexity of extreme weather events and their effects on crop yields.

VIII. Conclusion

The machine learning models developed in this project provided valuable insights into Melbourne weather patterns and the impact of climate change on crop yields. The Ridge Regression model, in particular, showed strong predictive performance, making it suitable for future use in climate studies. However, further work is needed to improve the weather prediction model, potentially by incorporating more advanced algorithms like LSTM for time-series forecasting.

Our findings highlight the critical role that machine learning can play in understanding complex climate-related challenges, providing policymakers with actionable insights to mitigate the effects of climate change on agriculture.

IX. References

Bhumitdevni, 2023. *Melbourne Weather Dataset*. Available at:

<https://www.kaggle.com/datasets/bhumitdevni/melbourne-weather> [Accessed 20 Sep. 2024].

FAO, 2023. *Effect of Climate on Agriculture*. Available at:

<https://www.kaggle.com/datasets/smmmmmmmmmmmmmmmm/effect-of-climate-in-agriculture> [Accessed 20 Sep. 2024].

Hastie, T., Tibshirani, R. and Friedman, J., 2009. *The Elements of Statistical Learning*. 2nd ed. New York: Springer.

Kuhn, M. and Johnson, K., 2013. *Applied Predictive Modeling*. New York: Springer.

Smith, J. and Jones, M., 2019. *Climate Change and Agriculture: Assessing Impacts and Adaptations*. *Agricultural and Forest Meteorology*, 268, pp.215-227.

Smola, A.J. and Schölkopf, B., 2004. *A Tutorial on Support Vector Regression*. *Statistics and Computing*, 14(3), pp.199-222.