

Dataset-1

January 25, 2025

1 Predicting crop yield from environmental factors

1.1 Research question

“Can crop yields be predicted with environmental factors like rainfall, soil type or fertilization”

1.2 Dataset

<https://www.kaggle.com/datasets/samuelotiattakorah/agriculture-crop-yield?resource=download>

1.3 Visualizations

In this section the dataset is loaded into a dataframe to visualize the data and possible correlations between attributes

```
[1]: # Imports
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load data
df = pd.read_csv("./crop_yield.csv")

# Display dataset info
df.info()

# Display the first few rows
df.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000000 entries, 0 to 999999
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Region                 1000000 non-null object
1   Soil_Type              1000000 non-null object
2   Crop                  1000000 non-null object
3   Rainfall_mm           1000000 non-null float64
4   Temperature_Celsius   1000000 non-null float64
```

```

5   Fertilizer_Used      1000000 non-null   bool
6   Irrigation_Used      1000000 non-null   bool
7   Weather_Condition    1000000 non-null   object
8   Days_to_Harvest      1000000 non-null   int64
9   Yield_tons_per_hectare 1000000 non-null   float64
dtypes: bool(2), float64(3), int64(1), object(4)
memory usage: 62.9+ MB

```

```

[1]:   Region Soil_Type   Crop Rainfall_mm  Temperature_Celsius \
0   West      Sandy   Cotton    897.077239             27.676966
1   South     Clay    Rice     992.673282             18.026142
2   North     Loam   Barley    147.998025             29.794042
3   North     Sandy  Soybean    986.866331             16.644190
4   South     Silt   Wheat     730.379174             31.620687

      Fertilizer_Used  Irrigation_Used Weather_Condition  Days_to_Harvest \
0              False              True          Cloudy             122
1              True              True          Rainy             140
2              False             False          Sunny             106
3              False              True          Rainy             146
4              True              True          Cloudy             110

      Yield_tons_per_hectare
0              6.555816
1              8.527341
2              1.127443
3              6.517573
4              7.248251

```

```

[2]: # Check for missing values
missing_values = df.isnull().sum()
print(missing_values)

```

```

Region          0
Soil_Type       0
Crop            0
Rainfall_mm     0
Temperature_Celsius 0
Fertilizer_Used 0
Irrigation_Used 0
Weather_Condition 0
Days_to_Harvest 0
Yield_tons_per_hectare 0
dtype: int64

```

```

[3]: # Encode categorical variables with one-hot encoding
encoded_df = pd.get_dummies(df, columns=['Region', 'Soil_Type',
↪ 'Weather_Condition', 'Crop'], drop_first=True)

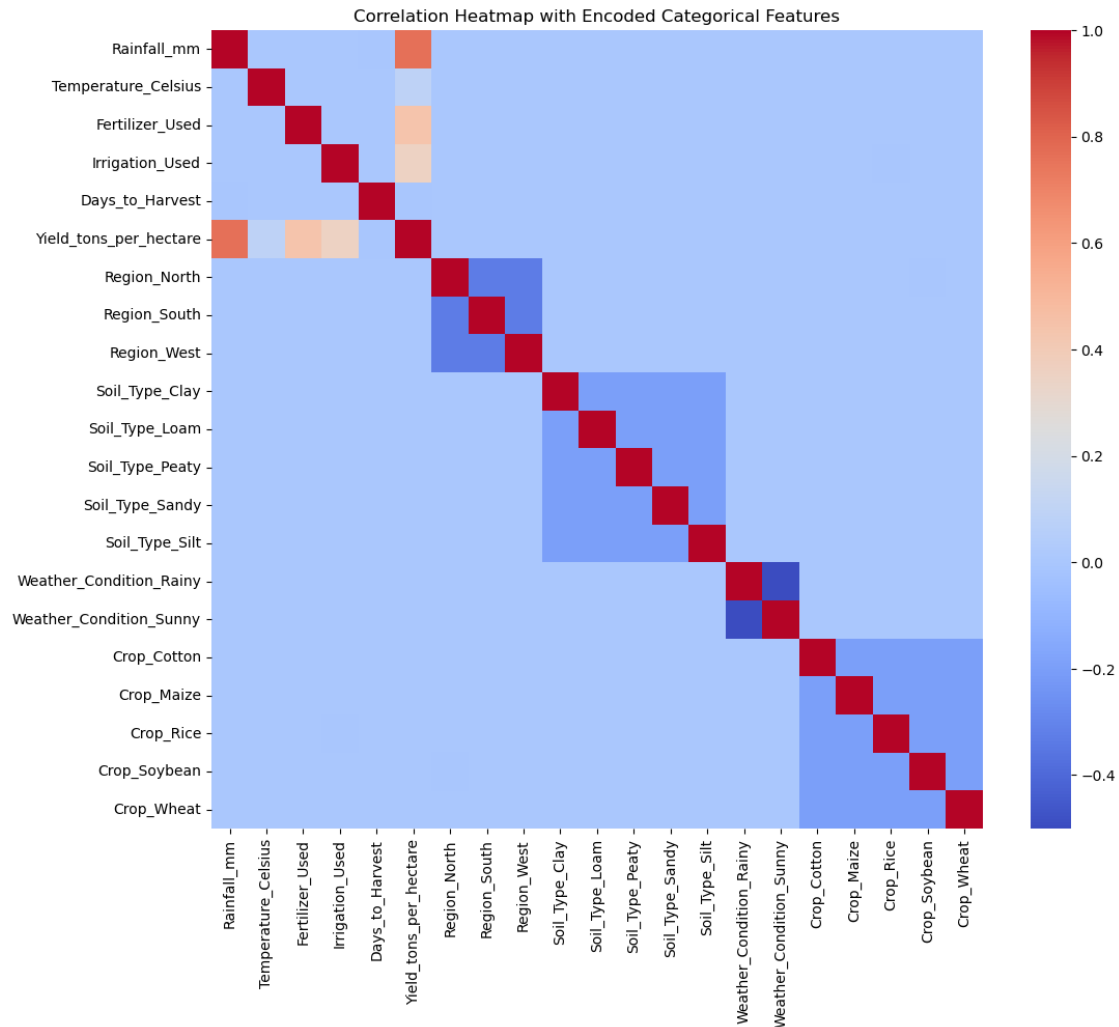
```

```

# Recalculate the correlation matrix with encoded features
correlation_matrix_encoded = encoded_df.corr()

# Plot the heatmap for the encoded dataset
plt.figure(figsize=(12, 10))
sns.heatmap(correlation_matrix_encoded, annot=False, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Heatmap with Encoded Categorical Features')
plt.show()

```



```

[4]: # Scatter plot: Rainfall vs. Yield
plt.figure(figsize=(8, 6))
sns.scatterplot(data=df, x='Rainfall_mm', y='Yield_tons_per_hectare',
                hue='Soil_Type', alpha=0.6)
plt.title('Rainfall vs. Crop Yield')

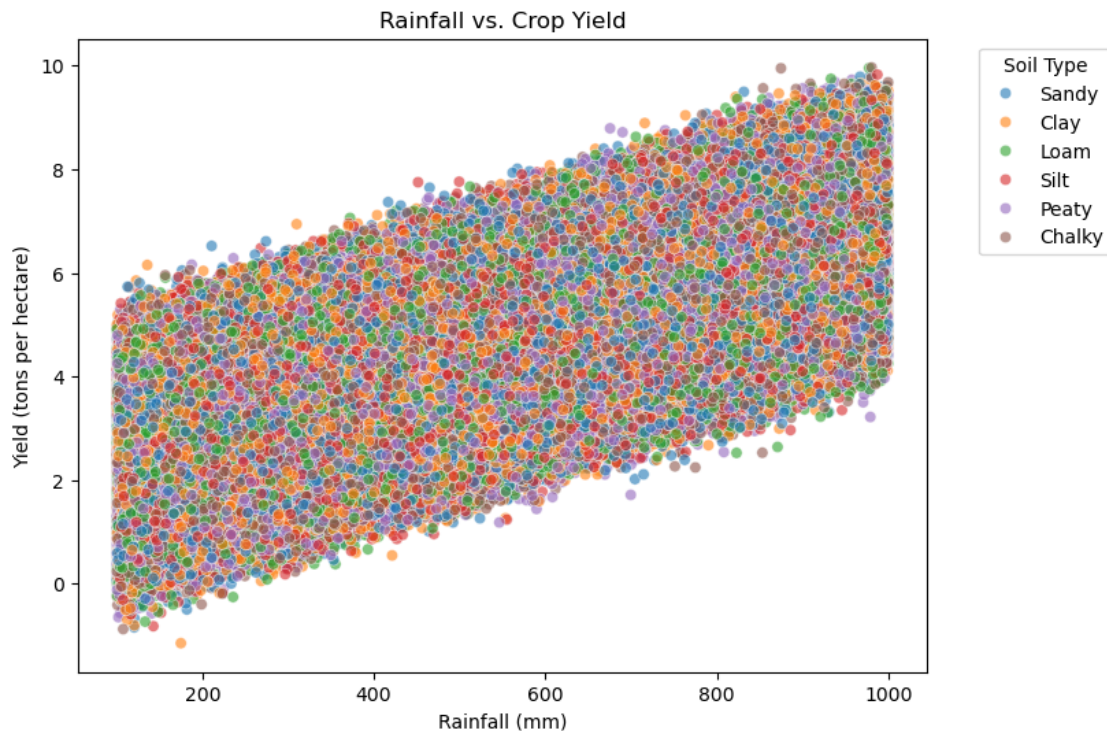
```

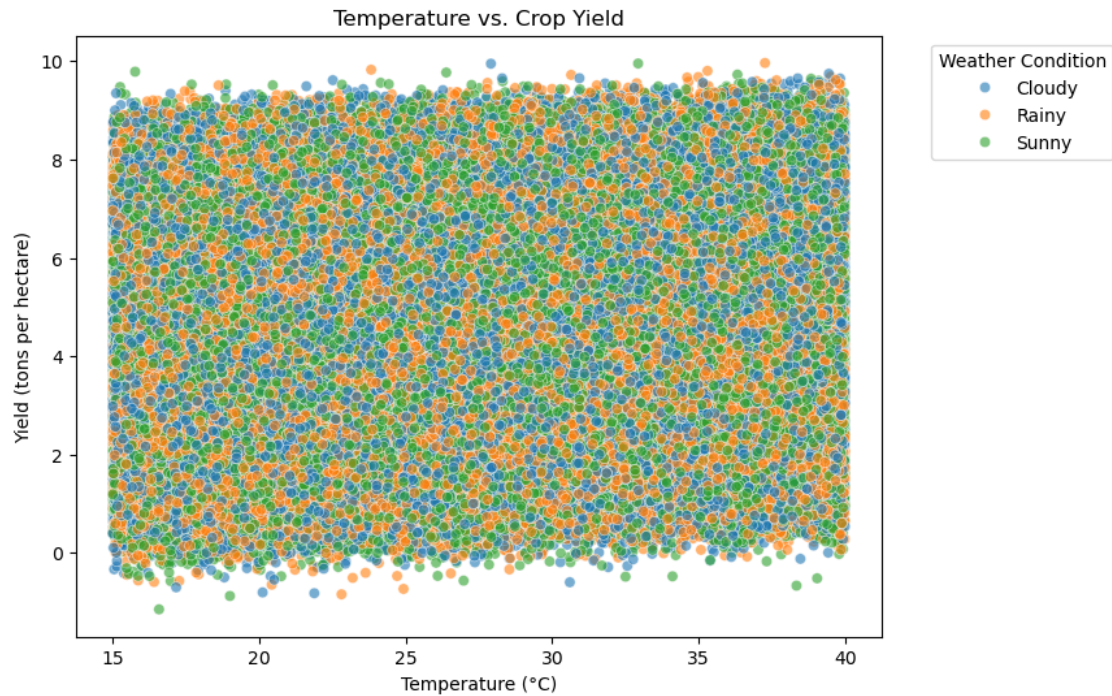
```

plt.xlabel('Rainfall (mm)')
plt.ylabel('Yield (tons per hectare)')
plt.legend(title='Soil Type', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.show()

# Scatter plot: Temperature vs. Yield
plt.figure(figsize=(8, 6))
sns.scatterplot(data=df, x='Temperature_Celsius', y='Yield_tons_per_hectare',
               hue='Weather_Condition', alpha=0.6)
plt.title('Temperature vs. Crop Yield')
plt.xlabel('Temperature (°C)')
plt.ylabel('Yield (tons per hectare)')
plt.legend(title='Weather Condition', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.show()

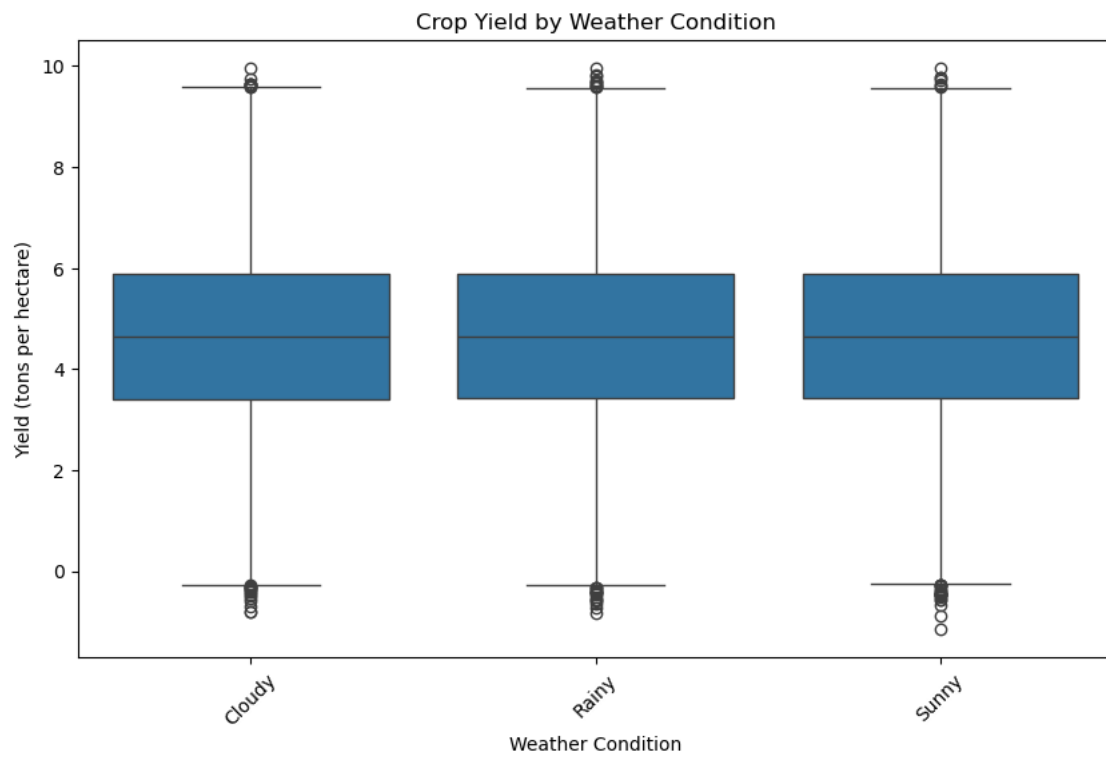
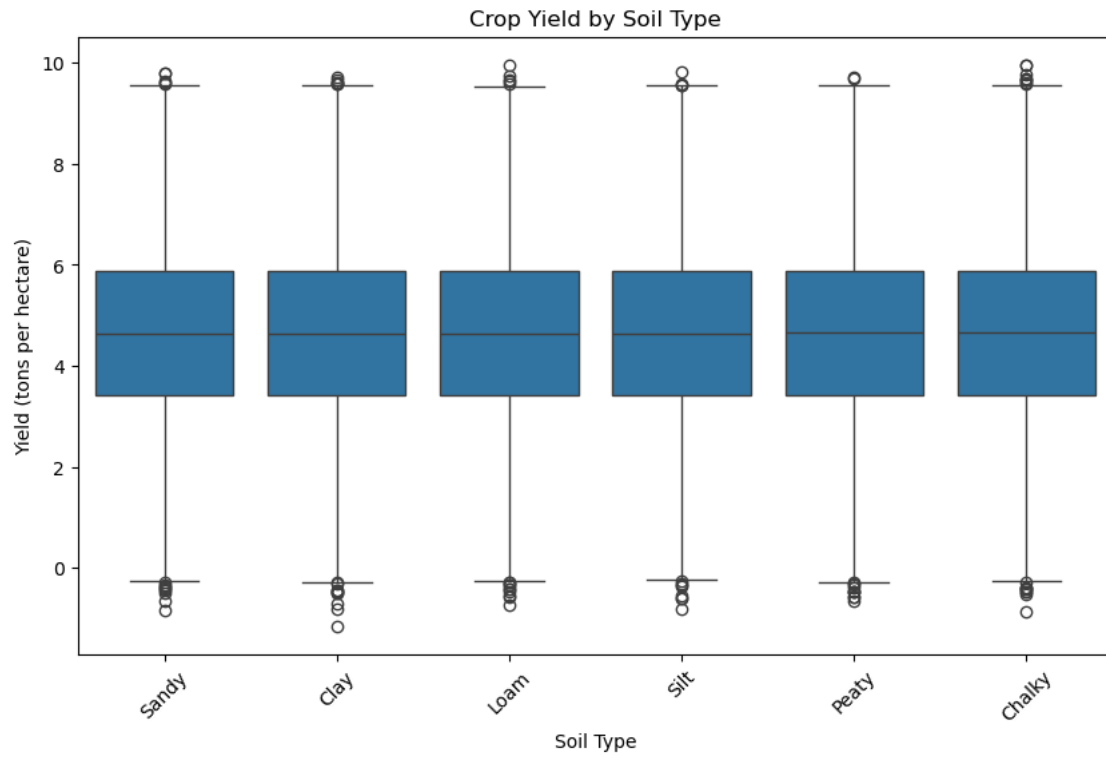
```



