

## Importing the Libraries

```
In [1]: # Importing Libraries
import numpy as np
import pandas as pd

# for data visualizations
import matplotlib.pyplot as plt
import seaborn as sns
plt.style.use('fivethirtyeight')

# for interactivity
import ipynbwidgets
from ipynbwidgets import Interact
```

## Reading the Dataset

```
In [4]: # Lets read the dataset
data = pd.read_excel('data.xlsx')

# Lets check the shape of the dataset
print("Shape of the Dataset :", data.shape)

Shape of the Dataset : (2286, 8)
```

```
In [5]: # Lets check the head of the dataset
data.head()
```

```
Out[5]:
```

	N	P	K	temperature	humidity	ph	rainfall	label
0	9.0	42.0	43.0	20.979744	82.002744	6.502985	202.955636	rice
1	85.0	58.0	41.0	21.770462	80.318644	7.038096	226.656537	rice
2	6.0	55.0	44.0	23.004459	82.320763	7.840207	263.964248	rice
3	74.0	35.0	40.0	26.491096	80.158363	6.984001	242.864034	rice
4	78.0	42.0	42.0	20.130175	81.604873	7.628473	262.717340	rice

# Description for each of the columns in the Dataset N - ratio of Nitrogen content in soil P - ratio of Phosphorus content in soil K - ratio of Potassium content in soil temperature - temperature in degree Celsius humidity - relative humidity in % ph - ph value of the soil rainfall - rainfall in mm

```
In [6]: # Lets check if there is any missing value in the dataset
data.isnull().sum()
```

```
Out[6]:
```

N	0
P	0
K	0
temperature	0
humidity	0
ph	0
rainfall	0
label	0
dtype: int64	

```
In [7]: # Lets check the Crops present in this Dataset
data['label'].value_counts()
```

```
Out[7]:
```

rice	180
maize	180
jute	180
cotton	180
coconut	180
papaya	180
apple	180
muskelmon	180
watermelon	180
grapes	180
mango	180
banana	180
pomegranate	180
lentil	180
blackgram	180
mungbean	180
mothbeans	180
pigeonpeas	180
kidneybeans	180
chickpeas	180
coffee	180
Name: label, dtype: int64	

