

# Configuration Manual

A Machine Learning Framework for Predicting Crop Production  
in Support of SDG13 Climate Action

MSc Research Project  
Data Analytics

Amrit Laxmanasa Shidling  
Student ID: x21198951@student.ncirl.ie

School of Computing  
National College of Ireland

Supervisor: Dr. Paul Stynes



**National College of Ireland  
Project Submission Sheet  
School of Computing**



<b>Student Name:</b>	Amrit Laxmanasa Shidling
<b>Student ID:</b>	x21198951@student.ncirl.ie
<b>Programme:</b>	Data Analytics
<b>Year:</b>	2023
<b>Module:</b>	MSc Research Project
<b>Supervisor:</b>	Dr. Paul Stynes
<b>Submission Due Date:</b>	14/12/2023
<b>Project Title:</b>	Configuration Manual
<b>Word Count:</b>	961
<b>Page Count:</b>	11

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** Internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use another author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

<b>Signature:</b>	Amrit Laxmanasa Shidling
<b>Date:</b>	31st January 2024

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

Attach a completed copy of this sheet to each project (including multiple copies).	<input checked="" type="checkbox"/>
<b>Attach a Moodle submission receipt of the online project submission</b> , to each project (including multiple copies).	<input checked="" type="checkbox"/>
<b>You must ensure that you retain a HARD COPY of the project</b> , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input checked="" type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

<b>Office Use Only</b>	
Signature:	
Date:	
Penalty Applied (if applicable):	

# Configuration Manual

Amrit Laxmanasa Shidling  
x21198951@student.ncirl.ie

## 1 Introduction

This configuration manual provides thorough information about hardware, and software requirements to run the project "A Machine Learning Framework for Predicting Crop Production in Support of SDG13 Climate Action". The documents also provide the step by step instructions on how to execute the project code.

## 2 System Requirement

The following section lists the minimum requirements for hardware and software. The local configuration of the development system is shown in 1.

### 2.1 Hardware Requirement

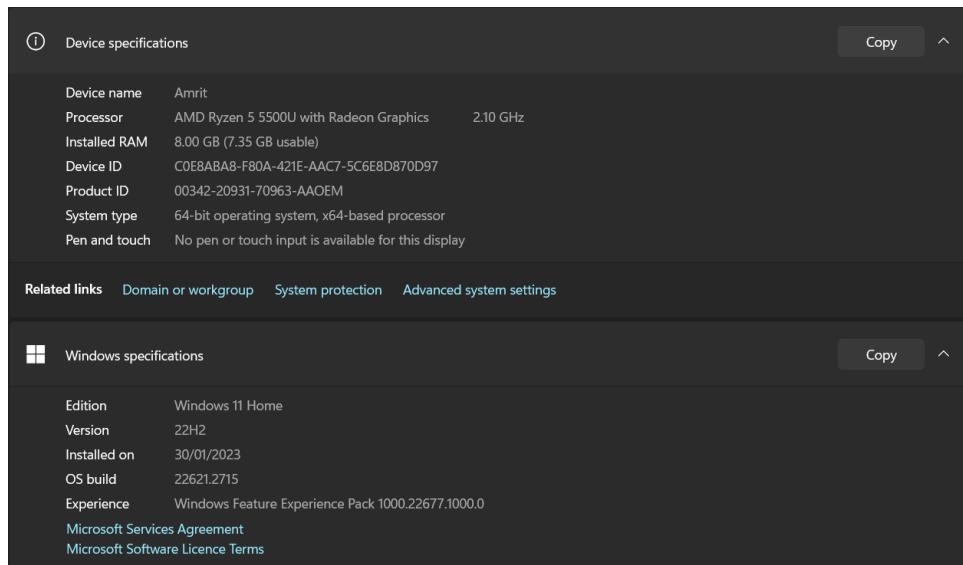


Figure 1: System Specification

- Processor: AMD Ryzen 5 5500U with Radeon Graphics 2.10 GHz or Intel
- Installed RAM: 8.00 GB
- System type: 64-bitoperating system, x64-based processor
- Storage: 500 GB HDD

## 2.2 Software Requirement

- **Microsoft Excel:** The datasets are downloaded and stored as comma separated value(csv) files.
- **Tableau:** Tableau <sup>1</sup> version 2023.1 is used for the extraction of data from pdf files (as shown in Figure 2), especially the crop production information which is downloaded from the USDA website.