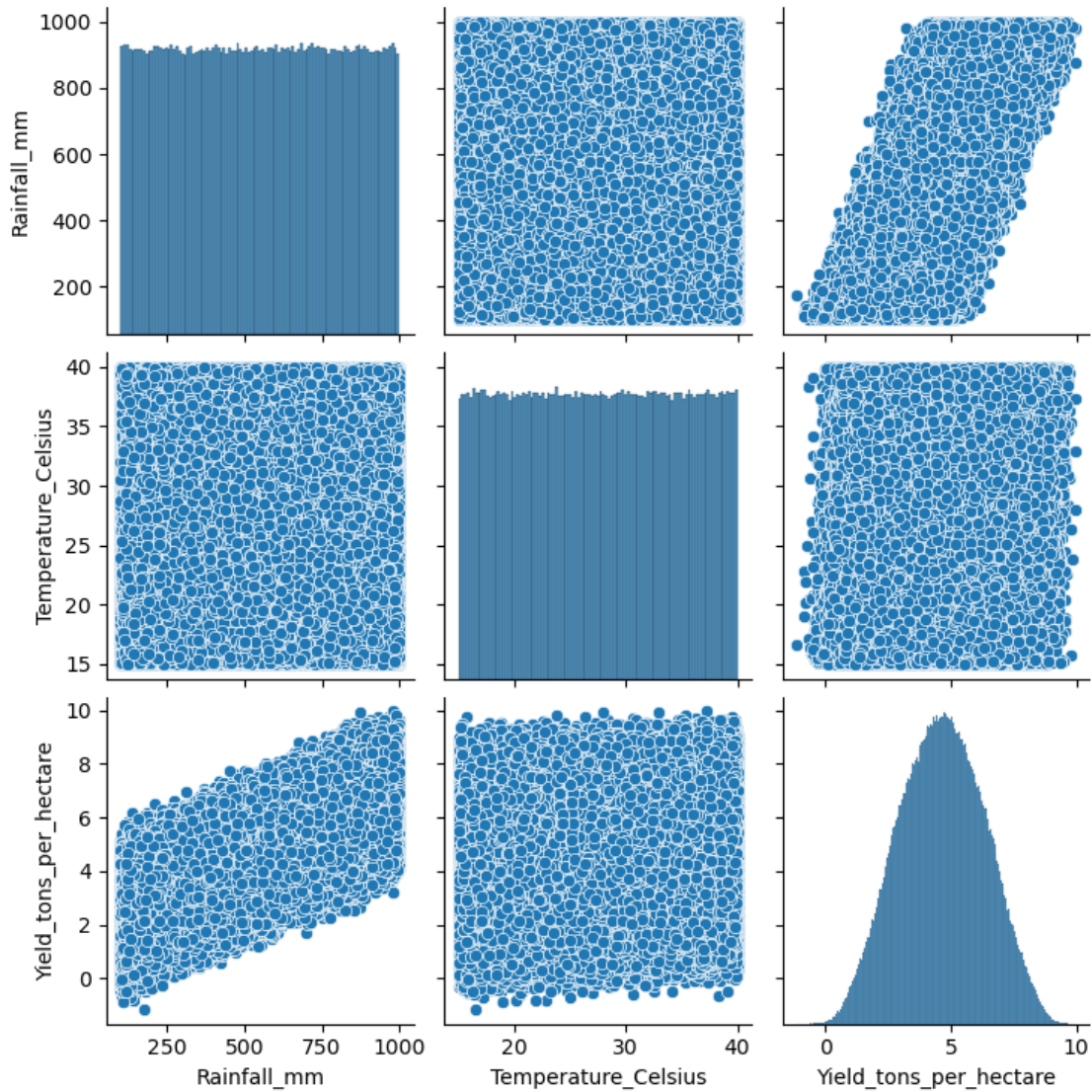


```
[5]: # Boxplot: Soil Type vs. Yield
plt.figure(figsize=(10, 6))
sns.boxplot(data=df, x='Soil_Type', y='Yield_tons_per_hectare')
plt.title('Crop Yield by Soil Type')
plt.xlabel('Soil Type')
plt.ylabel('Yield (tons per hectare)')
plt.xticks(rotation=45)
plt.show()

# Boxplot: Weather Condition vs. Yield
plt.figure(figsize=(10, 6))
sns.boxplot(data=df, x='Weather_Condition', y='Yield_tons_per_hectare')
plt.title('Crop Yield by Weather Condition')
plt.xlabel('Weather Condition')
plt.ylabel('Yield (tons per hectare)')
plt.xticks(rotation=45)
plt.show()
```



```
[6]: # Pairplot for selected features
selected_features = ['Rainfall_mm', 'Temperature_Celsius', 'Yield_tons_per_hectare']
sns.pairplot(df[selected_features])
plt.show()
```



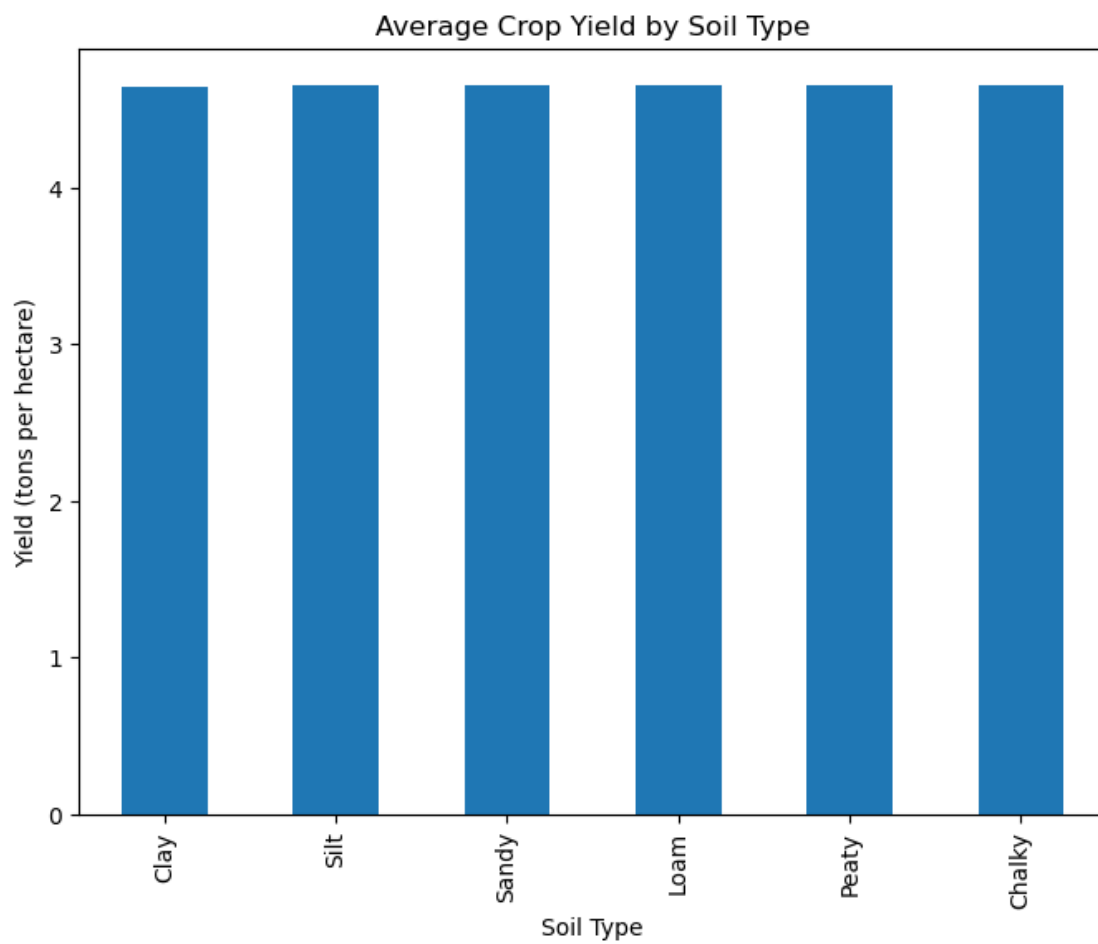
```
[7]: # Average yield by soil type
soil_yield = df.groupby('Soil_Type')['Yield_tons_per_hectare'].mean().
    sort_values()
soil_yield.plot(kind='bar', figsize=(8, 6))
plt.title('Average Crop Yield by Soil Type')
```

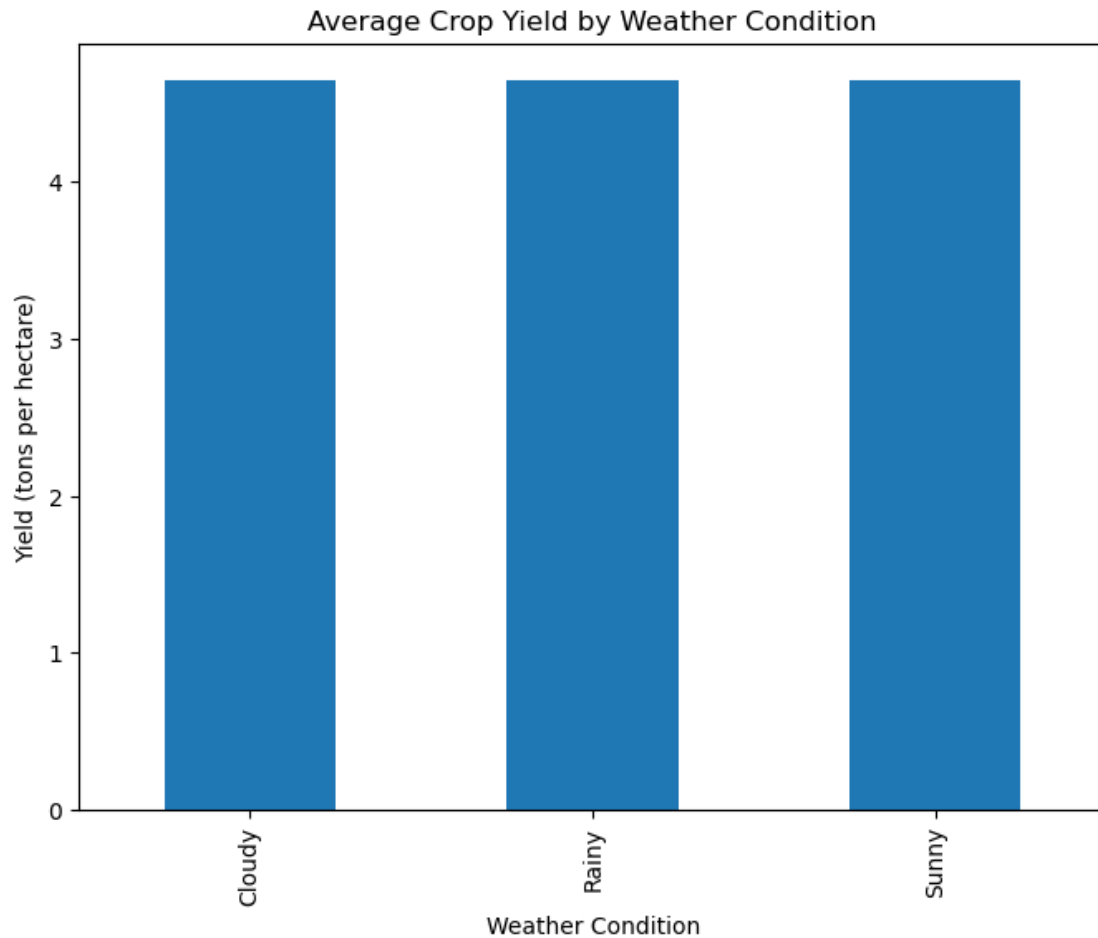
```

plt.ylabel('Yield (tons per hectare)')
plt.xlabel('Soil Type')
plt.show()

# Average yield by weather condition
weather_yield = df.groupby('Weather_Condition')['Yield_tons_per_hectare'].mean().
    ↪sort_values()
weather_yield.plot(kind='bar', figsize=(8, 6))
plt.title('Average Crop Yield by Weather Condition')
plt.ylabel('Yield (tons per hectare)')
plt.xlabel('Weather Condition')
plt.show()