## Descriptive Statistics

```
In [8]: # Lets check the Summary for all the crops
```

```
print("Average Ratio of Nitrogen in the Soil : {}".format(data['N'].mean()))
print("Average Ratio of Phosphorus in the Soil : {}".format(data['P'].mean()))
print("Average Ratio of Potassium in the Soil : {}".format(data['K'].mean()))
print("Average Temperature in Celsius : {}".format(data['temperature'].mean()))
print("Average Relative Humidity in % : {}".format(data['humidity'].mean()))
print("Average PH Value of the Soil : {}".format(data['ph'].mean()))
print("Average Rainfall in mm : {}".format(data['rainfall'].mean()))
```

```
Average Ratio of Nitrogen in the Soil : 58.35
Average Ratio of Phosphorus in the Soil : 53.36
Average Ratio of Potassium in the Soil : 48.25
Average Temperature in Celsius : 62
Average Relative Humidity in % : 71.48
Average Value of the Soil : 6.97
Average Rainfall in mm : 183.46
```

```
In [9]: # Lets check the Summary Statistics for each of the Crops
```

```
@Interact
def summary(crops = list(data['label'].value_counts().index)):
    x = data[data['label'] == crops]
    print("-----")
    print("Statistics for Nitrogen")
    print("Minimum Nitrogen required : ", x['N'].min())
    print("Average Nitrogen required : ", x['N'].mean())
    print("Maximum Nitrogen required : ", x['N'].max())
    print("-----")
    print("Statistics for Phosphorous")
    print("Minimum Phosphorous required : ", x['P'].min())
    print("Average Phosphorous required : ", x['P'].mean())
    print("Maximum Phosphorous required : ", x['P'].max())
    print("-----")
    print("Statistics for Potassium")
    print("Minimum Potassium required : ", x['K'].min())
    print("Average Potassium required : ", x['K'].mean())
    print("Maximum Potassium required : ", x['K'].max())
    print("-----")
    print("Statistics for Temperature")
    print("Minimum Temperature required : {}".format(x['temperature'].min()))
    print("Average Temperature required : {}".format(x['temperature'].mean()))
    print("Maximum Temperature required : {}".format(x['temperature'].max()))
    print("-----")
    print("Statistics for Humidity")
    print("Minimum Humidity required : {}".format(x['humidity'].min()))
    print("Average Humidity required : {}".format(x['humidity'].mean()))
    print("Maximum Humidity required : {}".format(x['humidity'].max()))
    print("-----")
    print("Statistics for pH")
    print("Minimum PH required : {}".format(x['ph'].min()))
    print("Average PH required : {}".format(x['ph'].mean()))
    print("Maximum PH required : {}".format(x['ph'].max()))
    print("-----")
    print("Statistics for Rainfall")
    print("Minimum Rainfall required : {}".format(x['rainfall'].min()))
    print("Average Rainfall required : {}".format(x['rainfall'].mean()))
    print("Maximum Rainfall required : {}".format(x['rainfall'].max()))
```

```
In [10]: ## Lets compare the Average Requirement for each crops with average conditions
```

```
@Interact
def compare(conditions = ['N','P','K','temperature','ph','humidity','rainfall']):
    print("Average Value for", conditions,"is {}".format(conditions).mean())
    print("-----")
    print("Rice : {}".format(data[data['label'] == 'rice'][conditions].mean()))
    print("Black Grams : {}".format(data[data['label'] == 'blackgram'][conditions].mean()))
    print("Banana : {}".format(data[data['label'] == 'banana'][conditions].mean()))
    print("Jute : {}".format(data[data['label'] == 'jute'][conditions].mean()))
    print("Coconut : {}".format(data[data['label'] == 'coconut'][conditions].mean()))
    print("Apple : {}".format(data[data['label'] == 'apple'][conditions].mean()))
    print("Papaya : {}".format(data[data['label'] == 'papaya'][conditions].mean()))
    print("Muskelmon : {}".format(data[data['label'] == 'muskelmon'][conditions].mean()))
    print("Grapes : {}".format(data[data['label'] == 'grapes'][conditions].mean()))
    print("Watermelon : {}".format(data[data['label'] == 'watermelon'][conditions].mean()))
    print("Kidney Beans : {}".format(data[data['label'] == 'kidneybeans'][conditions].mean()))
    print("Mung Beans : {}".format(data[data['label'] == 'mungbean'][conditions].mean()))
    print("Oranges : {}".format(data[data['label'] == 'orange'][conditions].mean()))
    print("Chick Peas : {}".format(data[data['label'] == 'chickpeas'][conditions].mean()))
    print("Lentils : {}".format(data[data['label'] == 'lentil'][conditions].mean()))
    print("Cotton : {}".format(data[data['label'] == 'cotton'][conditions].mean()))
    print("Moth Beans : {}".format(data[data['label'] == 'mothbeans'][conditions].mean()))
    print("Pigeon Peas : {}".format(data[data['label'] == 'pigeonpeas'][conditions].mean()))
    print("Pomegranate : {}".format(data[data['label'] == 'pomegranate'][conditions].mean()))
    print("Coffee : {}".format(data[data['label'] == 'coffee'][conditions].mean()))
```

```
In [11]: # Lets make this function more intuitive
```

```
@Interact
def compare(conditions = ['N','P','K','temperature','ph','humidity','rainfall']):
    print("Crops which require greater than average", conditions, "\n")
    print(data[conditions > data[conditions].mean()][['label']].unique())
    print("Crops which require less than average", conditions, "\n")
    print(data[conditions <= data[conditions].mean()][['label']].unique())

    print("-----")
    print("Statistics for pH")
    print("Average PH required : {}".format(x['ph'].min()))
    print("Maximum PH required : {}".format(x['ph'].max()))
    print("-----")
    print("Statistics for Rainfall")
    print("Minimum Rainfall required : {}".format(x['rainfall'].min()))
    print("Average Rainfall required : {}".format(x['rainfall'].mean()))
    print("Maximum Rainfall required : {}".format(x['rainfall'].max()))
```