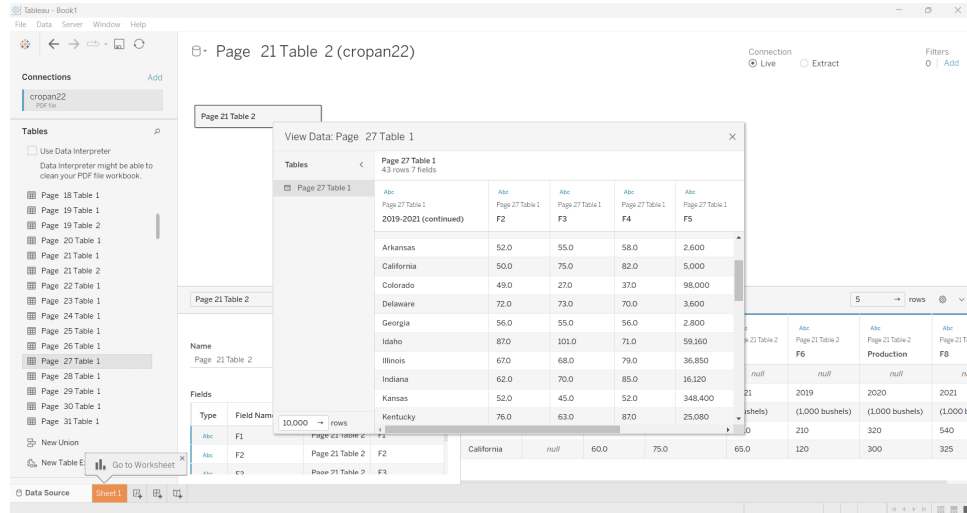


Figure 2: Tableau for extracting crop data from PDF

- **Jupyter Notebook:** Anaconda3 <sup>2</sup> Navigator which provides different software to perform project execution. Jupyter Notebook is installed through Anaconda and It is used as an IDE for the project implementation.

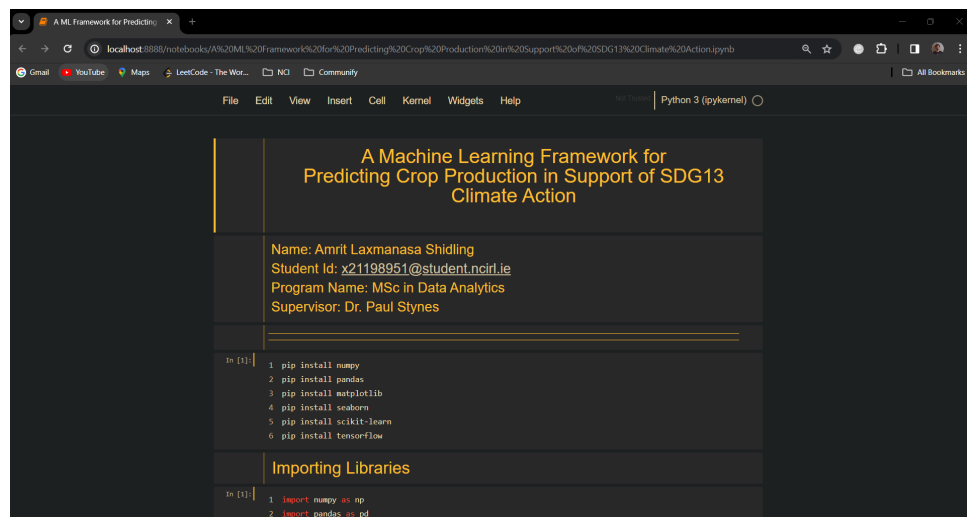
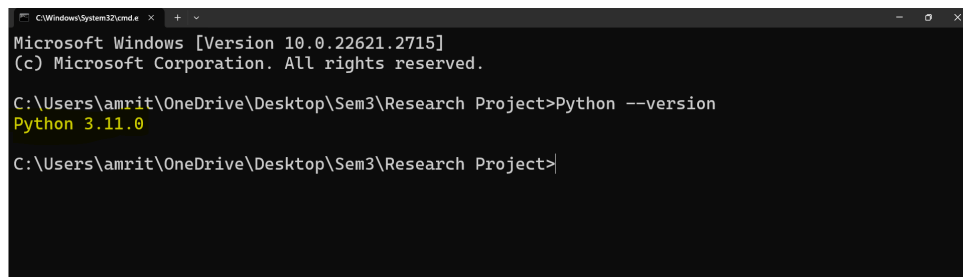