```

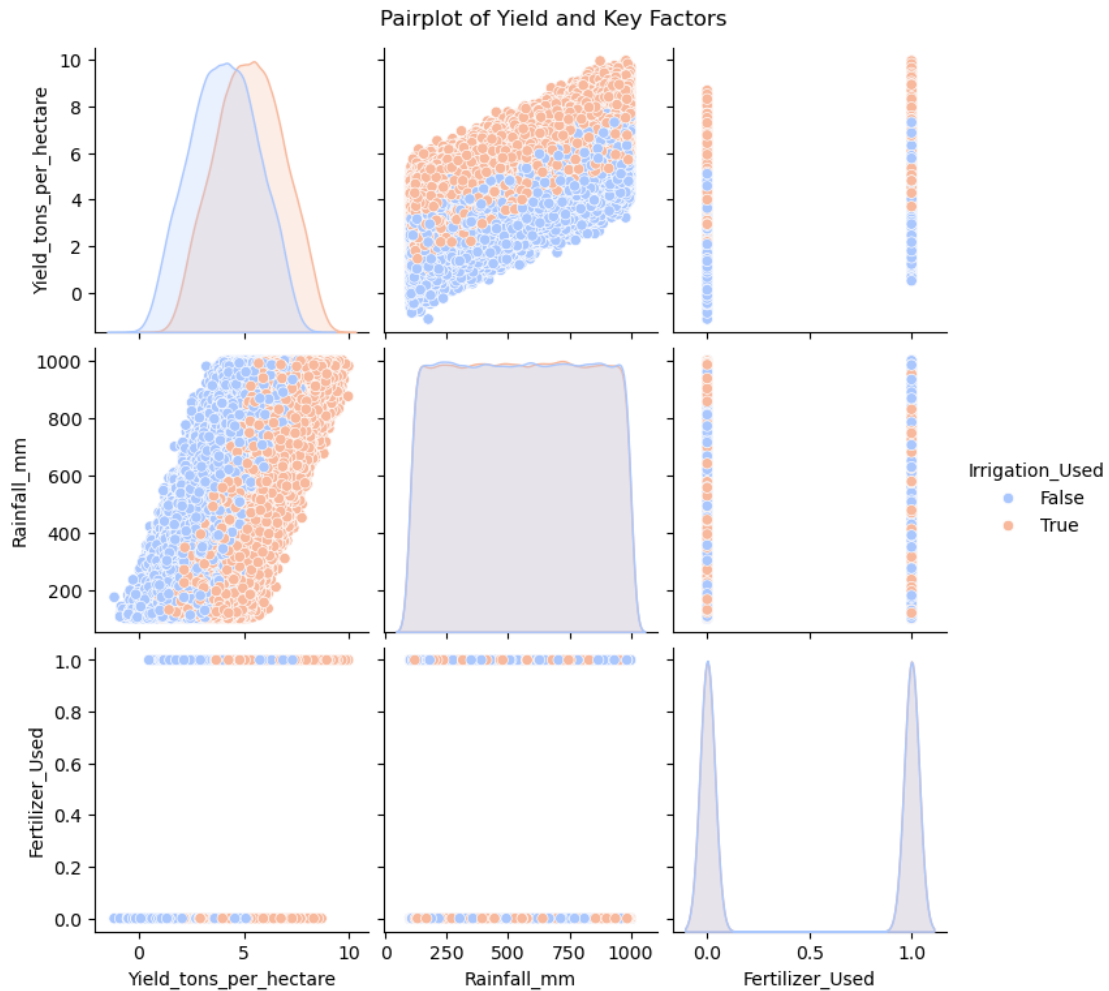




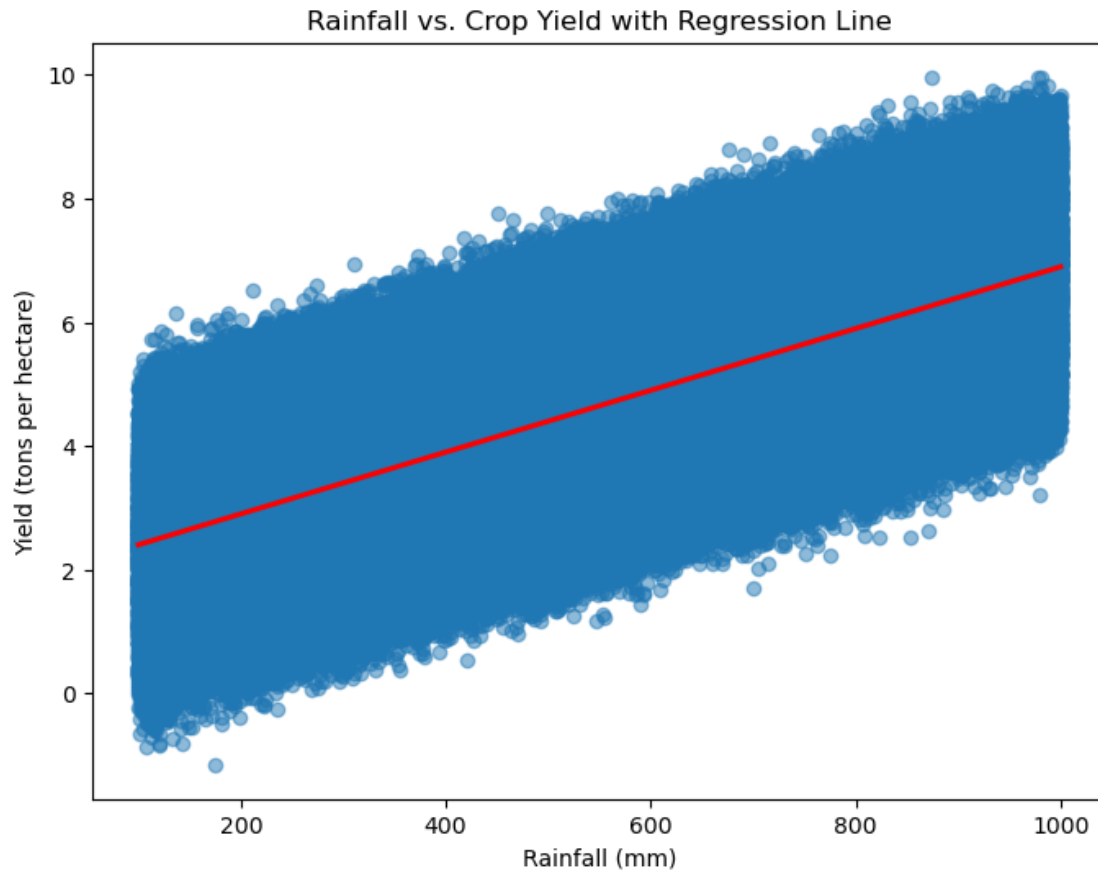
1.4 Deeper Visualization on key factors

```
[8]: # Pairplot for the key features
import seaborn as sns
import matplotlib.pyplot as plt

# Select relevant columns
selected_features = ['Yield_tons_per_hectare', 'Rainfall_mm', 'Fertilizer_Used',
                    ↪ 'Irrigation_Used']
sns.pairplot(df[selected_features], diag_kind='kde', hue='Irrigation_Used',
            ↪ palette='coolwarm')
plt.suptitle('Pairplot of Yield and Key Factors', y=1.02)
plt.show()
```



```
[9]: # Scatter plot with regression line for Rainfall vs. Yield
plt.figure(figsize=(8, 6))
sns.regplot(data=df, x='Rainfall_mm', y='Yield_tons_per_hectare',
            scatter_kws={'alpha':0.5}, line_kws={'color':'red'})
plt.title('Rainfall vs. Crop Yield with Regression Line')
plt.xlabel('Rainfall (mm)')
plt.ylabel('Yield (tons per hectare)')
plt.show()
```



```
[11]: # Boxplot for Fertilizer Used
plt.figure(figsize=(8, 6))
sns.boxplot(data=df, x='Fertilizer_Used', y='Yield_tons_per_hectare',
            palette='Set2')
plt.title('Effect of Fertilizer Usage on Yield')
plt.xlabel('Fertilizer Used')
plt.ylabel('Yield (tons per hectare)')
plt.xticks([0, 1], ['No', 'Yes'])
plt.show()

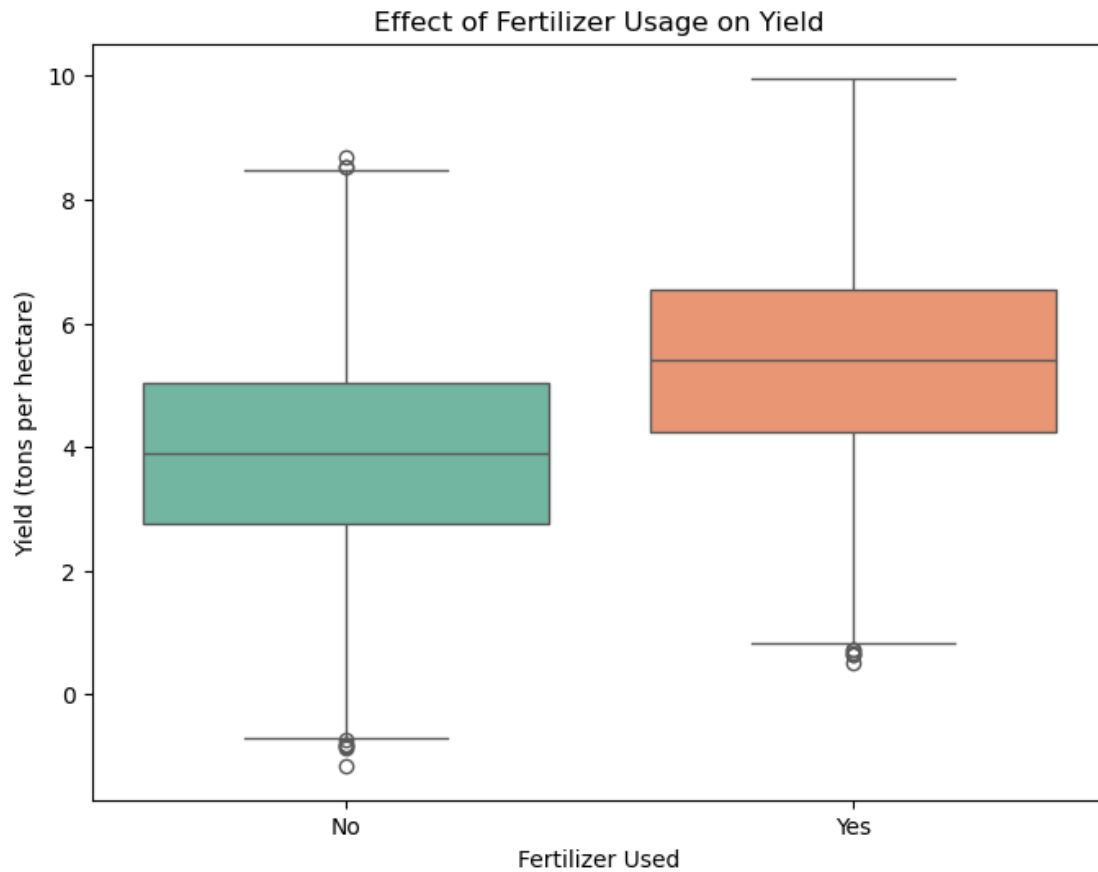
# Boxplot for Irrigation Used
plt.figure(figsize=(8, 6))
sns.boxplot(data=df, x='Irrigation_Used', y='Yield_tons_per_hectare',
            palette='Set3')
plt.title('Effect of Irrigation Usage on Yield')
plt.xlabel('Irrigation Used')
plt.ylabel('Yield (tons per hectare)')
plt.xticks([0, 1], ['No', 'Yes'])
plt.show()
```

```
C:\Users\Gregor\AppData\Local\Temp\ipykernel_11332\1144540635.py:3:
```

```
FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be removed in  
v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same  
effect.
```

```
sns.boxplot(data=df, x='Fertilizer_Used', y='Yield_tons_per_hectare',  
palette='Set2')
```

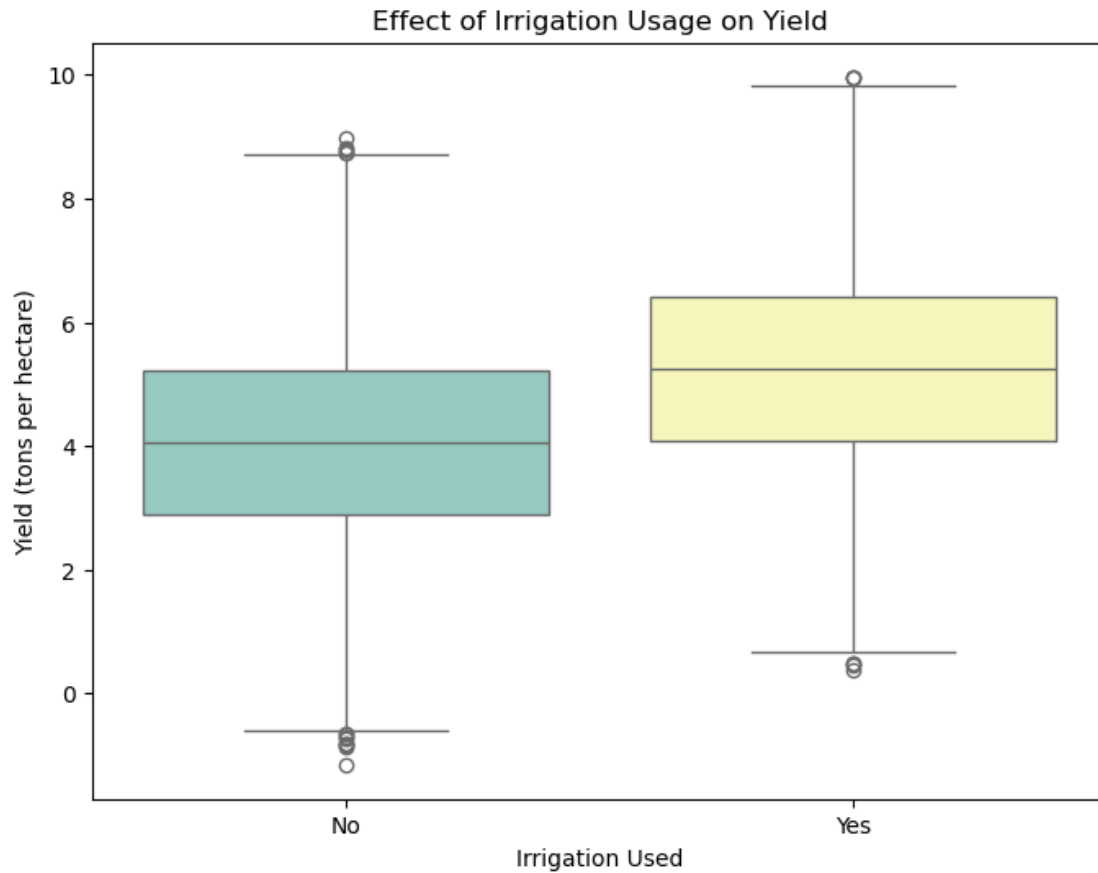


```
C:\Users\Gregor\AppData\Local\Temp\ipykernel_11332\1144540635.py:12:
```

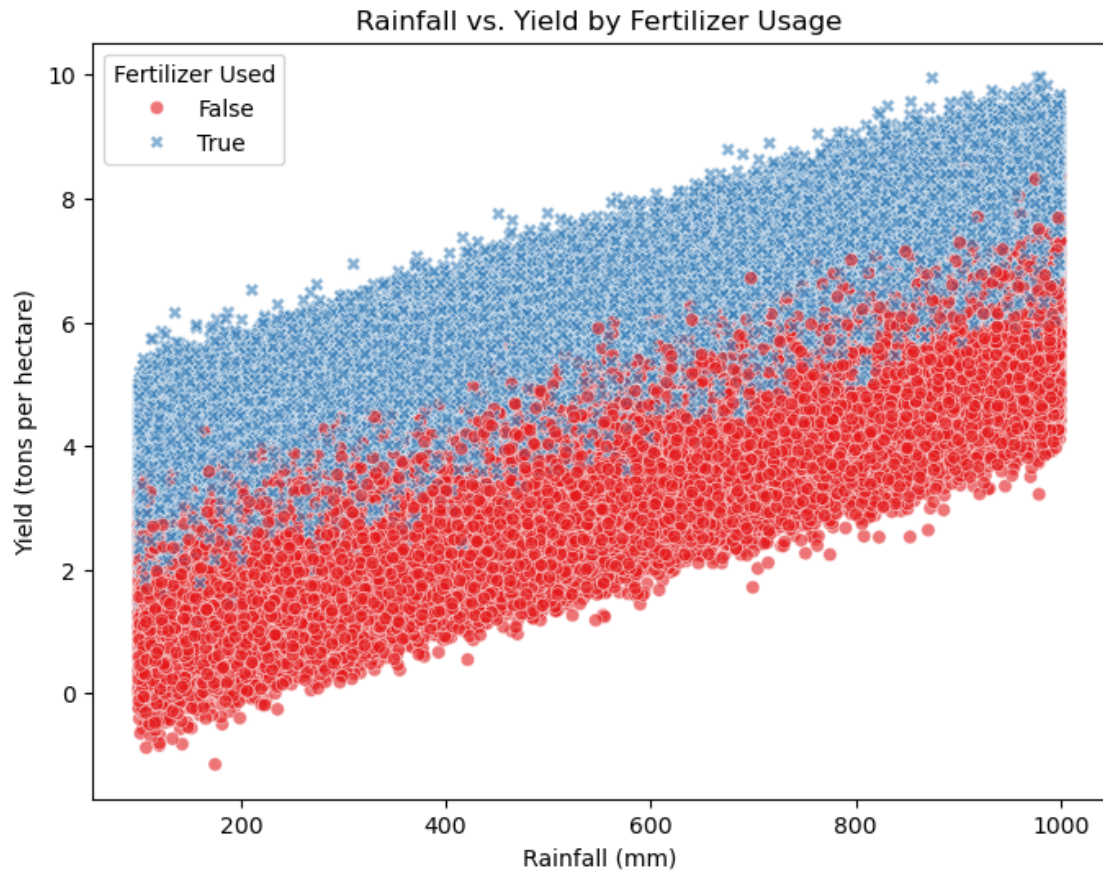
```
FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be removed in  
v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same  
effect.
```

```
sns.boxplot(data=df, x='Irrigation_Used', y='Yield_tons_per_hectare',  
palette='Set3')
```



```
[12]: # Interaction plot: Rainfall vs Yield by Fertilizer Used
plt.figure(figsize=(8, 6))
sns.scatterplot(data=df, x='Rainfall_mm', y='Yield_tons_per_hectare',
               hue='Fertilizer_Used', style='Fertilizer_Used', alpha=0.6, palette='Set1')
plt.title('Rainfall vs. Yield by Fertilizer Usage')
plt.xlabel('Rainfall (mm)')
plt.ylabel('Yield (tons per hectare)')
plt.legend(title='Fertilizer Used', loc='best')
plt.show()
```



```
[13]: # Correlation values for selected features
correlation_values = df[['Yield_tons_per_hectare', 'Rainfall_mm',
    ↳ 'Fertilizer_Used', 'Irrigation_Used']].corr()
print("Correlation Matrix:\n", correlation_values)
```