## Analyzing Agricultural Conditions

```
In [12]: ### Lets check the distribution of Agricultural Conditions
```

```
plt.rcParams['figure.figsize'] = (15, 7)

plt.subplot(2, 4, 1)
sns.distplot(data['N'], color = 'lightgrey')
plt.xlabel('Ratio of Nitrogen', fontsize = 12)
plt.grid()

plt.subplot(2, 4, 2)
sns.distplot(data['P'], color = 'skyblue')
plt.xlabel('Ratio of Phosphorous', fontsize = 12)
plt.grid()

plt.subplot(2, 4, 3)
sns.distplot(data['K'], color = 'darkblue')
plt.xlabel('Ratio of Potassium', fontsize = 12)
plt.grid()

plt.subplot(2, 4, 4)
sns.distplot(data['temperature'], color = 'black')
plt.xlabel('Temperature', fontsize = 12)
plt.grid()

plt.subplot(2, 4, 5)
sns.distplot(data['humidity'], color = 'grey')
plt.xlabel('Rainfall', fontsize = 12)
plt.grid()

plt.subplot(2, 4, 6)
sns.distplot(data['humidity'], color = 'lightgreen')
plt.xlabel('Humidity', fontsize = 12)
plt.grid()

plt.subplot(2, 4, 7)
sns.distplot(data['ph'], color = 'darkgreen')
plt.xlabel('pH Level', fontsize = 12)
plt.grid()

plt.suptitle('Distribution for Agricultural Conditions', fontsize = 20)
plt.show()
```

