```python
In [1]:   # import the essential libraries
          import numpy as np
          import pandas as pd
          import matplotlib.pyplot as plt
          import seaborn as sns
          import warnings
          warnings.filterwarnings('ignore')
```

```python
In [2]:   # load the dataset
          df=pd.read_csv('UM/Crop Production data.csv')
          df
```

Out[2]:

| | State_Name | District_Name | Crop_Year | Season | Crop | Area | Production |
|---|---|---|---|---|---|---|---|
| 0 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Arecanut | 1254.0 | 2000.0 |
| 1 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Other Kharif pulses | 2.0 | 1.0 |
| 2 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Rice | 102.0 | 321.0 |
| 3 | Andaman and Nicobar Islands | NICOBARS | 2000 | Whole Year | Banana | 176.0 | 641.0 |
| 4 | Andaman and Nicobar Islands | NICOBARS | 2000 | Whole Year | Cashewnut | 720.0 | 165.0 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 246086 | West Bengal | PURULIA | 2014 | Summer | Rice | 306.0 | 801.0 |
| 246087 | West Bengal | PURULIA | 2014 | Summer | Sesamum | 627.0 | 463.0 |
| 246088 | West Bengal | PURULIA | 2014 | Whole Year | Sugarcane | 324.0 | 16250.0 |
| 246089 | West Bengal | PURULIA | 2014 | Winter | Rice | 279151.0 | 597899.0 |
| 246090 | West Bengal | PURULIA | 2014 | Winter | Sesamum | 175.0 | 88.0 |

246091 rows × 7 columns

```python
In [3]:   # check the brief info of the dataset
          df.head()
```

Out[3]:

| | State_Name | District_Name | Crop_Year | Season | Crop | Area | Production |
|---|---|---|---|---|---|---|---|
| 0 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Arecanut | 1254.0 | 2000.0 |
| 1 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Other Kharif pulses | 2.0 | 1.0 |
| 2 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Rice | 102.0 | 321.0 |
| 3 | Andaman and Nicobar Islands | NICOBARS | 2000 | Whole Year | Banana | 176.0 | 641.0 |
| 4 | Andaman and Nicobar Islands | NICOBARS | 2000 | Whole Year | Cashewnut | 720.0 | 165.0 |

```python
In [4]:   # shape of the dataset
          df.shape
```

Out[4]:   (246091, 7)

```python
In [7]:   # index of the datset
          df.index
```

Out[7]:   RangeIndex(start=0, stop=246091, step=1)

```python
In [6]:   # columns of the dataset
          df.columns
```

Out[6]:   Index(['State_Name', 'District_Name', 'Crop_Year', 'Season', 'Crop', 'Area',
                 'Production'],
                dtype='object')

```python
In [8]:   # check the basic info about the dataset
          df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 246091 entries, 0 to 246090
Data columns (total 7 columns):
 #   Column         Non-Null Count   Dtype
---  ------         --------------   -----
 0   State_Name     246091 non-null  object
 1   District_Name  246091 non-null  object
 2   Crop_Year      246091 non-null  int64
 3   Season         246091 non-null  object
 4   Crop           246091 non-null  object
 5   Area           246091 non-null  float64
 6   Production     242361 non-null  float64
dtypes: float64(2), int64(1), object(4)
memory usage: 13.1+ MB
```

```python
In [9]:   # separete the list of features based on their data types
```

```
continuous_features=[]
categorical_features=[]
either_continuous_or_discrete_count=[]
for features in df.columns:
    if df[features].dtypes=='float64':
        continuous_features.append(features)
    elif df[features].dtypes=='object':
        categorical_features.append(features)
    else:
        either_continuous_or_discrete_count.append(features)
print('continuous features:',continuous_features)
print('categorical features:',categorical_features)
print('continuous or discrete count:',either_continuous_or_discrete_count)
```

```
continuous features: ['Area', 'Production']
categorical features: ['State_Name', 'District_Name', 'Season', 'Crop']
continuous or discrete count: ['Crop_Year']
```

In [10]:
```python
#  check is there any null values
df.isnull().sum()
```

Out[10]:
```
State_Name         0
District_Name      0
Crop_Year          0
Season             0
Crop               0
Area               0
Production      3730
dtype: int64
```

In [11]:
```python
#  check the duplicated records
df.duplicated().sum()
```

Out[11]: 0

In [12]:
```python
# It is a huse dataset so remove the null values
df.dropna(inplace=True)
```

In [13]:
```python
# after removing the null values
df.isnull().sum()
```

Out[13]:
```
State_Name       0
District_Name    0
Crop_Year        0
Season           0
Crop             0
Area             0
Production       0
dtype: int64
```

In [14]:
```python
# after removing the null values shape of the dataset
df.shape
```

Out[14]: (242361, 7)

In [15]:
```python
# value counts
df['State_Name'].value_counts()  # 33 states
```

```
Uttar Pradesh                     33189
Madhya Pradesh                    22604
Karnataka                         21079
Bihar                             18874
Assam                             14622
Odisha                            13524
Tamil Nadu                        13266
Maharashtra                       12496
Rajasthan                         12066
Chhattisgarh                      10368
West Bengal                        9597
Andhra Pradesh                     9561
Gujarat                            8365
Telangana                          5591
Uttarakhand                        4825
Haryana                            4540
Kerala                             4003
Nagaland                           3904
Punjab                             3143
Meghalaya                          2867
Arunachal Pradesh                  2545
Himachal Pradesh                   2456
Jammu and Kashmir                  1632
Tripura                            1412
Manipur                            1266
Jharkhand                          1266
Mizoram                             954
Puducherry                          872
Sikkim                              714
Dadra and Nagar Haveli              263
Goa                                 207
Andaman and Nicobar Islands         201
Chandigarh                           89
Name: State_Name, dtype: int64
```