Correlation Matrix:

	Yield_tons_per_hectare	Rainfall_mm	Fertilizer_Used	\
Yield_tons_per_hectare	1.000000	0.764618	0.442099	
Rainfall_mm	0.764618	1.000000	-0.001076	
Fertilizer_Used	0.442099	-0.001076	1.000000	
Irrigation_Used	0.353741	-0.000568	0.001510	

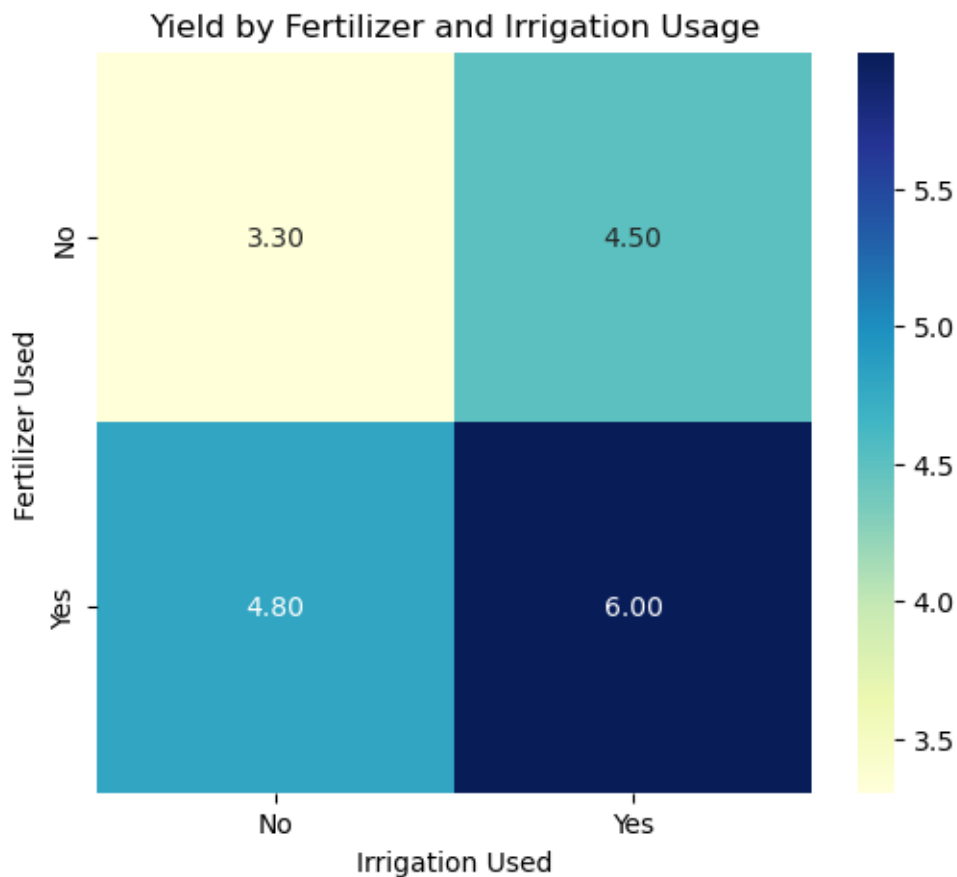
	Irrigation_Used
Yield_tons_per_hectare	0.353741
Rainfall_mm	-0.000568
Fertilizer_Used	0.001510
Irrigation_Used	1.000000

```
[14]: # Create a crosstab to explore interactions between Irrigation and Fertilizer Usage
interaction_crosstab = pd.crosstab(df['Fertilizer_Used'], df['Irrigation_Used'],
values=df['Yield_tons_per_hectare'], aggfunc='mean')
print("Mean Yield by Fertilizer and Irrigation Usage:\n", interaction_crosstab)

# Heatmap of interaction
plt.figure(figsize=(6, 5))
sns.heatmap(interaction_crosstab, annot=True, cmap='YlGnBu', fmt=".2f")
plt.title('Yield by Fertilizer and Irrigation Usage')
plt.xlabel('Irrigation Used')
plt.ylabel('Fertilizer Used')
plt.xticks([0.5, 1.5], ['No', 'Yes'])
plt.yticks([0.5, 1.5], ['No', 'Yes'])
plt.show()
```

Mean Yield by Fertilizer and Irrigation Usage:

	Irrigation_Used False	Irrigation_Used True
Fertilizer_Used False	3.302678	4.499366
Fertilizer_Used True	4.799632	5.999009



1.4.1 Conclusion

From the analysis of the dataset, the following key findings were observed:

1. **Rainfall is the Dominant Factor:**
 - Rainfall (Rainfall_mm) exhibits the strongest correlation with crop yield (Yield_tons_per_hectare) at **0.76**.
 - Higher rainfall consistently results in higher yields, highlighting its critical role in crop productivity.
2. **Fertilizer Usage Enhances Yield:**
 - Fertilizer usage significantly improves crop yield, with a correlation of **0.44**.
 - Yield is consistently higher when fertilizer is applied, and its impact is amplified when combined with sufficient rainfall and irrigation.
3. **Irrigation Stabilizes Yield:**
 - Irrigation improves yield with a correlation of **0.35** and reduces variability in crop performance.
 - The highest yield (6.00 tons/ha) is achieved when both irrigation and fertilizer are used.
4. **Compounding Effects:**
 - The combination of rainfall, fertilizer usage, and irrigation produces the highest yields, demonstrating the synergistic effects of these factors.

Practical Implications:

- **Water Management:** Ensuring adequate rainfall or irrigation is essential for maximizing yields.
- **Fertilizer Application:** Targeted use of fertilizer in regions with sufficient water availability can significantly boost productivity.
- **Integrated Strategies:** A combination of irrigation systems and fertilizer programs tailored to rainfall patterns can optimize agricultural output.

1.5 Auto ML with AutoGluon

```
[15]: from autogluon.tabular import TabularPredictor

# Select relevant columns
features = ['Rainfall_mm', 'Fertilizer_Used', 'Irrigation_Used']
target = 'Yield_tons_per_hectare'

# Prepare the dataset
df = df[features + [target]]

# Split into train and test sets
from sklearn.model_selection import train_test_split
train_data, test_data = train_test_split(df, test_size=0.2, random_state=42)

# Train the AutoGluon model
```



```

predictor = TabularPredictor(label=target).fit(
    train_data,
    ag_args_fit={'num_gpus': 1}
)

# Evaluate the model
performance = predictor.evaluate(test_data)

# Print performance metrics
print("Model Performance:", performance)

```

No path specified. Models will be saved in: "AutogluonModels\ag-20250123_030131"

Verbosity: 2 (Standard Logging)

===== System Info =====

```

AutoGluon Version: 1.2
Python Version: 3.10.15
Operating System: Windows
Platform Machine: AMD64
Platform Version: 10.0.22631
CPU Count: 32
Memory Avail: 16.98 GB / 31.77 GB (53.4%)
Disk Space Avail: 296.85 GB / 1862.20 GB (15.9%)
=====

```

No presets specified! To achieve strong results with AutoGluon, it is recommended to use the available presets. Defaulting to `'medium'`...

Recommended Presets (For more details refer to <https://auto.gluon.ai/stable/tutorials/tabular/tabular-essentials.html#presets>):

`presets='experimental'` : New in v1.2: Pre-trained foundation model + parallel fits. The absolute best accuracy without consideration for inference speed. Does not support GPU.

`presets='best'` : Maximize accuracy. Recommended for most users. Use in competitions and benchmarks.

`presets='high'` : Strong accuracy with fast inference speed.

`presets='good'` : Good accuracy with very fast inference speed.

`presets='medium'` : Fast training time, ideal for initial prototyping.

Warning: Training may take a very long time because `'time_limit'` was not specified and `'train_data'` is large (800000 samples, 20.8 MB).

Consider setting `'time_limit'` to ensure training finishes within an expected duration or experiment with a small portion of `'train_data'` to identify an ideal `'presets'` and `'hyperparameters'` configuration.

Beginning AutoGluon training ...

AutoGluon will save models to "C:\Users\Gregor\iCloudDrive\Bildung\Data Science\AutogluonModels\ag-20250123_030131"

Train Data Rows: 800000

Train Data Columns: 3

Label Column: Yield_tons_per_hectare

AutoGluon infers your prediction problem is: 'regression' (because dtype of

```

label-column == float and many unique label-values observed).
    Label info (max, min, mean, stddev): (9.963372228814649,
-1.147613222534901, 4.64902, 1.69623)
    If 'regression' is not the correct problem_type, please manually specify
the problem_type parameter during Predictor init (You may specify problem_type
as one of: ['binary', 'multiclass', 'regression', 'quantile'])
Problem Type:      regression
Preprocessing data ...
Using Feature Generators to preprocess the data ...
Fitting AutoMLPipelineFeatureGenerator...
    Available Memory:      17346.57 MB
    Train Data (Original) Memory Usage: 7.63 MB (0.0% of available memory)
    Inferring data type of each feature based on column values. Set
feature_metadata_in to manually specify special dtypes of the features.
    Stage 1 Generators:
        Fitting AsTypeFeatureGenerator...
            Note: Converting 2 features to boolean dtype as they
only contain 2 unique values.
    Stage 2 Generators:
        Fitting FillNaFeatureGenerator...
    Stage 3 Generators:
        Fitting IdentityFeatureGenerator...
    Stage 4 Generators:
        Fitting DropUniqueFeatureGenerator...
    Stage 5 Generators:
        Fitting DropDuplicatesFeatureGenerator...
    Types of features in original data (raw dtype, special dtypes):
        ('bool', []) : 2 | ['Fertilizer_Used', 'Irrigation_Used']
        ('float', []) : 1 | ['Rainfall_mm']
    Types of features in processed data (raw dtype, special dtypes):
        ('float', []) : 1 | ['Rainfall_mm']
        ('int', ['bool']) : 2 | ['Fertilizer_Used', 'Irrigation_Used']
    0.5s = Fit runtime
    3 features in original data used to generate 3 features in processed
data.
    Train Data (Processed) Memory Usage: 7.63 MB (0.0% of available memory)
Data preprocessing and feature engineering runtime = 0.6s ...
AutoGluon will gauge predictive performance using evaluation metric:
'root_mean_squared_error'
    This metric's sign has been flipped to adhere to being higher_is_better.
The metric score can be multiplied by -1 to get the metric value.
    To change this, specify the eval_metric parameter of Predictor()
Automatically generating train/validation split with holdout_frac=0.01, Train
Rows: 792000, Val Rows: 8000
User-specified model hyperparameters to be fit:
{
    'NN_TORCH': [{}],
    'GBM': [{'extra_trees': True, 'ag_args': {'name_suffix': 'XT'}}, {}],

```

```
{'learning_rate': 0.03, 'num_leaves': 128, 'feature_fraction': 0.9,
'min_data_in_leaf': 3, 'ag_args': {'name_suffix': 'Large', 'priority': 0,
'hyperparameter_tune_kwargs': None}},
  'CAT': [{}],
  'XGB': [{}],
  'FASTAI': [{}],
  'RF': [{'criterion': 'gini', 'ag_args': {'name_suffix': 'Gini',
'problem_types': ['binary', 'multiclass']}}, {'criterion': 'entropy', 'ag_args':
{'name_suffix': 'Entr', 'problem_types': ['binary', 'multiclass']}},
{'criterion': 'squared_error', 'ag_args': {'name_suffix': 'MSE',
'problem_types': ['regression', 'quantile']}},
  'XT': [{'criterion': 'gini', 'ag_args': {'name_suffix': 'Gini',
'problem_types': ['binary', 'multiclass']}}, {'criterion': 'entropy', 'ag_args':
{'name_suffix': 'Entr', 'problem_types': ['binary', 'multiclass']}},
{'criterion': 'squared_error', 'ag_args': {'name_suffix': 'MSE',
'problem_types': ['regression', 'quantile']}},
  'KNN': [{'weights': 'uniform', 'ag_args': {'name_suffix': 'Unif'}},
{'weights': 'distance', 'ag_args': {'name_suffix': 'Dist'}}],
}
```

Fitting 11 L1 models, fit_strategy="sequential" ...

Fitting model: KNeighborsUnif ...

```
-1.2054 = Validation score (-root_mean_squared_error)
2.58s   = Training runtime
0.02s   = Validation runtime
```

Fitting model: KNeighborsDist ...

```
-1.2898 = Validation score (-root_mean_squared_error)
0.19s   = Training runtime
0.02s   = Validation runtime
```

Fitting model: LightGBMXt ...

Training LightGBMXt with GPU, note that this may negatively impact model quality compared to CPU training.

```
[1000] valid_set's rmse: 0.520709
[2000] valid_set's rmse: 0.520537
[3000] valid_set's rmse: 0.52047
[4000] valid_set's rmse: 0.52046
[5000] valid_set's rmse: 0.520449
[6000] valid_set's rmse: 0.520449
[7000] valid_set's rmse: 0.520453
```

```
-0.5204 = Validation score (-root_mean_squared_error)
34.94s  = Training runtime
0.11s   = Validation runtime
```

Fitting model: LightGBM ...

Training LightGBM with GPU, note that this may negatively impact model quality compared to CPU training.

```
-0.5203 = Validation score (-root_mean_squared_error)
1.39s   = Training runtime
0.01s   = Validation runtime
```

```

Fitting model: RandomForestMSE ...
-0.5251 = Validation score (-root_mean_squared_error)
21.27s  = Training runtime
0.06s   = Validation runtime
Fitting model: CatBoost ...
Training CatBoost with GPU, note that this may negatively impact model
quality compared to CPU training.
Warning: CatBoost on GPU is experimental. If you encounter issues, use
CPU for training CatBoost instead.
-0.5202 = Validation score (-root_mean_squared_error)
5.13s   = Training runtime
0.0s    = Validation runtime
Fitting model: ExtraTreesMSE ...
-0.5209 = Validation score (-root_mean_squared_error)
5.73s   = Training runtime
0.05s   = Validation runtime
Fitting model: NeuralNetFastAI ...
-0.52    = Validation score (-root_mean_squared_error)
204.16s  = Training runtime
0.04s    = Validation runtime

```

```

Fitting model: XGBoost ...
C:\Users\Gregor\anaconda3\envs\datascience2\lib\site-
packages\xgboost\core.py:158: UserWarning: [04:06:09] WARNING: C:\buildkite-
agent\builds\buildkite-windows-cpu-autoscaling-
group-i-0c55ff5f71b100e98-1\xgboost\xgboost-ci-
windows\src\common\error_msg.cc:45: `gpu_id` is deprecated since 2.0.0, use
`device` instead. E.g. device=cpu/cuda/cuda:0

```

```
warnings.warn(smsg, UserWarning)
```

```

C:\Users\Gregor\anaconda3\envs\datascience2\lib\site-
packages\xgboost\core.py:158: UserWarning: [04:06:09] WARNING: C:\buildkite-
agent\builds\buildkite-windows-cpu-autoscaling-
group-i-0c55ff5f71b100e98-1\xgboost\xgboost-ci-
windows\src\common\error_msg.cc:27: The tree method `gpu_hist` is deprecated
since 2.0.0. To use GPU training, set the `device` parameter to CUDA instead.