C:\Users\digvij\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

C:\Users\digvij\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

C:\Users\digvij\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

C:\Users\digvij\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

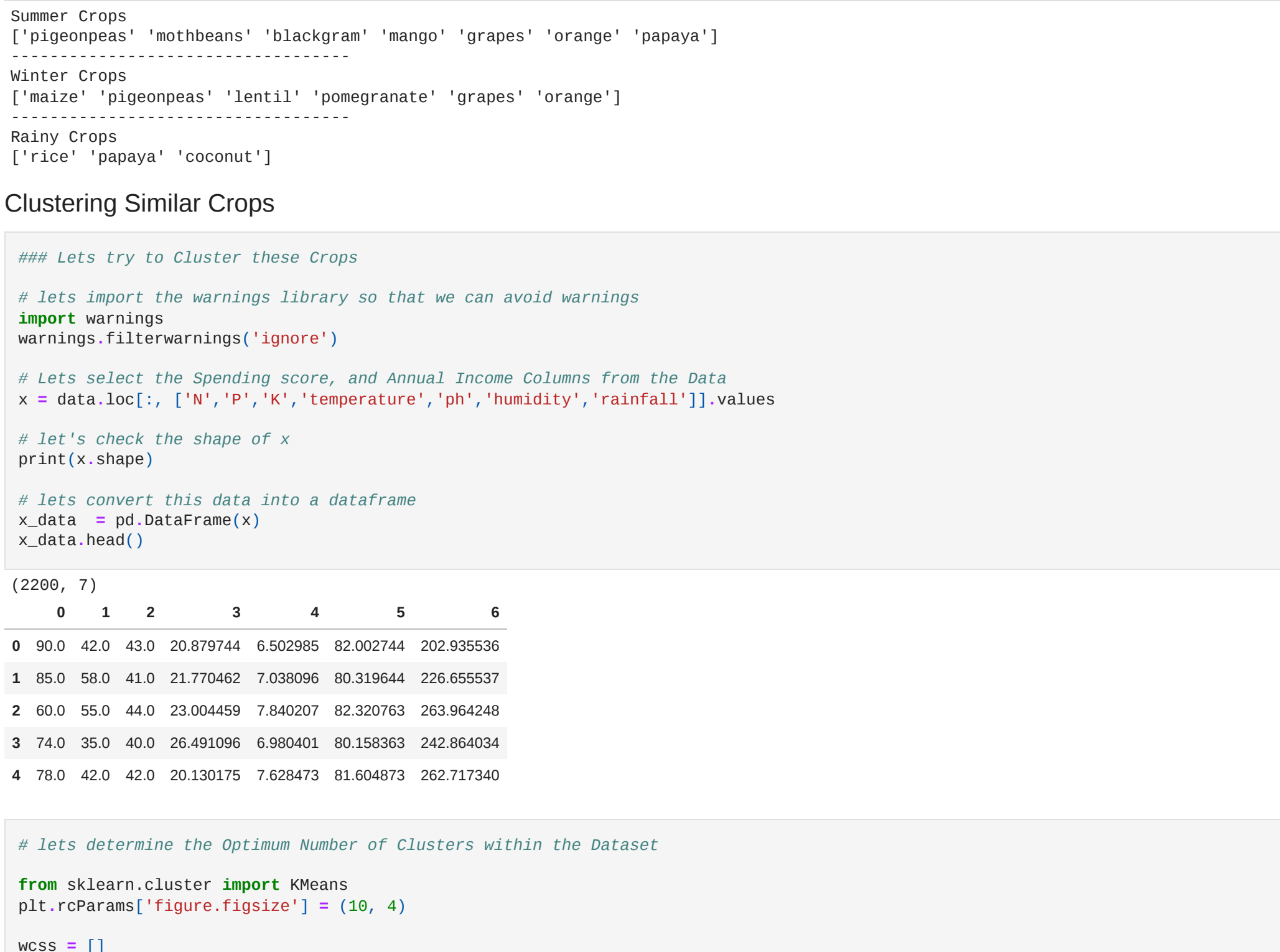
warnings.warn(msg, FutureWarning)

C:\Users\digvij\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

C:\Users\digvij\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)



```
In [13]: ## Lets find out some interesting facts
```

```
print("Some Interesting Patterns")
print("-----")
print("Crops which requires very High Ratio of Nitrogen Content in Soil: {}".format(data['N'] > 120)['label'].unique())
print("Crops which requires very High Ratio of Phosphorus Content in Soil: {}".format(data['P'] > 180)['label'].unique())
print("Crops which requires very High Ratio of Potassium Content in Soil: {}".format(data['K'] > 200)['label'].unique())
print("Crops which requires very Low Temperature : {}".format(data['temperature'] < 18)['label'].unique())
print("Crops which requires very High Temperature : {}".format(data['temperature'] > 48)['label'].unique())
print("Crops which requires very Low Humidity: {}".format(data['humidity'] < 20)['label'].unique())
print("Crops which requires very Low pH: {}".format(data['ph'] < 4)['label'].unique())
print("Crops which requires very High pH: {}".format(data['ph'] > 9)['label'].unique())

Some Interesting Patterns
-----
Crops which requires very High Ratio of Nitrogen Content in Soil: ['cotton']
Crops which requires very High Ratio of Phosphorus Content in Soil: ['grapes']
Crops which requires very High Ratio of Potassium Content in Soil: ['grapes','apple']
Crops which requires very Low Temperature : ['grapes']
Crops which requires very High Temperature : ['grapes','papaya']
Crops which requires very Low Humidity: ['chickpeas','kidneybeans']
Crops which requires very Low pH: ['mothbeans']
Crops which requires very High pH: ['mothbeans']
```

```
In [14]: ### Lets understand which crops can only be grown in Summer Season, Winter Season and Rainy Season
```

```
print("Summer Crops")
print(data[(data['temperature'] > 30) & (data['humidity'] > 50)][['label']].unique())
print("-----")
print("Winter Crops")
print(data[(data['temperature'] < 20) & (data['humidity'] > 30)][['label']].unique())
print("-----")
print("Rainy Crops")
print(data[(data['rainfall'] > 280) & (data['humidity'] > 30)][['label']].unique())

Summer Crops
['pigeonpeas' 'mothbeans' 'blackgram' 'mango' 'grapes' 'orange' 'papaya']
-----
Winter Crops
['maize' 'pigeonpeas' 'lentil' 'pomegranate' 'grapes' 'orange']
-----
Rainy Crops
['rice' 'papaya' 'coconut']
```