Figure 3: Jupyter notebook for project code

<sup>1</sup><https://www.tableau.com/>

<sup>2</sup><https://www.anaconda.com/>

- **Python:** Python version 3.11.0 (as shown in Figure 4 is used throughout the project including Exploratory Data Analysis(EDA), Data Preprocessing, Model Training and Evaluation.



```

C:\Windows\System32\cmd.exe
Microsoft Windows [Version 10.0.22621.2715]
(c) Microsoft Corporation. All rights reserved.

C:\Users\amrit\OneDrive\Desktop\Sem3\Research Project>Python --version
Python 3.11.0

C:\Users\amrit\OneDrive\Desktop\Sem3\Research Project>

```

Figure 4: Python Version

### 3 Package Requirement

The Python packages were installed using pip in Jupyter Notebook. The packages required are listed below:

- Pandas: Data manipulation and analysis
- NumPy: Numerical operations and handling arrays
- Matplotlib: Visualisation Library
- Seaborn: Data visualization library based on Matplotlib.
- scikit-learn: Machine Learning tasks, including data splitting, preprocessing, model evaluation, and ensemble methods.
- Tensorflow: Building and training deep learning models, particularly a Sequential model.

The packages can be installed using the following commands:

```

!pip install numpy
!pip install pandas
!pip install matplotlib
!pip install seaborn
!pip install scikit-learn
!pip install tensorflow

```

### 4 Dataset Description

There are 3 different datasets are used in the project.

1. Soil Moisture and Drought Dataset: This dataset is obtained from obtained from the NASA Langley Research Center (LaRC) POWER Project funded through the NASA Earth Science/Applied Science Program and published in Kaggle <sup>3</sup>

---

<sup>3</sup><https://www.kaggle.com/datasets/cdminix/us-drought-meteorological-data>

Variable	Description	Units
WS10M.MIN	Minimum Wind Speed at 10 Meters	m/s
QV2M	Specific Humidity at 2 Meters	g/kg
T2M_RANGE	Temperature Range at 2 Meters	°C
WS10M	Wind Speed at 10 Meters	m/s
T2M	Temperature at 2 Meters	°C
WS50M.MIN	Minimum Wind Speed at 50 Meters	m/s
T2M_MAX	Maximum Temperature at 2 Meters	°C
WS50M	Wind Speed at 50 Meters	m/s
TS	Earth Skin Temperature	°C
WS50M_RANGE	Wind Speed Range at 50 Meters	m/s
WS50M_MAX	Maximum Wind Speed at 50 Meters	m/s
WS10M_MAX	Maximum Wind Speed at 10 Meters	m/s
WS10M_RANGE	Wind Speed Range at 10 Meters	m/s
PS	Surface Pressure	kPa
T2MDEW	Dew/Frost Point at 2 Meters	°C
T2M_MIN	Minimum Temperature at 2 Meters	°C
T2MWET	Wet Bulb Temperature at 2 Meters	°C
PRECTOT	Precipitation	mm day <sup>-1</sup>
state_name	State Name	—
production	Crop Production	in 1,000 bushels

Figure 5: Data Description of Soil Moisture and Drought Dataset

2. Crop Production Dataset: This dataset is obtained from the United States Department of Agriculture website <sup>4</sup>. The description of the dataset is as shown in 1

Table 1: Dataset Description

Column	Description
state_name	State name
yield_per_acre_2019	Yield per acre in 2019
yield_per_acre_2020	Yield per acre in 2020
yield_per_acre_2021	Yield per acre in 2021
production_2019	Production in 2019
production_2020	Production in 2020
production_2021	Production in 2021

3. Fips code dataset: This dataset is obtained from the GitHub repository <sup>5</sup>. The description of this dataset is shown in 2.

Table 2: Country Fips code dataset

Column	Description
fips	FIPS code
county_name	County name
state_abbr	State abbreviation
state_name	State name
long_name	Long name

<sup>4</sup><https://www.usda.gov/>

<sup>5</sup>[https://github.com/kjhealy/fips-codes/blob/master/state\\_and\\_county\\_fips\\_master.csv](https://github.com/kjhealy/fips-codes/blob/master/state_and_county_fips_master.csv)

## 5 Folder Structure

This section explains the local folder structure for executing the project. The project is executed with 2 different .ipynb (Python Notebook) files. The first one with the name "experiment Drought Prediction Classification.ipynb" contains code for state-of-the-art experiment execution. The second file named "Prediction of crop production.ipynb" contains the code for all other experiments.