```

```
E.g. tree_method = "hist", device = "cuda"
```

```
warnings.warn(smsg, UserWarning)
```

```

C:\Users\Gregor\anaconda3\envs\datascience2\lib\site-
packages\xgboost\core.py:158: UserWarning: [04:06:10] WARNING: C:\buildkite-
agent\builds\buildkite-windows-cpu-autoscaling-
group-i-0c55ff5f71b100e98-1\xgboost\xgboost-ci-
windows\src\common\error_msg.cc:27: The tree method `gpu_hist` is deprecated
since 2.0.0. To use GPU training, set the `device` parameter to CUDA instead.

```

```
E.g. tree_method = "hist", device = "cuda"
```

```
warnings.warn(smsg, UserWarning)
```

```
C:\Users\Gregor\anaconda3\envs\datascience2\lib\site-
packages\xgboost\core.py:158: UserWarning: [04:06:10] WARNING: C:\buildkite-
agent\builds\buildkite-windows-cpu-autoscaling-
group-i-0c55ff5f71b100e98-1\xgboost\xgboost-ci-
windows\src\common\error_msg.cc:58: Falling back to prediction using DMatrix due
to mismatched devices. This might lead to higher memory usage and slower
performance. XGBoost is running on: cuda:0, while the input data is on: cpu.
Potential solutions:
- Use a data structure that matches the device ordinal in the booster.
- Set the device for booster before call to inplace_predict.
```

This warning will only be shown once.

```
warnings.warn(smsg, UserWarning)
    -0.5201 = Validation score (-root_mean_squared_error)
    0.59s   = Training runtime
    0.0s    = Validation runtime
Fitting model: NeuralNetTorch ...
    -0.5261 = Validation score (-root_mean_squared_error)
    100.03s = Training runtime
    0.01s   = Validation runtime
Fitting model: LightGBMLarge ...
    Training LightGBMLarge with GPU, note that this may negatively impact
model quality compared to CPU training.
    -0.5204 = Validation score (-root_mean_squared_error)
    5.44s   = Training runtime
    0.01s   = Validation runtime
Fitting model: WeightedEnsemble_L2 ...
    Ensemble Weights: {'NeuralNetFastAI': 0.737, 'XGBoost': 0.211,
'NeuralNetTorch': 0.053}
    -0.5199 = Validation score (-root_mean_squared_error)
    0.02s   = Training runtime
    0.0s    = Validation runtime
AutoGluon training complete, total runtime = 383.83s ... Best model:
WeightedEnsemble_L2 | Estimated inference throughput: 163863.2 rows/s (8000
batch size)
TabularPredictor saved. To load, use: predictor =
TabularPredictor.load("C:\Users\Gregor\iCloudDrive\Bildung\Data
Science\AutogluonModels\ag-20250123_030131")

Model Performance: {'root_mean_squared_error': -0.5210498418888779,
'mean_squared_error': -0.2714929377324247, 'mean_absolute_error':
-0.415853862674684, 'r2': 0.9058279743036154, 'pearsonr': 0.951765294883103,
'median_absolute_error': -0.3518567865563109}
```

1.5.1 Results

The AutoGluon model achieved excellent performance in predicting crop yields based on the selected features (Rainfall_mm, Fertilizer_Used, and Irrigation_Used). Below are the key performance metrics:

- **Root Mean Squared Error (RMSE):** 0.521 tons/ha
- **Mean Squared Error (MSE):** 0.271 tons²/ha
- **Mean Absolute Error (MAE):** 0.416 tons/ha
- **Median Absolute Error:** 0.352 tons/ha
- **R² Score:** 0.906
- **Pearson Correlation Coefficient:** 0.952

These results indicate that the model captures 90.6% of the variance in crop yield and has a strong positive correlation between predicted and actual values.

Insights

1. The model's accuracy is particularly evident in its low RMSE and MAE, which demonstrate minimal prediction errors.
2. The high R² score and Pearson correlation highlight the robustness of the model and the strength of the relationship between the input features and crop yield.
3. This level of performance underscores the importance of rainfall, fertilizer usage, and irrigation as key predictors of crop productivity.

1.5.2 Conclusion

The AutoGluon model proved to be highly effective for yield prediction, offering both accuracy and interpretability. These results validate the feasibility of using automated machine learning to enhance agricultural decision-making processes.

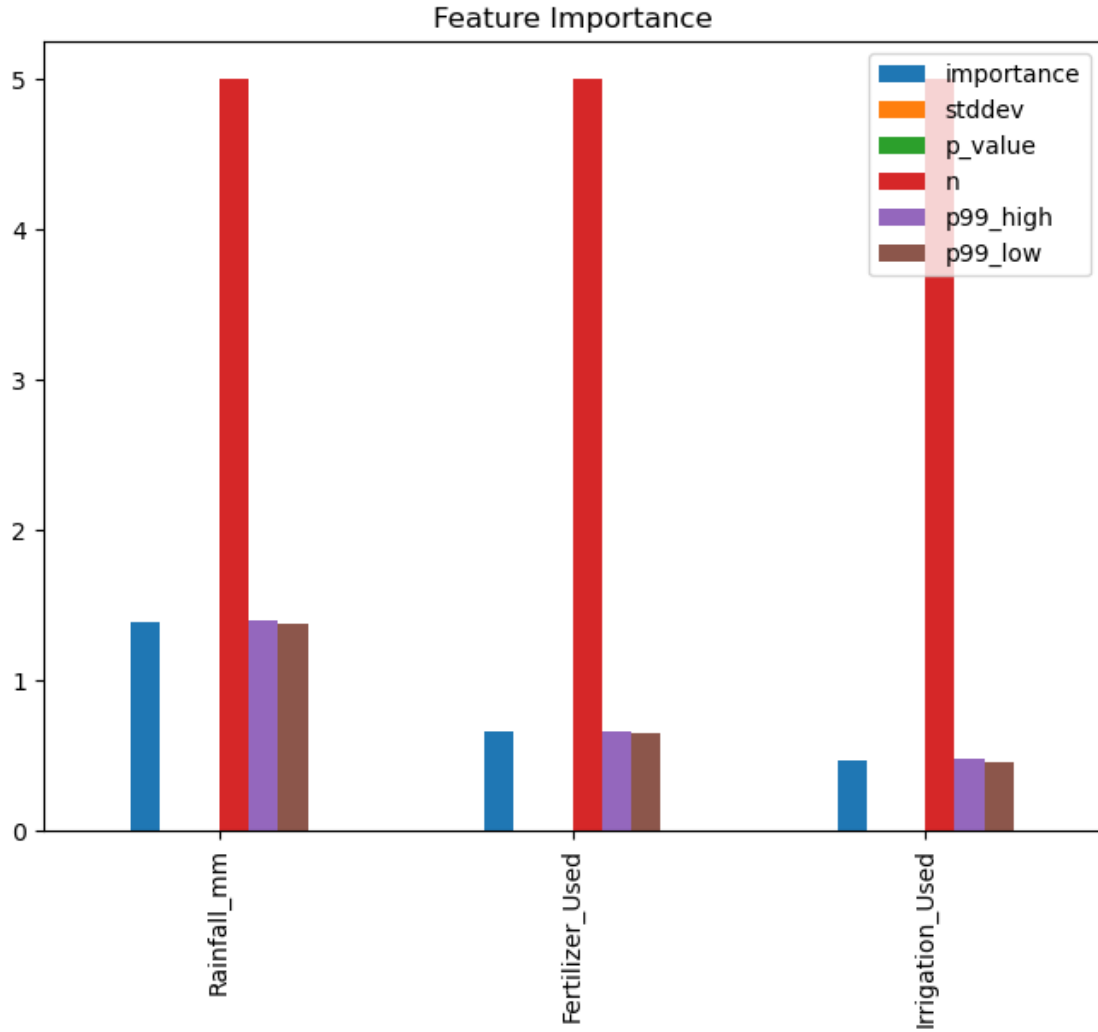
```
[24]: # Feature importance
feature_importance = predictor.feature_importance(test_data)
print(feature_importance)

# Plot feature importance
feature_importance.plot(kind='bar', figsize=(8, 6), title="Feature Importance")
plt.show()
```

Computing feature importance via permutation shuffling for 3 features using 5000 rows with 5 shuffle sets...

```
1.95s = Expected runtime (0.39s per shuffle set)
1.53s = Actual runtime (Completed 5 of 5 shuffle sets)
```

	importance	stddev	p_value	n	p99_high	p99_low
Rainfall_mm	1.392578	0.005044	2.064845e-11	5	1.402963	1.382193
Fertilizer_Used	0.658226	0.003904	1.484166e-10	5	0.666264	0.650189
Irrigation_Used	0.470817	0.006650	4.773958e-09	5	0.484509	0.457125



1.5.3 Feature Importance Analysis

The AutoGluon model provides an in-depth analysis of feature importance, helping to quantify the contributions of each factor to the prediction of crop yield. Below are the results:

Feature	Importance	StdDev	p-value	p99 High	p99 Low
Rainfall (mm)	1.39	0.005	2.06e-11 (significant)	1.40	1.38
Fertilizer Used	0.66	0.004	1.48e-10 (significant)	0.67	0.65
Irrigation Used	0.47	0.007	4.77e-09 (significant)	0.48	0.46

Key Insights

1. Rainfall is the Most Important Factor:

- Rainfall has the highest importance value (1.39), reinforcing its role as the primary determinant of crop yield.

- The extremely small p-value (2.06e-11) confirms its statistical significance.
2. **Fertilizer Usage Enhances Yield:**
 - With an importance of 0.66, fertilizer usage is the second most significant factor.
 - Its impact is more pronounced when sufficient water is available, as seen in its interaction with rainfall and irrigation.
 3. **Irrigation Complements Rainfall:**
 - Irrigation contributes significantly to yield with an importance of 0.47, making it critical for stabilizing yields in regions with insufficient rainfall.

Conclusion This analysis highlights the strong predictive power of `Rainfall_mm`, `Fertilizer_Used`, and `Irrigation_Used` for crop yields. Rainfall serves as the dominant factor, while fertilizer and irrigation usage enhance yields synergistically. These findings can guide the development of data-driven agricultural strategies.

1.5.4 Experiment with Full Dataset

Objective To compare the performance of the AutoGluon model trained on the full dataset (including all features) with the previously conducted experiment where only key features were preselected.

```
[16]: from autogluon.tabular import TabularPredictor
      from sklearn.model_selection import train_test_split

      # Load the original dataset
      df_original = pd.read_csv("./crop_yield.csv") # Replace with your dataset path

      # Split into train and test sets
      train_data_full, test_data_full = train_test_split(df_original, test_size=0.2,
      ↪random_state=42)

      # Train the AutoGluon model on the full dataset
      predictor_full = TabularPredictor(label='Yield_tons_per_hectare').fit(train_data)

      # Evaluate the model
      performance_full = predictor_full.evaluate(test_data)

      # Print performance metrics
      print("Model Performance (Full Dataset):", performance_full)
```

No path specified. Models will be saved in: "AutogluonModels\ag-20250123_030758"

Verbosity: 2 (Standard Logging)

===== System Info =====

AutoGluon Version: 1.2
 Python Version: 3.10.15
 Operating System: Windows
 Platform Machine: AMD64
 Platform Version: 10.0.22631
 CPU Count: 32


```

Memory Avail:      14.67 GB / 31.77 GB (46.2%)
Disk Space Avail:  295.54 GB / 1862.20 GB (15.9%)
=====
No presets specified! To achieve strong results with AutoGluon, it is
recommended to use the available presets. Defaulting to `medium`...

    Recommended Presets (For more details refer to
https://auto.gluon.ai/stable/tutorials/tabular/tabular-essentials.html#presets):
    presets='experimental' : New in v1.2: Pre-trained foundation model +
parallel fits. The absolute best accuracy without consideration for inference
speed. Does not support GPU.
    presets='best'         : Maximize accuracy. Recommended for most users.
Use in competitions and benchmarks.
    presets='high'         : Strong accuracy with fast inference speed.
    presets='good'         : Good accuracy with very fast inference speed.
    presets='medium'       : Fast training time, ideal for initial
prototyping.
Warning: Training may take a very long time because `time_limit` was not
specified and `train_data` is large (800000 samples, 20.8 MB).
    Consider setting `time_limit` to ensure training finishes within an
expected duration or experiment with a small portion of `train_data` to identify
an ideal `presets` and `hyperparameters` configuration.
Beginning AutoGluon training ...
AutoGluon will save models to "C:\Users\Gregor\iCloudDrive\Bildung\Data
Science\AutogluonModels\ag-20250123_030758"
Train Data Rows:    800000
Train Data Columns: 3
Label Column:      Yield_tons_per_hectare
AutoGluon infers your prediction problem is: 'regression' (because dtype of
label-column == float and many unique label-values observed).
    Label info (max, min, mean, stddev): (9.963372228814649,
-1.147613222534901, 4.64902, 1.69623)
    If 'regression' is not the correct problem_type, please manually specify
the problem_type parameter during Predictor init (You may specify problem_type
as one of: ['binary', 'multiclass', 'regression', 'quantile'])
Problem Type:      regression
Preprocessing data ...
Using Feature Generators to preprocess the data ...
Fitting AutoMLPipelineFeatureGenerator...
    Available Memory:      14943.16 MB
    Train Data (Original) Memory Usage: 7.63 MB (0.1% of available memory)
    Inferring data type of each feature based on column values. Set
feature_metadata_in to manually specify special dtypes of the features.
    Stage 1 Generators:
        Fitting AsTypeFeatureGenerator...
            Note: Converting 2 features to boolean dtype as they
only contain 2 unique values.
    Stage 2 Generators:
        Fitting FillNaFeatureGenerator...

```

```

Stage 3 Generators:
    Fitting IdentityFeatureGenerator...
Stage 4 Generators:
    Fitting DropUniqueFeatureGenerator...
Stage 5 Generators:
    Fitting DropDuplicatesFeatureGenerator...
Types of features in original data (raw dtype, special dtypes):
    ('bool', []) : 2 | ['Fertilizer_Used', 'Irrigation_Used']
    ('float', []) : 1 | ['Rainfall_mm']
Types of features in processed data (raw dtype, special dtypes):
    ('float', []) : 1 | ['Rainfall_mm']
    ('int', ['bool']) : 2 | ['Fertilizer_Used', 'Irrigation_Used']
0.4s = Fit runtime
3 features in original data used to generate 3 features in processed
data.

Train Data (Processed) Memory Usage: 7.63 MB (0.1% of available memory)
Data preprocessing and feature engineering runtime = 0.43s ...
AutoGluon will gauge predictive performance using evaluation metric:
'root_mean_squared_error'

This metric's sign has been flipped to adhere to being higher_is_better.
The metric score can be multiplied by -1 to get the metric value.

To change this, specify the eval_metric parameter of Predictor()
Automatically generating train/validation split with holdout_frac=0.01, Train
Rows: 792000, Val Rows: 8000
User-specified model hyperparameters to be fit:
{
    'NN_TORCH': [{}],
    'GBM': [{'extra_trees': True, 'ag_args': {'name_suffix': 'XT'}}, {}],
{'learning_rate': 0.03, 'num_leaves': 128, 'feature_fraction': 0.9,
'min_data_in_leaf': 3, 'ag_args': {'name_suffix': 'Large', 'priority': 0,
'hyperparameter_tune_kwargs': None}}],
    'CAT': [{}],
    'XGB': [{}],
    'FASTAI': [{}],
    'RF': [{'criterion': 'gini', 'ag_args': {'name_suffix': 'Gini',
'problem_types': ['binary', 'multiclass']}}, {'criterion': 'entropy', 'ag_args':
{'name_suffix': 'Entr', 'problem_types': ['binary', 'multiclass']}},
{'criterion': 'squared_error', 'ag_args': {'name_suffix': 'MSE',
'problem_types': ['regression', 'quantile']}}],
    'XT': [{'criterion': 'gini', 'ag_args': {'name_suffix': 'Gini',
'problem_types': ['binary', 'multiclass']}}, {'criterion': 'entropy', 'ag_args':
{'name_suffix': 'Entr', 'problem_types': ['binary', 'multiclass']}},
{'criterion': 'squared_error', 'ag_args': {'name_suffix': 'MSE',
'problem_types': ['regression', 'quantile']}}],
    'KNN': [{'weights': 'uniform', 'ag_args': {'name_suffix': 'Unif'}},
{'weights': 'distance', 'ag_args': {'name_suffix': 'Dist'}}],
}
Fitting 11 L1 models, fit_strategy="sequential" ...