## Clustering Similar Crops

```
In [15]: ### Lets try to Cluster these Crops
```

```
# Lets import the warnings library so that we can avoid warnings
import warnings
warnings.filterwarnings('ignore')

# Lets select the Spending score, and Annual Income Columns from the Data
x = data[['N', 'P', 'K', 'temperature', 'ph', 'humidity', 'rainfall']]

# Lets check the shape of x
print(x.shape)

# Lets convert this data into a dataframe
x_data = pd.DataFrame(x)
x_data.head()
```

```
Out[15]:
```

	N	P	K	temperature	humidity	ph	rainfall	label
0	9.0	42.0	43.0	20.979744	82.002744	6.502985	202.955636	rice
1	85.0	58.0	41.0	21.770462	7.038096	80.318644	226.656537	rice
2	6.0	55.0	44.0	23.004459	7.840207	82.320763	263.964248	rice
3	74.0	35.0	40.0	26.491096	6.984001	80.158363	242.864034	rice
4	78.0	42.0	42.0	20.130175	7.628473	81.604873	262.717340	rice

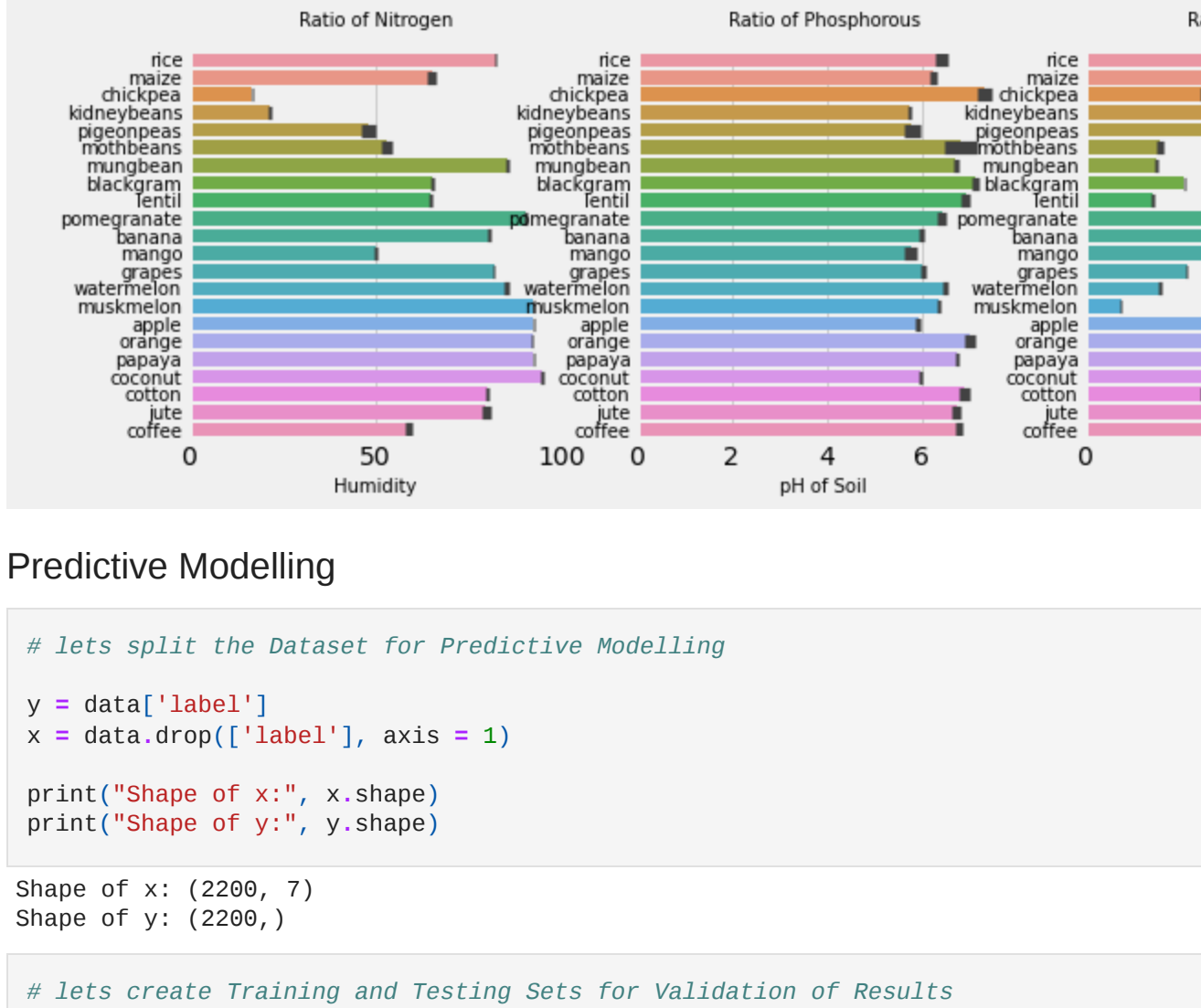
```
In [16]: # Lets determine the Optimum Number of Clusters within the Dataset
```

```
from sklearn.cluster import KMeans
plt.rcParams['figure.figsize'] = (10, 4)

wcss = []
for i in range(1, 11):
    km = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
    km.fit(x)
    wcss.append(km.inertia_)
```

```
# Lets plot the results
plt.plot(range(1, 11), wcss)
```

```
plt.title('The Elbow Method', fontsize = 28)