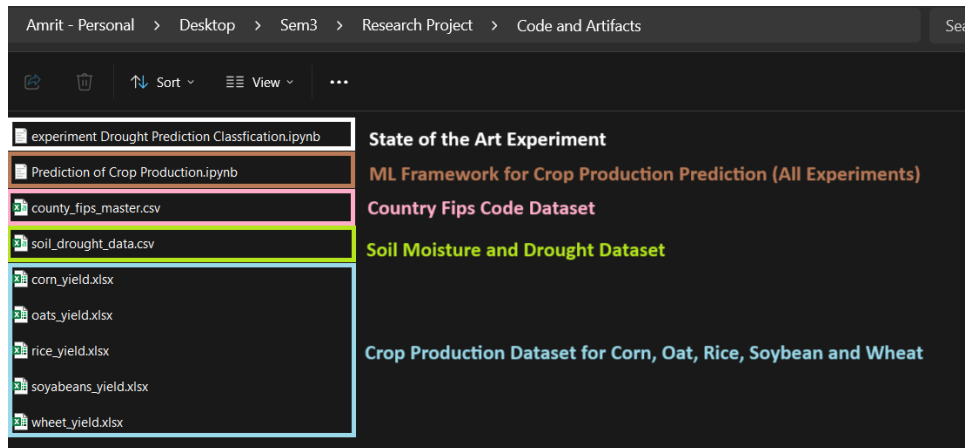


Figure 6: Filenames and Folder Structure

All the files should be in the same folder as the code reads these datafiles from the folder where code is present.

## 6 Downloading and Execution of code

As there is a restriction on the size of the artefact that can be uploaded (100MB), The dataset and code artefacts are uploaded to one drive.

Here is a link to code and data: [ClickForCodeData](#).

If the previous link doesn't work due to an encoding issue, the short URL can be used: [ClickForCodeData](#). Clicking on this link will open the folder as shown in the Figure 7

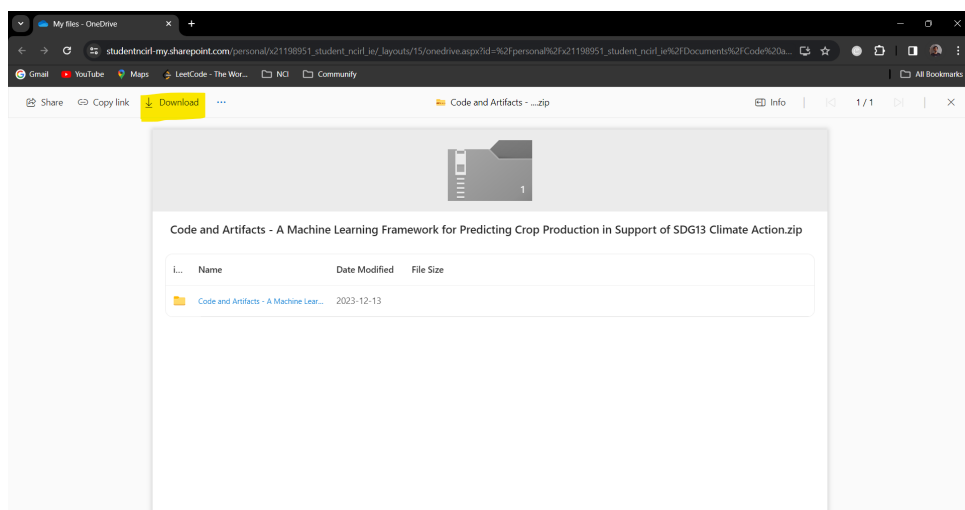


Figure 7: Zipped Folder Containing Code and Data Files



By clicking on the download as shown in Figure 7 the code and data will be downloaded. This is a zipped file, which needs to be unzipped.

To run the code the Jupyter Notebook should be started from either Anaconda or from the command prompt by navigating the folder - opening the cmd prompt - running command "Python -m notebook". It will open a new tab in the default browser as shown in Figure 8.