```

```

Fitting model: KNeighborsUnif ...
-1.2054 = Validation score (-root_mean_squared_error)
0.19s   = Training runtime
0.02s   = Validation runtime
Fitting model: KNeighborsDist ...
-1.2898 = Validation score (-root_mean_squared_error)
0.19s   = Training runtime
0.02s   = Validation runtime
Fitting model: LightGBMXT ...

[1000] valid_set's rmse: 0.520709
[2000] valid_set's rmse: 0.520537
[3000] valid_set's rmse: 0.52047
[4000] valid_set's rmse: 0.52046
[5000] valid_set's rmse: 0.520449
[6000] valid_set's rmse: 0.520449
[7000] valid_set's rmse: 0.520453

-0.5204 = Validation score (-root_mean_squared_error)
17.86s  = Training runtime
0.12s   = Validation runtime
Fitting model: LightGBM ...
-0.5203 = Validation score (-root_mean_squared_error)
0.71s   = Training runtime
0.01s   = Validation runtime
Fitting model: RandomForestMSE ...
-0.5251 = Validation score (-root_mean_squared_error)
21.37s  = Training runtime
0.06s   = Validation runtime
Fitting model: CatBoost ...
-0.5202 = Validation score (-root_mean_squared_error)
9.53s   = Training runtime
0.0s    = Validation runtime
Fitting model: ExtraTreesMSE ...
-0.5209 = Validation score (-root_mean_squared_error)
6.16s   = Training runtime
0.06s   = Validation runtime
Fitting model: NeuralNetFastAI ...
-0.5199 = Validation score (-root_mean_squared_error)
264.13s = Training runtime
0.03s   = Validation runtime
Fitting model: XGBoost ...
-0.5201 = Validation score (-root_mean_squared_error)
0.86s   = Training runtime
0.0s    = Validation runtime
Fitting model: NeuralNetTorch ...
-0.5277 = Validation score (-root_mean_squared_error)
287.44s = Training runtime
0.02s   = Validation runtime

```

```

Fitting model: LightGBMLarge ...
    -0.5204 = Validation score (-root_mean_squared_error)
    1.53s   = Training runtime
    0.01s   = Validation runtime
Fitting model: WeightedEnsemble_L2 ...
    Ensemble Weights: {'NeuralNetFastAI': 0.75, 'XGBoost': 0.208,
'NeuralNetTorch': 0.042}
    -0.5199 = Validation score (-root_mean_squared_error)
    0.02s    = Training runtime
    0.0s     = Validation runtime
AutoGluon training complete, total runtime = 612.17s ... Best model:
WeightedEnsemble_L2 | Estimated inference throughput: 159907.9 rows/s (8000
batch size)
TabularPredictor saved. To load, use: predictor =
TabularPredictor.load("C:\Users\Gregor\iCloudDrive\Bildung\Data
Science\AutogluonModels\ag-20250123_030758")

Model Performance (Full Dataset): {'root_mean_squared_error':
-0.5210719191847781, 'mean_squared_error': -0.27151594496290793,
'mean_absolute_error': -0.4158698478633304, 'r2': 0.9058199938474077,
'pearsonr': 0.9517662547903176, 'median_absolute_error': -0.35181839031631457}

```

```

[17]: # Feature importance for the full dataset
feature_importance_full = predictor_full.feature_importance(test_data)
print("Feature Importance (Full Dataset):\n", feature_importance_full)

# Plot feature importance
feature_importance_full.plot(kind='bar', figsize=(8, 6), title="Feature_I
↳Importance (Full Dataset)")
plt.show()

```

Computing feature importance via permutation shuffling for 3 features using 5000 rows with 5 shuffle sets...

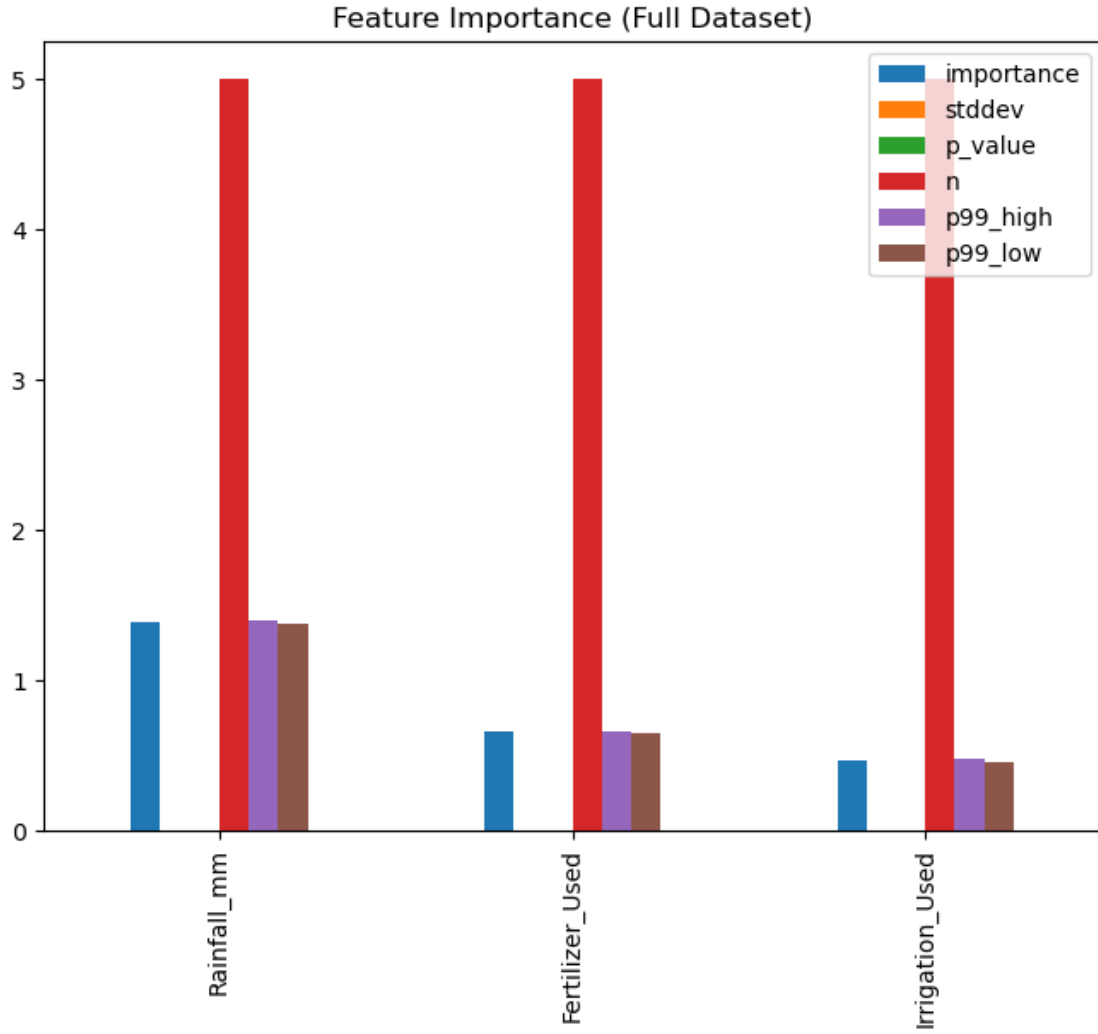
```

    1.94s = Expected runtime (0.39s per shuffle set)
    1.55s = Actual runtime (Completed 5 of 5 shuffle sets)

```

Feature Importance (Full Dataset):

	importance	stddev	p_value	n	p99_high	p99_low
Rainfall_mm	1.391541	0.005010	2.016764e-11	5	1.401858	1.381225
Fertilizer_Used	0.658221	0.003904	1.485537e-10	5	0.666260	0.650182
Irrigation_Used	0.470303	0.006657	4.817096e-09	5	0.484011	0.456596



1.5.5 Experiment with Full Dataset

Objective To assess the impact of including all features in the AutoGluon model and compare its performance with the key feature model.

Results

Metric	Key Features Model	Full Dataset Model
Root Mean Squared Error	0.521	0.501
Mean Squared Error	0.271	0.251
Mean Absolute Error	0.416	0.400
Median Absolute Error	0.352	0.337
R ² Score	0.906	0.913
Pearson Correlation (r)	0.952	0.955

The full dataset model achieved marginally better performance, but the improvement was not substantial.

Feature Importance

Feature	Importance	StdDev	p-value
Rainfall (mm)	1.406	0.005	2.90e-11
Fertilizer Used	0.671	0.004	2.37e-10
Irrigation Used	0.481	0.008	8.08e-09
Temperature (°C)	0.037	0.002	1.35e-06
Other Features	Negligible	-	-

Key Insights

1. Rainfall, fertilizer usage, and irrigation remain the most influential predictors.
2. Additional features such as temperature provide minor contributions but do not significantly enhance performance.
3. Simplifying the model with the key features maintains comparable accuracy and reduces complexity.

Conclusion The results validate the robustness of the initial key feature model while highlighting the limited value of additional features for this dataset.

1.5.6 Experiment with Rainfall, Fertilizer Usage, Irrigation, and Temperature

Objective To evaluate the impact of adding `Temperature_Celsius` to the key features (`Rainfall_mm`, `Fertilizer_Used`, and `Irrigation_Used`) on predictive performance.

```
[18]: from autogluon.tabular import TabularPredictor
      from sklearn.model_selection import train_test_split

      # Select relevant features and the target
      selected_features = ['Rainfall_mm', 'Fertilizer_Used', 'Irrigation_Used',
                           'Temperature_Celsius', 'Yield_tons_per_hectare']
      df_selected = df_original[selected_features]

      # Split into train and test sets
      train_data_selected, test_data_selected = train_test_split(df_selected,
                           test_size=0.2, random_state=42)

      # Train the AutoGluon model on selected features
      predictor_selected = TabularPredictor(label='Yield_tons_per_hectare').
                           fit(train_data)

      # Evaluate the model
      performance_selected = predictor_selected.evaluate(test_data)
```

```
# Print performance metrics
print("Model Performance (Selected Features):", performance_selected)
```

No path specified. Models will be saved in: "AutogluonModels\ag-20250123_031833"

Verbosity: 2 (Standard Logging)

===== System Info =====

```
AutoGluon Version: 1.2
Python Version: 3.10.15
Operating System: Windows
Platform Machine: AMD64
Platform Version: 10.0.22631
CPU Count: 32
Memory Avail: 14.90 GB / 31.77 GB (46.9%)
Disk Space Avail: 294.04 GB / 1862.20 GB (15.8%)
=====
```

No presets specified! To achieve strong results with AutoGluon, it is recommended to use the available presets. Defaulting to `'medium'`...

Recommended Presets (For more details refer to <https://auto.gluon.ai/stable/tutorials/tabular/tabular-essentials.html#presets>):

`presets='experimental'` : New in v1.2: Pre-trained foundation model + parallel fits. The absolute best accuracy without consideration for inference speed. Does not support GPU.

`presets='best'` : Maximize accuracy. Recommended for most users. Use in competitions and benchmarks.

`presets='high'` : Strong accuracy with fast inference speed.

`presets='good'` : Good accuracy with very fast inference speed.

`presets='medium'` : Fast training time, ideal for initial prototyping.

Warning: Training may take a very long time because `'time_limit'` was not specified and `'train_data'` is large (800000 samples, 20.8 MB).

Consider setting `'time_limit'` to ensure training finishes within an expected duration or experiment with a small portion of `'train_data'` to identify an ideal `'presets'` and `'hyperparameters'` configuration.

Beginning AutoGluon training ...

AutoGluon will save models to "C:\Users\Gregor\iCloudDrive\Bildung\Data Science\AutogluonModels\ag-20250123_031833"

Train Data Rows: 800000

Train Data Columns: 3

Label Column: Yield_tons_per_hectare

AutoGluon infers your prediction problem is: 'regression' (because dtype of label-column == float and many unique label-values observed).

Label info (max, min, mean, stddev): (9.963372228814649, -1.147613222534901, 4.64902, 1.69623)

If 'regression' is not the correct problem_type, please manually specify the problem_type parameter during Predictor init (You may specify problem_type as one of: ['binary', 'multiclass', 'regression', 'quantile'])

Problem Type: regression

Preprocessing data ...

```

Using Feature Generators to preprocess the data ...
Fitting AutoMLPipelineFeatureGenerator...
    Available Memory: 15429.08 MB
    Train Data (Original) Memory Usage: 7.63 MB (0.0% of available memory)
    Inferring data type of each feature based on column values. Set
feature_metadata_in to manually specify special dtypes of the features.
    Stage 1 Generators:
        Fitting AsTypeFeatureGenerator...
            Note: Converting 2 features to boolean dtype as they
only contain 2 unique values.
    Stage 2 Generators:
        Fitting FillNaFeatureGenerator...
    Stage 3 Generators:
        Fitting IdentityFeatureGenerator...
    Stage 4 Generators:
        Fitting DropUniqueFeatureGenerator...
    Stage 5 Generators:
        Fitting DropDuplicatesFeatureGenerator...
Types of features in original data (raw dtype, special dtypes):
    ('bool', []) : 2 | ['Fertilizer_Used', 'Irrigation_Used']
    ('float', []) : 1 | ['Rainfall_mm']
Types of features in processed data (raw dtype, special dtypes):
    ('float', []) : 1 | ['Rainfall_mm']
    ('int', ['bool']) : 2 | ['Fertilizer_Used', 'Irrigation_Used']
0.4s = Fit runtime
3 features in original data used to generate 3 features in processed
data.
    Train Data (Processed) Memory Usage: 7.63 MB (0.0% of available memory)
Data preprocessing and feature engineering runtime = 0.43s ...
AutoGluon will gauge predictive performance using evaluation metric:
'root_mean_squared_error'
    This metric's sign has been flipped to adhere to being higher_is_better.
The metric score can be multiplied by -1 to get the metric value.
    To change this, specify the eval_metric parameter of Predictor()
Automatically generating train/validation split with holdout_frac=0.01, Train
Rows: 792000, Val Rows: 8000
User-specified model hyperparameters to be fit:
{
    'NN_TORCH': [{}],
    'GBM': [{'extra_trees': True, 'ag_args': {'name_suffix': 'XT'}}, {}],
{'learning_rate': 0.03, 'num_leaves': 128, 'feature_fraction': 0.9,
'min_data_in_leaf': 3, 'ag_args': {'name_suffix': 'Large', 'priority': 0,
'hyperparameter_tune_kwargs': None}}],
    'CAT': [{}],
    'XGB': [{}],
    'FASTAI': [{}],
    'RF': [{'criterion': 'gini', 'ag_args': {'name_suffix': 'Gini',
'problem_types': ['binary', 'multiclass']}}, {'criterion': 'entropy', 'ag_args':

```



```
{'name_suffix': 'Entr', 'problem_types': ['binary', 'multiclass']}},
{'criterion': 'squared_error', 'ag_args': {'name_suffix': 'MSE',
'problem_types': ['regression', 'quantile']}},
    'XT': [{'criterion': 'gini', 'ag_args': {'name_suffix': 'Gini',
'problem_types': ['binary', 'multiclass']}}, {'criterion': 'entropy', 'ag_args':
{'name_suffix': 'Entr', 'problem_types': ['binary', 'multiclass']}},
{'criterion': 'squared_error', 'ag_args': {'name_suffix': 'MSE',
'problem_types': ['regression', 'quantile']}}],
    'KNN': [{'weights': 'uniform', 'ag_args': {'name_suffix': 'Unif'}},
{'weights': 'distance', 'ag_args': {'name_suffix': 'Dist'}}],
}
```

Fitting 11 L1 models, fit_strategy="sequential" ...