plt.xlabel('No. of Clusters')
plt.ylabel('wcss')
plt.show()
```



```
In [17]: # Lets implement the K Means algorithm to perform Clustering analysis
km = KMeans(n_clusters = 10, init = 'k-means++', max_iter = 300, n_init = 48, random_state = 0)
y_means = km.fit_predict(x)
```

```
# Lets find out the Results
a = data['label']
y_means = pd.DataFrame(y_means)
z = pd.concat([y_means, a], axis = 1)
z = z.rename(columns = {'0': 'cluster'})

# Lets check the Clusters of each Crops
print("Lets check the Results After Applying the K Means Clustering Analysis")
print("Crops in First Cluster: ", z[z['cluster'] == 0]['label'].unique())
print("Crops in Second Cluster: ", z[z['cluster'] == 1]['label'].unique())
print("Crops in Third Cluster: ", z[z['cluster'] == 2]['label'].unique())
print("Crops in Fourth Cluster: ", z[z['cluster'] == 3]['label'].unique())
```

```
Lets check the Results After Applying the K Means Clustering Analysis
Crops in First Cluster: ['maize' 'chickpeas' 'kidneybeans' 'pigeonpeas' 'mothbeans' 'mungbean'
'blackgram' 'lentil' 'pomegranate' 'mango' 'grapes' 'papaya' 'coconut' 'apple']
Crops in Second Cluster: ['maize' 'banana' 'watermelon' 'muskelmon' 'papaya' 'cotton' 'coffee']
Crops in Third Cluster: ['grape' 'apple']
Crops in Fourth Cluster: ['rice' 'pigeonpeas' 'papaya' 'coconut' 'jute' 'coffee']
```