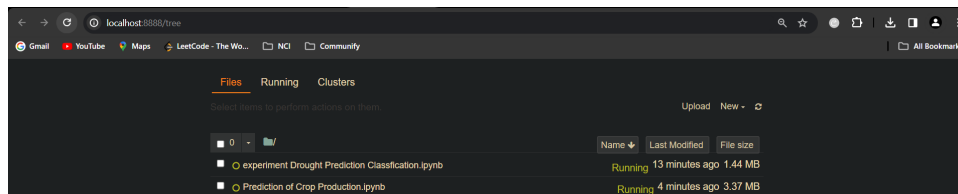


Figure 8: Running the code

The next step is to open the file and run the code by using the option Kernel -> Restart and Run All as shown in Figure 9. It will run all the code. The code covers all the sections from importing libraries to evaluating models. while running the code, some of the files are generated for storing intermediate results.

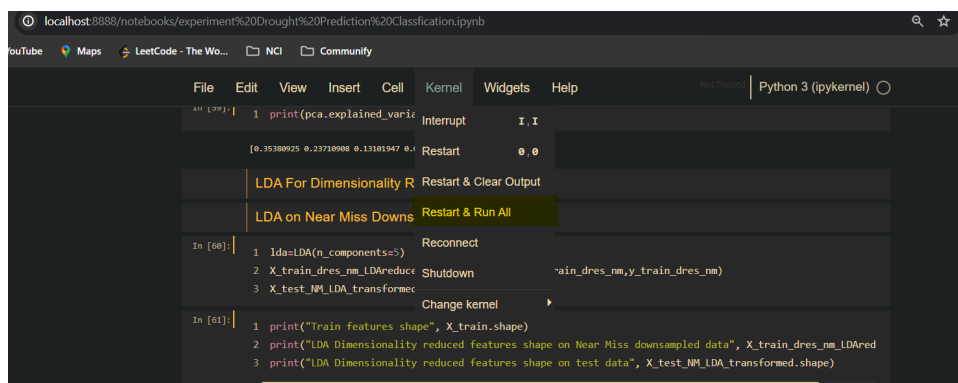


Figure 9: Running the code

The modular code is written for each of the functionality including Exploratory Data Analysis, Data Manipulation, Model Training and Evaluation. Each of the sections is highlighted with the details. For each crop type, 4 models are trained and there are 5 types of crops, a total of 20 models and each model evaluation is provided in the code.

## 7 Running Code

The below sections give an overview of some of the sample codes from the project.

- Importing Libraries: All required libraries are imported as shown in Figure 10

```
Importing Libraries

In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5
        6 from sklearn.model_selection import train_test_split, cross_val_score, KFold
        7 from sklearn.preprocessing import StandardScaler, MinMaxScaler
        8 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
        9 from sklearn.linear_model import LinearRegression
       10 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
       11 from sklearn.inspection import permutation_importance
       12
       13 import tensorflow as tf
       14 from tensorflow.keras.models import Sequential
       15 from tensorflow.keras.layers import LSTM, Dense
```

Figure 10: Importing Libraries

- Combining Dataset: The datasets are combined based on the state FIPs code and Date as shown in Figure 11

```
Function for Combining Soil Moisture and Crop data

In [20]: 1 def prepare_data_for_model(drought_soil_df, crop_production_file_path):
        2     yield_df = pd.read_csv(crop_production_file_path)
        3
        4     melted_df = pd.melt(yield_df, id_vars=['state_name'], var_name='variable', value_name='value')
        5     melted_df['year'] = melted_df['variable'].str.extract(r'(\d+)')
        6     melted_df['type'] = melted_df['variable'].str.extract(r'(\.+)\d+')
        7     result_df = melted_df.pivot_table(index=['state_name', 'year'], columns=['type'], values='value', a
        8     result_df.columns = ['state_name', 'year', 'yield_per_acre', 'production']
        9
       10     yield_df = result_df
       11     soil_drought_yield_df = drought_soil_df
       12     soil_drought_yield_df['date'] = pd.to_datetime(soil_drought_yield_df['date'])
       13     soil_drought_yield_df['year'] = soil_drought_yield_df['date'].dt.year
       14     yield_df['year'] = yield_df['year'].astype(int)
       15     soil_drought_yield_df = pd.merge(soil_drought_yield_df, yield_df, on=['state_name', 'year'], how='l
       16
       17     soil_drought_yield_df.head()
       18     soil_drought_yield_df[soil_drought_yield_df['production'].notnull()].head()
       19     soil_drought_yield_df['production'] = pd.to_numeric(soil_drought_yield_df['production'].replace('(N
       20     soil_drought_yield_df = soil_drought_yield_df[soil_drought_yield_df['production'].notnull()]
       21
       22     soil_drought_yield_df = soil_drought_yield_df[soil_drought_yield_df['score'].notnull()]
       23
       24     return soil_drought_yield_df
```