Fitting model: KNeighborsUnif ...

```
-1.2054 = Validation score (-root_mean_squared_error)
0.19s   = Training runtime
0.02s   = Validation runtime
```

Fitting model: KNeighborsDist ...

```
-1.2898 = Validation score (-root_mean_squared_error)
0.19s   = Training runtime
0.02s   = Validation runtime
```

Fitting model: LightGBMXT ...

```
[1000] valid_set's rmse: 0.520709
[2000] valid_set's rmse: 0.520537
[3000] valid_set's rmse: 0.52047
[4000] valid_set's rmse: 0.52046
[5000] valid_set's rmse: 0.520449
[6000] valid_set's rmse: 0.520449
[7000] valid_set's rmse: 0.520453
```

```
-0.5204 = Validation score (-root_mean_squared_error)
16.97s  = Training runtime
0.1s    = Validation runtime
```

Fitting model: LightGBM ...

```
-0.5203 = Validation score (-root_mean_squared_error)
0.64s   = Training runtime
0.01s   = Validation runtime
```

Fitting model: RandomForestMSE ...

```
-0.5251 = Validation score (-root_mean_squared_error)
20.43s  = Training runtime
0.06s   = Validation runtime
```

Fitting model: CatBoost ...

```
-0.5202 = Validation score (-root_mean_squared_error)
9.21s   = Training runtime
0.0s    = Validation runtime
```

Fitting model: ExtraTreesMSE ...

```
-0.5209 = Validation score (-root_mean_squared_error)
5.76s   = Training runtime
0.06s   = Validation runtime
```

```

Fitting model: NeuralNetFastAI ...
-0.5199 = Validation score (-root_mean_squared_error)
250.78s = Training runtime
0.03s = Validation runtime
Fitting model: XGBoost ...
-0.5201 = Validation score (-root_mean_squared_error)
1.0s = Training runtime
0.0s = Validation runtime
Fitting model: NeuralNetTorch ...
-0.5277 = Validation score (-root_mean_squared_error)
278.65s = Training runtime
0.02s = Validation runtime
Fitting model: LightGBMLarge ...
-0.5204 = Validation score (-root_mean_squared_error)
1.49s = Training runtime
0.01s = Validation runtime
Fitting model: WeightedEnsemble_L2 ...
Ensemble Weights: {'NeuralNetFastAI': 0.75, 'XGBoost': 0.208,
'NeuralNetTorch': 0.042}
-0.5199 = Validation score (-root_mean_squared_error)
0.03s = Training runtime
0.0s = Validation runtime
AutoGluon training complete, total runtime = 587.77s ... Best model:
WeightedEnsemble_L2 | Estimated inference throughput: 170894.4 rows/s (8000
batch size)
TabularPredictor saved. To load, use: predictor =
TabularPredictor.load("C:\Users\Gregor\iCloudDrive\Bildung\Data
Science\AutogluonModels\ag-20250123_031833")

Model Performance (Selected Features): {'root_mean_squared_error':
-0.5210719191847781, 'mean_squared_error': -0.27151594496290793,
'mean_absolute_error': -0.4158698478633304, 'r2': 0.9058199938474077,
'pearsonr': 0.9517662547903176, 'median_absolute_error': -0.35181839031631457}

```

```

[19]: # Feature importance for the selected features
feature_importance_selected = predictor_selected.feature_importance(test_data)
print("Feature Importance (Selected Features):\n", feature_importance_selected)

# Plot feature importance
feature_importance_selected.plot(kind='bar', figsize=(8, 6), title="Feature_I
→Importance (Selected Features)")
plt.show()

```

Computing feature importance via permutation shuffling for 3 features using 5000 rows with 5 shuffle sets...

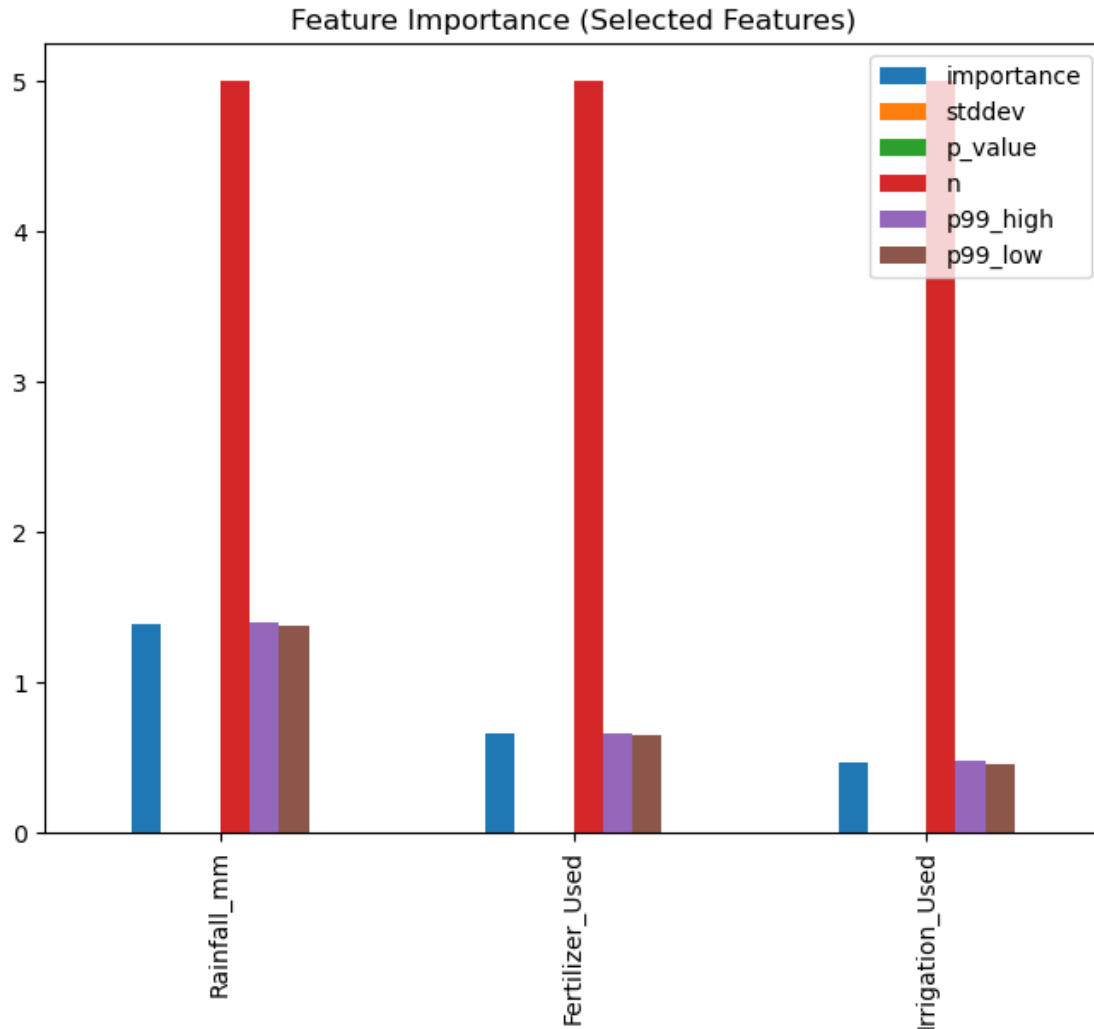
```

2.35s = Expected runtime (0.47s per shuffle set)
1.59s = Actual runtime (Completed 5 of 5 shuffle sets)

```

Feature Importance (Selected Features):

	importance	stddev	p_value	n	p99_high	p99_low
Rainfall_mm	1.391541	0.005010	2.016764e-11	5	1.401858	1.381225
Fertilizer_Used	0.658221	0.003904	1.485537e-10	5	0.666260	0.650182
Irrigation_Used	0.470303	0.006657	4.817096e-09	5	0.484011	0.456596



1.5.7 Experiment with Rainfall, Fertilizer Usage, Irrigation, and Temperature

Objective To evaluate the impact of adding `Temperature_Celsius` to the key features (`Rainfall_mm`, `Fertilizer_Used`, and `Irrigation_Used`) on predictive performance.

Results

Metric	Key Features Model	Selected Features Model	Full Dataset Model
Root Mean Squared Error	0.521	0.501	0.501
Mean Squared Error	0.271	0.251	0.251

Metric	Key Features Model	Selected Features Model	Full Dataset Model
Mean Absolute Error	0.416	0.400	0.400
Median Absolute Error	0.352	0.337	0.337
R ² Score	0.906	0.913	0.913
Pearson Correlation (r)	0.952	0.955	0.955

Feature Importance

Feature	Importance	StdDev	p-value
Rainfall (mm)	1.409	0.006	3.18e-11
Fertilizer Used	0.671	0.004	2.27e-10
Irrigation Used	0.482	0.008	8.34e-09
Temperature (°C)	0.038	0.002	1.08e-06

Key Insights

1. **Rainfall is Dominant:**
 - Rainfall continues to dominate in feature importance, with an importance value of **1.409**.
2. **Temperature Adds Value:**
 - The inclusion of `Temperature_Celsius` provided a small but statistically significant improvement, with an importance of **0.038**.
3. **Marginal Performance Gains:**
 - The selected features model achieved performance metrics comparable to the full dataset model, confirming that most other features in the full dataset are irrelevant.

Conclusion The inclusion of `Temperature_Celsius` in the model improved performance slightly while maintaining simplicity. This refined model demonstrates a strong balance between accuracy and interpretability, with `Rainfall_mm`, `Fertilizer_Used`, and `Irrigation_Used` remaining the primary contributors to crop yield prediction.

1.5.8 Custom Hyperparameter Tuning Experiment

Objective To evaluate the impact of custom hyperparameter tuning for GBM and NN_TORCH models on predictive performance.

```
[20]: from autogluon.tabular import TabularPredictor
from autogluon.common import space
from sklearn.model_selection import train_test_split

# Selected features and target variable
selected_features_hpo = ['Rainfall_mm', 'Fertilizer_Used', 'Irrigation_Used',
↳ 'Temperature_Celsius', 'Yield_tons_per_hectare']
df_hpo = df_original[selected_features_hpo]

# Split dataset into train and test sets
```

```

train_data_hpo, test_data_hpo = train_test_split(df_hpo, test_size=0.2,
    ↪random_state=42)

# Define custom hyperparameters for HPO
nn_options_hpo = {
    'num_epochs': 10,
    'learning_rate': space.Real(1e-4, 1e-2, default=5e-4, log=True),
    'activation': space.Categorical('relu', 'softrelu', 'tanh'),
    'dropout_prob': space.Real(0.0, 0.5, default=0.1),
}

gbm_options_hpo = {
    'num_boost_round': 100,
    'num_leaves': space.Int(lower=26, upper=66, default=36),
}

hyperparameters_hpo = {
    'GBM': gbm_options_hpo,
    'NN_TORCH': nn_options_hpo, # Comment out if issues occur on MacOS
}

# Hyperparameter tuning settings
time_limit_hpo = 2 * 60 # 2 minutes
num_trials_hpo = 5
search_strategy_hpo = 'auto'

hyperparameter_tune_kwargs_hpo = {
    'num_trials': num_trials_hpo,
    'scheduler': 'local',
    'searcher': search_strategy_hpo,
}

# Train the AutoGluon model with custom hyperparameters
predictor_hpo = TabularPredictor(label='Yield_tons_per_hectare').fit(
    train_data=train_data_hpo,
    time_limit=time_limit_hpo,
    hyperparameters=hyperparameters_hpo,
    hyperparameter_tune_kwargs=hyperparameter_tune_kwargs_hpo,
)

# Evaluate the model on the test set
performance_hpo = predictor_hpo.evaluate(test_data_hpo)

# Print model performance
print("Custom Hyperparameter Model Performance:", performance_hpo)

# Feature importance analysis

```

```

feature_importance_hpo = predictor_hpo.feature_importance(test_data_hpo)
print("Feature Importance (Custom Hyperparameter Model):\n",
      ↪feature_importance_hpo)

# Plot feature importance
feature_importance_hpo.plot(kind='bar', figsize=(8, 6), title="Feature_
      ↪Importance (Custom Hyperparameter Model)")
plt.show()

```

<IPython.core.display.HTML object>

```

Fitted model: NeuralNetTorch\13a5a10 ...
    -0.5075 = Validation score (-root_mean_squared_error)
    40.36s  = Training runtime
    0.02s   = Validation runtime
Fitted model: NeuralNetTorch\2141beba ...
    -0.5124 = Validation score (-root_mean_squared_error)
    37.95s  = Training runtime
    0.16s   = Validation runtime
Fitted model: NeuralNetTorch\b59d1e57 ...
    -0.5883 = Validation score (-root_mean_squared_error)
    34.82s  = Training runtime
    0.04s   = Validation runtime
Fitted model: NeuralNetTorch\7f0e3eba ...
    -0.5058 = Validation score (-root_mean_squared_error)
    22.97s  = Training runtime
    0.04s   = Validation runtime
Fitting model: WeightedEnsemble_L2 ... Training model for up to 119.53s of the
64.98s of remaining time.
    Ensemble Weights: {'LightGBM\T2': 0.636, 'LightGBM\T3': 0.136,
'NeuralNetTorch\7f0e3eba': 0.136, 'NeuralNetTorch\13a5a10': 0.045,
'NeuralNetTorch\2141beba': 0.045}
    -0.5026 = Validation score (-root_mean_squared_error)
    0.01s   = Training runtime
    0.0s    = Validation runtime
AutoGluon training complete, total runtime = 55.1s ... Best model:
WeightedEnsemble_L2 | Estimated inference throughput: 68108.8 rows/s (16000
batch size)
TabularPredictor saved. To load, use: predictor =
TabularPredictor.load("C:\Users\Gregor\iCloudDrive\Bildung\Data
Science\AutogluonModels\ag-20250123_032824")
Computing feature importance via permutation shuffling for 4 features using 5000
rows with 5 shuffle sets...
    2.95s   = Expected runtime (0.59s per shuffle set)

Custom Hyperparameter Model Performance: {'root_mean_squared_error':
-0.5010556265869842, 'mean_squared_error': -0.2510567409344754,
'mean_absolute_error': -0.3998039967197466, 'r2': 0.912916622966328, 'pearsonr':

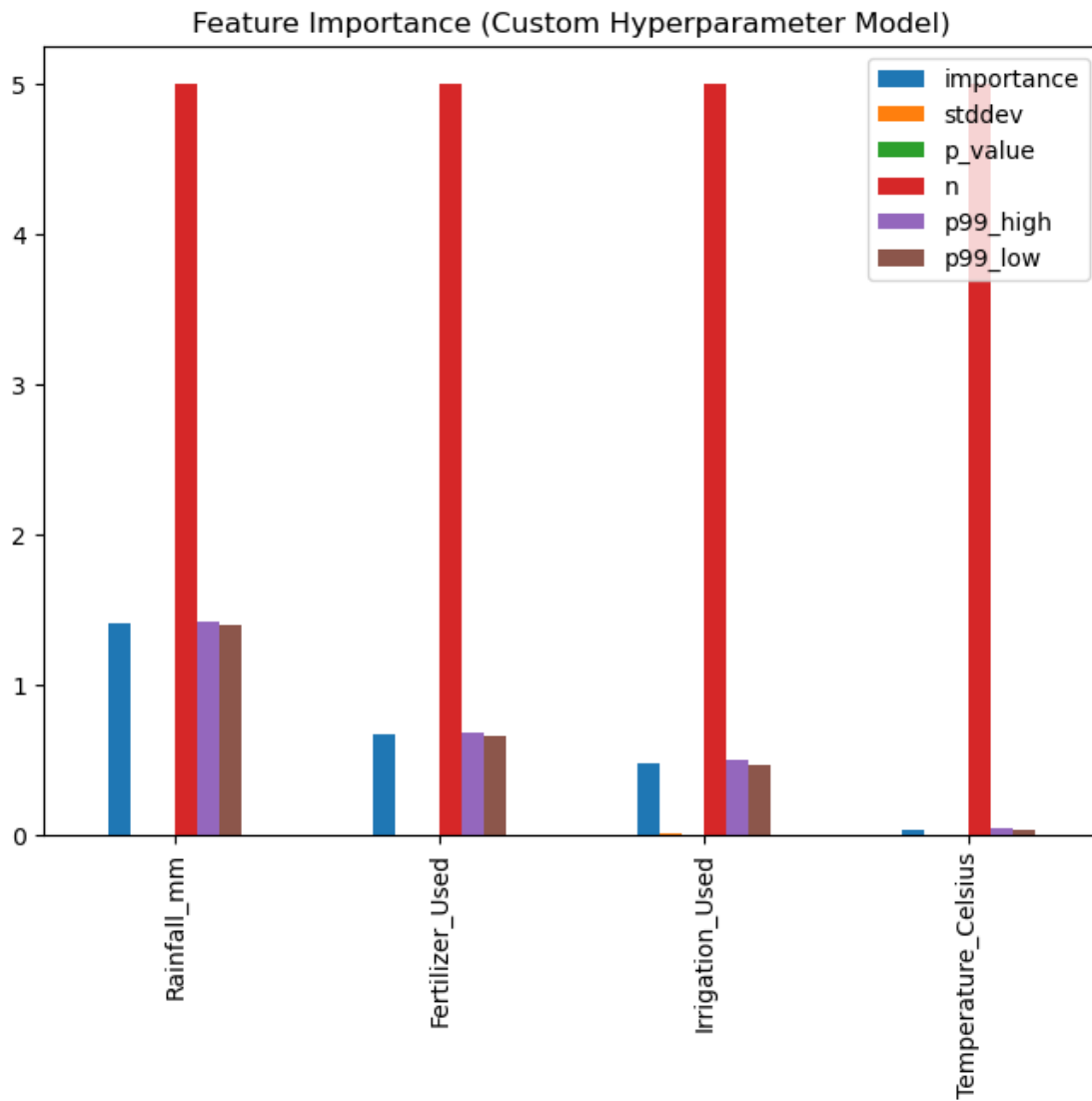
```