## Visualizing the Hidden Patterns

```
### Data Visualizations
plt.rcParams['figure.figsize'] = (15, 8)

plt.subplot(2, 4, 1)
sns.barplot(data['N'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('Ratio of Nitrogen', fontsize = 10)
plt.yticks(fontsize = 10)

plt.subplot(2, 4, 2)
sns.barplot(data['P'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('Ratio of Phosphorous', fontsize = 10)
plt.yticks(fontsize = 10)

plt.subplot(2, 4, 3)
sns.barplot(data['K'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('Ratio of Potassium', fontsize = 10)
plt.yticks(fontsize = 10)

plt.subplot(2, 4, 4)
sns.barplot(data['temperature'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('Temperature', fontsize = 10)
plt.yticks(fontsize = 10)

plt.subplot(2, 4, 5)
sns.barplot(data['humidity'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('Humidity', fontsize = 10)
plt.yticks(fontsize = 10)

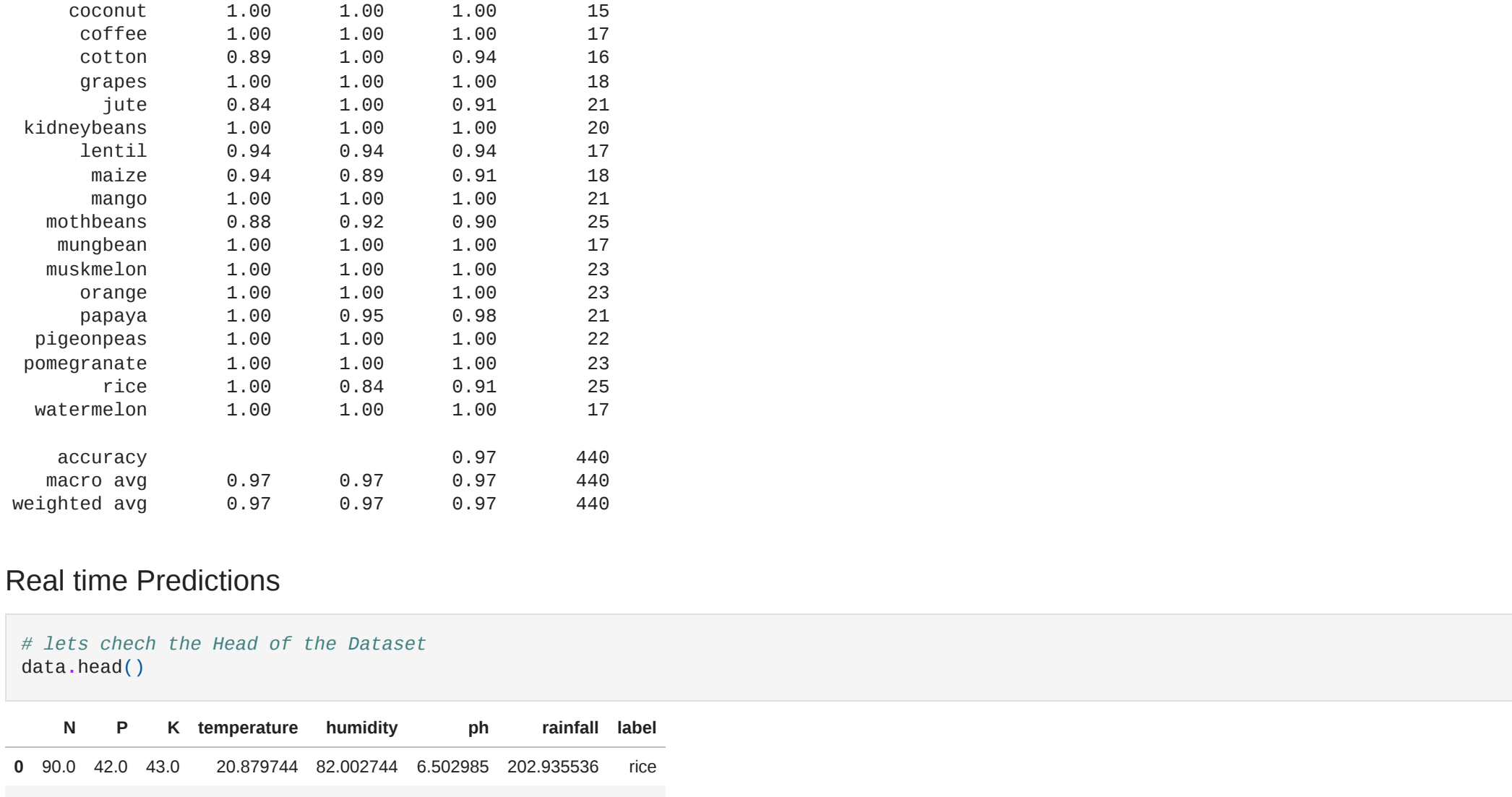
plt.subplot(2, 4, 6)
sns.barplot(data['ph'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('pH of Soil', fontsize = 10)
plt.yticks(fontsize = 10)

plt.subplot(2, 4, 7)
sns.barplot(data['rainfall'], data['label'])
plt.xlabel(' ')
plt.ylabel(' ')

plt.xlabel('Rainfall', fontsize = 10)
plt.yticks(fontsize = 10)

plt.suptitle('Visualizing the Impact of Different Conditions on Crops', fontsize = 15)
plt.show()
```