In [16]:
```python
# len of unique values
df['District_Name'].nunique()
```

Out[16]: 646

In [22]:
```python
# crop year unique values
df['Crop_Year'].value_counts()
```

Out[22]:
```
2003    17139
2002    16536
2007    14269
2008    14230
2006    13976
2004    13858
2010    13793
2011    13791
2009    13767
2000    13553
2005    13519
2013    13475
2001    13293
2012    13184
1999    12441
1998    11262
2014    10815
1997     8899
2015      561
Name: Crop_Year, dtype: int64
```

In [23]:
```python
# season value counts
df['Season'].value_counts()
```

Out[23]:
```
Kharif        94283
Rabi          66160
Whole Year    56127
Summer        14811
Winter         6050
Autumn         4930
Name: Season, dtype: int64
```

In [26]:
```python
# unique values of crop
df['Crop'].unique()
```

```
Out[26]: array(['Arecanut', 'Other Kharif pulses', 'Rice', 'Banana', 'Cashewnut',
        'Coconut ', 'Dry ginger', 'Sugarcane', 'Sweet potato', 'Tapioca',
        'Black pepper', 'Dry chillies', 'other oilseeds', 'Turmeric',
        'Maize', 'Moong(Green Gram)', 'Urad', 'Arhar/Tur', 'Groundnut',
        'Sunflower', 'Bajra', 'Castor seed', 'Cotton(lint)', 'Horse-gram',
        'Jowar', 'Korra', 'Ragi', 'Tobacco', 'Gram', 'Wheat', 'Masoor',
        'Sesamum', 'Linseed', 'Safflower', 'Onion', 'other misc. pulses',
        'Samai', 'Small millets', 'Coriander', 'Potato',
        'Other  Rabi pulses', 'Soyabean', 'Beans & Mutter(Vegetable)',
        'Bhindi', 'Brinjal', 'Citrus Fruit', 'Cucumber', 'Grapes', 'Mango',
        'Orange', 'other fibres', 'Other Fresh Fruits', 'Other Vegetables',
        'Papaya', 'Pome Fruit', 'Tomato', 'Mesta', 'Cowpea(Lobia)',
        'Lemon', 'Pome Granet', 'Sapota', 'Cabbage', 'Rapeseed &Mustard',
        'Peas  (vegetable)', 'Niger seed', 'Bottle Gourd', 'Varagu',
        'Garlic', 'Ginger', 'Oilseeds total', 'Pulses total', 'Jute',
        'Peas & beans (Pulses)', 'Blackgram', 'Paddy', 'Pineapple',
        'Barley', 'Sannhamp', 'Khesari', 'Guar seed', 'Moth',
        'Other Cereals & Millets', 'Cond-spcs other', 'Turnip', 'Carrot',
        'Redish', 'Arcanut (Processed)', 'Atcanut (Raw)',
        'Cashewnut Processed', 'Cashewnut Raw', 'Cardamom', 'Rubber',
        'Bitter Gourd', 'Drum Stick', 'Jack Fruit', 'Snak Guard', 'Tea',
        'Coffee', 'Cauliflower', 'Other Citrus Fruit', 'Water Melon',
        'Total foodgrain', 'Kapas', 'Colocosia', 'Lentil', 'Bean',
        'Jobster', 'Perilla', 'Rajmash Kholar', 'Ricebean (nagadal)',
        'Ash Gourd', 'Beet Root', 'Lab-Lab', 'Ribed Guard', 'Yam',
        'Pump Kin', 'Apple', 'Peach', 'Pear', 'Plums', 'Litchi', 'Ber',
        'Other Dry Fruit', 'Jute & mesta'], dtype=object)
```

In [27]:
```python
# after cleaning the dataset converted into new dataframe
df1=df.copy()
```

In [28]:
```python
# view the new dataframe
df1
```

Out[28]:

| | State_Name | District_Name | Crop_Year | Season | Crop | Area | Production |
|---|---|---|---|---|---|---|---|
| 0 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Arecanut | 1254.0 | 2000.0 |
| 1 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Other Kharif pulses | 2.0 | 1.0 |
| 2 | Andaman and Nicobar Islands | NICOBARS | 2000 | Kharif | Rice | 102.0 | 321.0 |
| 3 | Andaman and Nicobar Islands | NICOBARS | 2000 | Whole Year | Banana | 176.0 | 641.0 |
| 4 | Andaman and Nicobar Islands | NICOBARS | 2000 | Whole Year | Cashewnut | 720.0 | 165.0 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 246086 | West Bengal | PURULIA | 2014 | Summer | Rice | 306.0 | 801.0 |
| 246087 | West Bengal | PURULIA | 2014 | Summer | Sesamum | 627.0 | 463.0 |
| 246088 | West Bengal | PURULIA | 2014 | Whole Year | Sugarcane | 324.0 | 16250.0 |
| 246089 | West Bengal | PURULIA | 2014 | Winter | Rice | 279151.0 | 597899.0 |
| 246090 | West Bengal | PURULIA | 2014 | Winter | Sesamum | 175.0 | 88.0 |

242361 rows × 7 columns

In [29]:
```python
# check the null values of the new dataset
df1.isnull().sum()
```

Out[29]:
```
State_Name       0
District_Name    0
Crop_Year        0
Season           0
Crop             0
Area             0
Production       0
dtype: int64
```

In [30]:
```python
# check the duplicated values of the new dataframe
df1.duplicated().sum()
```

Out[30]: 0

In [31]:
```python
# converted into new file
df1.to_csv('Crop_producion data.csv',index=False)
```

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js