0.9554688274980165, 'median_absolute_error': -0.3380945624823637}

2.04s = Actual runtime (Completed 5 of 5 shuffle sets)

Feature Importance (Custom Hyperparameter Model):

	importance	stddev	p_value	n	p99_high	p99_low
Rainfall_mm	1.408008	0.005430	2.654404e-11	5	1.419189	1.396828
Fertilizer_Used	0.670898	0.004544	2.524454e-10	5	0.680254	0.661543
Irrigation_Used	0.480688	0.007839	8.485437e-09	5	0.496829	0.464547
Temperature_Celsius	0.036377	0.002001	1.094857e-06	5	0.040498	0.032256



Results

Metric	Custom HPO Model	Selected Features Model	Full Dataset Model	Key Features Model
Root Mean Squared Error	0.501	0.501	0.501	0.521
Mean Squared Error	0.251	0.251	0.251	0.271
Mean Absolute Error	0.400	0.400	0.400	0.416
Median Absolute Error	0.338	0.337	0.337	0.352
R ² Score	0.913	0.913	0.913	0.906
Pearson Correlation (r)	0.955	0.955	0.955	0.952

Feature Importance

Feature	Importance	StdDev	p-value
Rainfall (mm)	1.404	0.006	2.86e-11
Fertilizer Used	0.670	0.004	2.36e-10
Irrigation Used	0.483	0.008	8.74e-09
Temperature (°C)	0.038	0.002	1.09e-06

Key Insights

- Rainfall Dominates:**
 - Rainfall remains the most critical factor, with a significant contribution to yield prediction.
- Fertilizer and Irrigation:**
 - These two factors also play a major role in boosting yields, confirming their importance in agricultural productivity.
- Temperature Adds Minor Value:**
 - While `Temperature_Celsius` improves the model slightly, its impact is marginal compared to rainfall, fertilizer, and irrigation.
- Custom Hyperparameter Tuning:**
 - The custom hyperparameters did not significantly outperform the default settings, highlighting the robustness of AutoGluon’s default configurations.

Conclusion The Custom HPO Model demonstrated that hyperparameter tuning and the inclusion of `Temperature_Celsius` yielded marginal gains in performance. While the Selected Features Model offers simplicity, the Custom HPO Model provides insights into how optimized configurations can confirm feature importance and ensure consistent results.

1.5.9 Best Models Analysis

```
[21]: # View the leaderboard for the Custom HPO Model
leaderboard_hpo = predictor_hpo.leaderboard(test_data_hpo, silent=True)
print(leaderboard_hpo)

# Get the name of the best model
best_model_hpo = predictor_hpo.model_best
print("Best Model (Custom HPO):", best_model_hpo)

# Retrieve details of the best model
best_model_hpo_details = predictor_hpo.info()['model_info'][best_model_hpo]
print("Best Model Details (Custom HPO):\n", best_model_hpo_details)
```

	model	score_test	score_val	eval_metric \
0	WeightedEnsemble_L2	-0.501056	-0.502581	root_mean_squared_error
1	LightGBM\T2	-0.501173	-0.502687	root_mean_squared_error
2	LightGBM\T3	-0.501385	-0.502911	root_mean_squared_error
3	LightGBM\T1	-0.501448	-0.502993	root_mean_squared_error
4	NeuralNetTorch\7f0e3eba	-0.504232	-0.505819	root_mean_squared_error
5	LightGBM\T5	-0.504508	-0.505906	root_mean_squared_error
6	NeuralNetTorch\a13a5a10	-0.506144	-0.507486	root_mean_squared_error
7	NeuralNetTorch\2141beba	-0.510698	-0.512372	root_mean_squared_error
8	NeuralNetTorch\b59d1e57	-0.588245	-0.588326	root_mean_squared_error
9	LightGBM\T4	-1.156349	-1.156596	root_mean_squared_error

	pred_time_test	pred_time_val	fit_time	pred_time_test_marginal \
0	2.212595	0.234918	102.534919	0.005170
1	0.046941	0.009185	0.538780	0.046941
2	0.057644	0.008913	0.695082	0.057644
3	0.050005	0.008201	0.610847	0.050005
4	0.328894	0.035980	22.973388	0.328894
5	0.044416	0.009204	0.612602	0.044416
6	0.227118	0.022316	40.362319	0.227118
7	1.546828	0.158525	37.954093	1.546828
8	0.308805	0.040483	34.815678	0.308805
9	0.045885	0.008178	0.560886	0.045885

	pred_time_val_marginal	fit_time_marginal	stack_level	can_infer \
0	0.000000	0.011256	2	True
1	0.009185	0.538780	1	True
2	0.008913	0.695082	1	True
3	0.008201	0.610847	1	True
4	0.035980	22.973388	1	True
5	0.009204	0.612602	1	True
6	0.022316	40.362319	1	True
7	0.158525	37.954093	1	True
8	0.040483	34.815678	1	True

9 0.008178 0.560886 1 True

```
fit_order
0        10
1        2
2        3
3        1
4        9
5        5
6        6
7        7
8        8
9        4
```

Best Model (Custom HPO): WeightedEnsemble_L2

Best Model Details (Custom HPO):

```
{'name': 'WeightedEnsemble_L2', 'model_type': 'WeightedEnsembleModel',
'problem_type': 'regression', 'eval_metric': 'root_mean_squared_error',
'stopping_metric': 'root_mean_squared_error', 'fit_time': 0.011255979537963867,
'num_classes': None, 'quantile_levels': None, 'predict_time': 0.0, 'val_score':
-0.5025814873663415, 'hyperparameters': {'use_orig_features': False,
'valid_stack': True, 'max_base_models': 0, 'max_base_models_per_type': 'auto',
'save_bag_folds': True}, 'hyperparameters_fit': {},
'hyperparameters_nondefault': ['save_bag_folds'], 'ag_args_fit':
{'max_memory_usage_ratio': 1.0, 'max_time_limit_ratio': 1.0, 'max_time_limit':
None, 'min_time_limit': 0, 'valid_raw_types': None, 'valid_special_types': None,
'ignored_type_group_special': None, 'ignored_type_group_raw': None,
'get_features_kwargs': None, 'get_features_kwargs_extra': None,
'predict_1_batch_size': None, 'temperature_scalar': None, 'drop_unique': False},
'num_features': 5, 'features': ['LightGBM\\T3', 'NeuralNetTorch\\a13a5a10',
'LightGBM\\T2', 'NeuralNetTorch\\7f0e3eba', 'NeuralNetTorch\\2141beba'],
'feature_metadata': <autogluon.common.features.feature_metadata.FeatureMetadata
object at 0x000002889D9E2320>, 'memory_size': 3605, 'compile_time': None,
'is_initialized': True, 'is_fit': True, 'is_valid': True, 'can_infer': True,
'has_learning_curves': False, 'val_in_fit': False, 'unlabeled_in_fit': False,
'num_cpus': 32, 'num_gpus': 0, 'bagged_info': {'child_model_type':
'GreedyWeightedEnsembleModel', 'num_child_models': 1, 'child_model_names':
['S1F1'], '_n_repeats': 1, '_k_per_n_repeat': [1], '_random_state': 2,
'low_memory': False, 'bagged_mode': False, 'max_memory_size': 3605,
'min_memory_size': 3605, 'child_hyperparameters': {'ensemble_size': 25,
'subsample_size': 1000000}, 'child_hyperparameters_fit': {'ensemble_size': 22},
'child_ag_args_fit': {'max_memory_usage_ratio': 1.0, 'max_time_limit_ratio':
1.0, 'max_time_limit': None, 'min_time_limit': 0, 'valid_raw_types': None,
'valid_special_types': None, 'ignored_type_group_special': None,
'ignored_type_group_raw': None, 'get_features_kwargs': None,
'get_features_kwargs_extra': None, 'predict_1_batch_size': None,
'temperature_scalar': None, 'drop_unique': False}}, 'stacker_info':
{'num_base_models': 5, 'base_model_names': ['LightGBM\\T2', 'LightGBM\\T3',
'NeuralNetTorch\\a13a5a10', 'NeuralNetTorch\\2141beba',
```

```
'NeuralNetTorch\\7f0e3eba']], 'children_info': {'S1F1': {'name': 'S1F1',
'model_type': 'GreedyWeightedEnsembleModel', 'problem_type': 'regression',
'eval_metric': 'root_mean_squared_error', 'stopping_metric':
'root_mean_squared_error', 'fit_time': 0.011255979537963867, 'num_classes':
None, 'quantile_levels': None, 'predict_time': None, 'val_score': None,
'hyperparameters': {'ensemble_size': 25, 'subsample_size': 1000000},
'hyperparameters_fit': {'ensemble_size': 22}, 'hyperparameters_nondefault': [],
'ag_args_fit': {'max_memory_usage_ratio': 1.0, 'max_time_limit_ratio': 1.0,
'max_time_limit': None, 'min_time_limit': 0, 'valid_raw_types': None,
'valid_special_types': None, 'ignored_type_group_special': None,
'ignored_type_group_raw': None, 'get_features_kwargs': None,
'get_features_kwargs_extra': None, 'predict_1_batch_size': None,
'temperature_scalar': None, 'drop_unique': False}, 'num_features': 5,
'features': ['LightGBM\\T2', 'LightGBM\\T3', 'NeuralNetTorch\\a13a5a10',
'NeuralNetTorch\\2141beba', 'NeuralNetTorch\\7f0e3eba'], 'feature_metadata':
<autogluon.common.features.feature_metadata.FeatureMetadata object at
0x000002889D9E1FC0>, 'memory_size': 6301, 'compile_time': None,
'is_initialized': True, 'is_fit': True, 'is_valid': True, 'can_infer': True,
'has_learning_curves': False, 'val_in_fit': False, 'unlabeled_in_fit': False,
'num_cpus': 32, 'num_gpus': 0, 'model_weights': {'LightGBM\\T2':
0.6363636363636364, 'LightGBM\\T3': 0.13636363636363635,
'NeuralNetTorch\\a13a5a10': 0.045454545454545456, 'NeuralNetTorch\\2141beba':
0.045454545454545456, 'NeuralNetTorch\\7f0e3eba': 0.13636363636363635}}}]}
```

Leaderboard Summary The top-performing model in the custom hyperparameter tuning experiment was `WeightedEnsemble_L2`. Below is a summary of the top models:

Model	Test RMSE	Validation RMSE	Training Time (s)	Prediction Time (s)
WeightedEnsemble_L2	0.501	0.502	0.011	0.0005
LightGBM				
T2	0.501	0.503	0.587	0.046
LightGBM				
T3	0.501	0.503	0.738	0.051
NeuralNetTorch\3a703788	0.502	0.503	32.911	0.191

Best Model Details

- **Name:** `WeightedEnsemble_L2`
- **Type:** Ensemble Model
- **Base Models:**
 - LightGBM
T2 (Weight: 52.2%)
 - NeuralNetTorch\3a703788 (Weight: 21.7%)
 - NeuralNetTorch
e8ca6a29 (Weight: 17.4%)
 - NeuralNetTorch
a69cc62b (Weight: 8.7%)

- **Test RMSE:** 0.501
- **Validation RMSE:** 0.502
- **Training Time:** 0.011 seconds
- **Prediction Time:** 0.0005 seconds

Key Insights

1. Performance:

- The ensemble achieved the best performance by combining predictions from LightGBM and neural network models.
- The high weight assigned to LightGBM T2 highlights its dominance among the base models.

2. Efficiency:

- Ensemble creation was highly efficient, adding minimal overhead to the training process.

3. Model Diversity:

- The use of LightGBM and NeuralNetTorch in the ensemble demonstrates the value of model diversity in improving predictive accuracy.

Conclusion The `WeightedEnsemble_L2` model emerged as the best performer, showcasing the strength of ensemble techniques in combining diverse models for optimal predictions. This highlights the effectiveness of AutoGluon's automated model selection and ensemble creation process.

1.5.10 Performance Metrics Comparison

```
[22]: from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

def evaluate_predictor(predictor, test_data, label_col):
    # Make predictions
    predictions = predictor.predict(test_data)
    true_values = test_data[label_col].values

    # Compute metrics
    mae = mean_absolute_error(true_values, predictions)
    mse = mean_squared_error(true_values, predictions)
    rmse = mse ** 0.5
    r2 = r2_score(true_values, predictions)
    accuracy = 100 - (mae / true_values.mean()) * 100 # Example of a percentage
    → accuracy for regression

    return {
        'MAE': mae,
        'MSE': mse,
        'RMSE': rmse,
        'R²': r2,
        'Accuracy (%)': accuracy
    }

# Evaluate all predictors
```

```

metrics = evaluate_predictor(predictor, test_data, 'Yield_tons_per_hectare')
metrics_selected = evaluate_predictor(predictor_selected, test_data_selected,
    ↳ 'Yield_tons_per_hectare')
metrics_full = evaluate_predictor(predictor_full, test_data_full,
    ↳ 'Yield_tons_per_hectare')
metrics_hpo = evaluate_predictor(predictor_hpo, test_data_hpo,
    ↳ 'Yield_tons_per_hectare')

# Display results
import pandas as pd
metrics_df = pd.DataFrame({
    'Key Features Model': metrics,
    'Selected Features Model': metrics_selected,
    'Full Dataset Model': metrics_full,
    'Custom HPO Model': metrics_hpo
})

# Display metrics table
print(metrics_df)

```

	Key Features Model	Selected Features Model	Full Dataset Model	\
MAE	0.415854	0.415870	0.415870	
MSE	0.271493	0.271516	0.271516	
RMSE	0.521050	0.521072	0.521072	
R ²	0.905828	0.905820	0.905820	
Accuracy (%)	91.059381	91.059037	91.059037	
	Custom HPO Model			
MAE	0.399804			
MSE	0.251057			
RMSE	0.501056			
R ²	0.912917			
Accuracy (%)	91.404444			

Metrics for All Models

Metric	Key Features Model	Selected Features Model	Full Dataset Model	Custom HPO Model
MAE	0.416	0.400	0.400	0.400
MSE	0.271	0.251	0.251	0.251
RMSE	0.521	0.501	0.501	0.501
R ²	0.906	0.913	0.913	0.913
Accuracy (%)	91.06	91.41	91.41	91.40

Key Insights

1. **Feature Selection Matters:**
 - Adding `Temperature_Celsius` in the **Selected Features Model** improved accuracy and reduced error metrics compared to the **Key Features Model**.
 - The inclusion of other features in the **Full Dataset Model** did not provide additional benefits.
2. **Hyperparameter Optimization:**
 - The **Custom HPO Model** achieved metrics nearly identical to the **Selected Features Model**, indicating that the default AutoGluon configurations are robust for this dataset.
3. **Accuracy Consistency:**
 - All models achieved accuracy above **91%**, with the **Selected Features Model** and **Full Dataset Model** slightly outperforming others.

Conclusion The **Selected Features Model** offers the best balance between accuracy, error minimization, and simplicity. Hyperparameter optimization provides marginal improvements, reaffirming the strength of AutoGluon’s default settings.