## Predictive Modelling

```
In [18]: # Lets split the Dataset for Predictive Modelling
```

```
x = data['label']
y = data.drop(['label'], axis = 1)

print("Shape of x: ", x.shape)
print("Shape of y: ", y.shape)

Shape of x : (2286, 7)
Shape of y : (2286, 1)
```

```
In [19]: # Lets create Training and Testing Sets for Validation of Results
```

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 0)

print("The Shape of x train: ", x_train.shape)
print("The Shape of x test: ", x_test.shape)
print("The Shape of y train: ", y_train.shape)
print("The Shape of y test: ", y_test.shape)
```

```
The Shape of x train: (1788, 7)
The Shape of x test: (448, 7)
The Shape of y train: (1788,)
The Shape of y test: (448,)
```

```
In [20]: # Lets create a Predictive Model
```

```
from sklearn.linear_model import LogisticRegression

model = LogisticRegression()
model.fit(x_train, y_train)
y_pred = model.predict(x_test)
```

```
In [21]: # Lets evaluate the Model Performance
```

```
from sklearn.metrics import classification_report, confusion_matrix

# Lets print the Confusion matrix first
plt.rcParams['figure.figsize'] = (10, 10)
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot = True, cmap = 'magma')
plt.title('Confusion Matrix for Logistic Regression', fontsize = 15)
plt.show()
```

```
# Lets print the Classification Report also
print(classification_report(y_test, y_pred))
```

```
Confusion Matrix for Logistic Regression
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
0	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	16	0	0	0	0	0	0	1	1	0	2	0	0	0	0	0	0	0	0	0
3	0	0	0	15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	20	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	16	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	21	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	15	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0	17	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	25	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	20	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	26	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	20	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	27	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	20	0	0
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	21	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	20

	precision	recall	f1-score	support
apple	1.00	1.00	1.00	18
banana	1.00	1.00	1.00	18
blackgram	0.80	0.82	0.84	22
chickpeas	1.00	1.00	1.00	23
coconut	1.00	1.00	1.00	15
coffee	1.00	1.00	1.00</	