Figure 11: Combining Dataset

- Exploratory Data Analysis: The EDA is conducted on the data using modular function as shown in Figure 12

```
Exploratory Data Analysis

In [17]: 1 def display_basic_info(df):
2         print("Dataset Information:")
3         print(df.info())
4
5         def display_summary_statistics(df):
6             print("Summary Statistics:")
7             df.describe()
8
9         def visualize_missing_data(df):
10             plt.figure(figsize=(10, 6))
11             sns.heatmap(df.isnull(), cbar=False, cmap='viridis')
12             plt.title("Missing Data Visualization")
13             plt.show()
14
15         def visualize_correlation(df):
16             correlation_matrix = df.corr()
17             plt.figure(figsize=(12, 8))
18             sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f", linewidths=0.5, xticklabels
19                         = correlation_matrix.columns)
20             plt.title("Correlation Heatmap")
21             plt.xticks(rotation=90)
22             plt.yticks(rotation=0)
23             plt.show()
24
25         def explore_categorical_variables(df, categorical_columns, crop):
26             for column in categorical_columns:
27                 plt.figure(figsize=(8, 6))
28                 sns.boxplot(x=column, y='production', data=df)
29                 plt.title(f"Median {crop} Production across {column}")
30                 plt.xticks(rotation=90)
31                 plt.show()
```

Figure 12: Modular Code for EDA

- Correlation Analysis: The Correlation Analysis is conducted on the data as shown in Figure 13

```

Correlation Analysis

In [18]: 1 def correlation_analysis(data, strong_threshold=0.8, weak_threshold=0.2):
2         numeric_columns = data.select_dtypes(include=['float64', 'int64'])
3         correlation_matrix = numeric_columns.corr()
4
5         plt.figure(figsize=(12, 8))
6         sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f", linewidths=0.5)
7         plt.title("Correlation Heatmap")
8         plt.xticks(rotation=90)
9         plt.yticks(rotation=0)
10        plt.show()
11
12        upper = correlation_matrix.where(
13            np.triu(np.ones(correlation_matrix.shape), k=1).astype(np.bool)
14        )
15        to_drop_strong = [
16            column for column in upper.columns if
17            any((upper[column].abs() > strong_threshold) | (upper[column] < -strong_threshold))
18        ]
19
20        # Handle weak correlations
21        to_drop_weak = [
22            column for column in correlation_matrix.columns if
23            all((correlation_matrix[column].abs() < weak_threshold) & (correlation_matrix[column] > -weak_t
24        ]
25
26        to_drop = list(set(to_drop_strong + to_drop_weak))
27        print("Columns to drop:", to_drop)
28
29        return data.drop(to_drop, axis=1)

```

Figure 13: Correlation Analysis

- Applying Model: Function is written to accept the data and name of the model to be applied to perform training as shown in Figure 14

```

Function for Model Preparation and Training

In [21]: 1 def apply_model(soil_drought_yield_df, model_name):
2         features = soil_drought_yield_df.drop(columns=['production', 'date'])
3         target = soil_drought_yield_df['production']
4
5         X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2, random_state=42)
6         X_train = pd.get_dummies(X_train, columns=['state_name'])
7         X_test = pd.get_dummies(X_test, columns=['state_name'])
8
9         scaler = StandardScaler()
10        X_train_scaled = scaler.fit_transform(X_train)
11        X_test_scaled = scaler.transform(X_test)
12        if model_name == 'NeuralNetwork':
13            y_train = y_train.values.astype('float32')
14            model = tf.keras.Sequential([
15                tf.keras.layers.Dense(128, activation='relu', input_shape=(X_train_scaled.shape[1],)),
16                tf.keras.layers.Dense(64, activation='relu'),
17                tf.keras.layers.Dense(1) # Output Layer for regression
18            ])
19            model.compile(optimizer='adam', loss='mean_squared_error')
20            model.fit(X_train_scaled, y_train, epochs=50, batch_size=32, validation_split=0.2)
21            predictions = model.predict(X_test_scaled).flatten()
22            return evaluate_model(predictions, y_test)
23        else:
24            model = models.get(model_name)
25            model.fit(X_train_scaled, y_train)
26            predictions = model.predict(X_test_scaled)
27            return evaluate_model(predictions, y_test)

```

Figure 14: Model Training Method

- Model Evaluation: Methods are written to perform evaluation on of the model performance and plotting the visuals as shown in Figure 15

```

Functin for Model Evaluation

In [22]: 1 def evaluate_model(predictions, y_test):
2         mse = mean_squared_error(y_test, predictions)
3         mae = mean_absolute_error(y_test, predictions)
4         r2 = r2_score(y_test, predictions)
5         return {'mse': mse, 'mae': mae, 'r2': r2}

Displaying Result

In [23]: 1 def display_result_in_tabular(model_results, crop_name):
2         df_results = pd.DataFrame(model_results).T
3         plt.figure(figsize=(10, 6))
4         sns.heatmap(df_results, annot=True, cmap="YlGnBu", fmt='.4f')
5         plt.title(f'Model Comparison for {crop_name}')
6         plt.show()

Plotting the Model Performance

In [24]: 1 def plot_result(model_results, crop_name):
2         models = list(model_results.keys())
3         metrics = list(model_results[models[0]].keys())
4         for metric in metrics:
5             values = [model_results[model][metric] for model in models]
6             plt.bar(models, values, label=metric)
7
8         plt.xlabel('Models')
9         plt.ylabel('Metric Values')
10        plt.title(f'Model Comparison for {crop_name}')
11        plt.legend()
12        plt.show()

```

Figure 15: Model Evaluation

- Cross Validation: Methods are written to perform cross-validation on of the model performance as shown in Figure 16

```

In [19]: 1 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
2 def evaluate_model(predictions, true_values):
3     mse = mean_squared_error(true_values, predictions)
4     mae = mean_absolute_error(true_values, predictions)
5     r2 = r2_score(true_values, predictions)
6
7     print("Mean Squared Error:", mse)
8     print("Mean Absolute Error:", mae)
9     print("R-squared:", r2)
10
11 def cross_validation(model, X, y, n_splits=5):
12     kf = KFold(n_splits=n_splits, shuffle=True, random_state=42)
13     scores = cross_val_score(model, X, y, cv=kf, scoring='neg_mean_squared_error')
14     rmse_scores = np.sqrt(-scores)
15     return rmse_scores.mean()
16
17 def assess_feature_importance(model, X, y):
18     result = permutation_importance(model, X, y, n_repeats=10, random_state=42, n_jobs=-1)
19     feature_importance = result.importances_mean
20     return feature_importance
21
22 def evaluate_model_robustness(model, X, y):
23     perturbed_X = X + np.random.normal(scale=0.01, size=X.shape)
24     perturbed_predictions = model.predict(perturbed_X)
25     perturbed_rmse = np.sqrt(mean_squared_error(y, perturbed_predictions))
26     return perturbed_rmse
27
28 def apply_model_with_evaluation(soil_drought_yield_df, model_name):
29     features = soil_drought_yield_df.drop(columns=['production', 'date'])
30     target = soil_drought_yield_df['production']

```

Figure 16: Cross Validation Methods

- Calling Model Training Methods: To perform model training the modular method `apply_model` is called as shown in Figure 17

```

Training and Evaluating Linear Regression Model

In [65]: 1 corn_evaluation_result_lr = apply_model(corn_soil_drought_df, 'LinearRegression')
2 corn_evaluation_result_lr
3

{'mse': 2060104238.637036, 'mae': 27131.571965290033, 'r2': 0.9875402913965784}

Training and Evaluating Gradient Boost Regression Model

In [66]: 1 corn_evaluation_results_gb = apply_model(corn_soil_drought_df, 'GradientBoostingRegressor')
2 corn_evaluation_results_gb
3

{'mse': 1045506433.8708422, 'mae': 17504.58938481542, 'r2': 0.9936766765172759}

Training and Evaluating Random Forest Regression Model

In [67]: 1 corn_evaluation_result_rf = apply_model(corn_soil_drought_df, 'RandomForestRegressor')
2 corn_evaluation_result_rf
3

{'mse': 24.7918894164158, 'mae': 3.9740302462044492, 'r2': 0.9999999998500563}

Training and Evaluating FeedForward Nueral Network

In [70]: 1 corn_evaluation_result_nn = apply_model(corn_soil_drought_df, 'NueralNetwork')
2 corn_evaluation_result_nn

```

Figure 17: Calling Methods