

Load and read the data

```
1 # Import required libraries
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import shapefile
```

```
1 # Read the shapefiles
2 def read_shp(file_path):
3     sf = shapefile.Reader(file_path)
4     fields = [f[0] for f in sf.fields[1:]]
5     records = sf.records()
6     df = pd.DataFrame(records, columns=fields)
7     return df
8
9 district_shp = read_shp("SHAPEFILES/Uganda_Districts.shp")
10 subcounty_shp = read_shp("SHAPEFILES/Uganda_Subcounties.shp")
11 maize_shp = read_shp("SHAPEFILES/Crop_Type_Map_Maize.shp")
12 sorghum_shp = read_shp("SHAPEFILES/Crop_Type_Map_Sorghum.shp")
```

```
1 # Read the two tables
2 district_df = pd.read_csv("TABLES/Uganda_Karamoja_District_Crop_Yield_Population.csv")
3 subcounty_df = pd.read_csv("TABLES/Uganda_Karamoja_Subcounty_Crop_Yield_Population.csv")
```

Data Understanding

```
1 #Check sample data for districts
2 district_df.head()
```

	OBJECTID	NAME	POP	Area	S_Yield_Ha	M_Yield_Ha	Crop_Area_Ha	S_Area_Ha	M_Area_Ha	S_Prod_Tot	M_Prod_Tot
0	92	ABIM	90385	2771977106	449	1040	5470.068394	3277.295971	1848.621855	1471506	192256
1	96	AMUDAT	101790	1643582836	205	1297	5765.443719	2973.423860	2733.661014	609552	354555
2	20	KAABONG	627057	7373606003	279	945	28121.672530	20544.194960	7394.416334	5731830	698772
3	85	KOTIDO	243157	3641539808	331	1148	53032.649450	50247.443900	1751.372284	16631904	201057
4	5	MOROTO	127811	3570160948	128	355	5954.814048	4741.748776	1190.050606	606944	42246

```
1 #Check sample data for subcounty
2 subcounty_df.head()
```

	OBJECTID	SUBCOUNTY_NAME	DISTRICT_NAME	POP	Area	Karamoja	S_Yield_Ha	M_Yield_Ha	Crop_Area_Ha	S_Area_Ha	M_Area_Ha
0	263	KACHERI	KOTIDO	17244	1067176155	Y	354.207411	1137.467019	7023.533691	6434.342449	528
1	264	KOTIDO	KOTIDO	52771	597575188	Y	367.890523	1162.996687	13587.990760	12455.592640	824
2	265	KOTIDO TOWN COUNCIL	KOTIDO	27389	23972401	Y	369.314177	1167.005832	1656.531855	1520.322052	8
3	266	NAKAPERIMORU	KOTIDO	38775	419111591	Y	283.324569	852.366578	7087.823334	6761.488901	45
4	267	PANYANGARA	KOTIDO	65704	880955930	Y	373.836926	1283.859882	10398.249390	10111.198130	172

Both columns have similar fields and most of them being numbers.

```
1 # Check info about the dfs
2 district_df.info()
3 # Dataset has 7 districts with 11 columns of crop and population and metadata.No null values
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7 entries, 0 to 6
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   OBJECTID        7 non-null     int64
```

```

1  NAME          7 non-null    object
2  POP           7 non-null    int64
3  Area          7 non-null    int64
4  S_Yield_Ha    7 non-null    int64
5  M_Yield_Ha    7 non-null    int64
6  Crop_Area_Ha  7 non-null    float64
7  S_Area_Ha     7 non-null    float64
8  M_Area_Ha     7 non-null    float64
9  S_Prod_Tot    7 non-null    int64
10 M_Prod_Tot    7 non-null    int64
dtypes: float64(3), int64(7), object(1)
memory usage: 744.0+ bytes

```

Data entails 7 districts with 11 columns. No null values.

```

1 # Check info about the dfs
2 subcounty_df.info()
3 # We have 52 rows of sub-counties with 13 columns of crop, population and meta data. No null values.

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 52 entries, 0 to 51
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   OBJECTID              52 non-null    int64
1   SUBCOUNTY_NAME      52 non-null    object
2   DISTRICT_NAME        52 non-null    object
3   POP                  52 non-null    int64
4   Area                 52 non-null    int64
5   Karamoja             52 non-null    object
6   S_Yield_Ha           52 non-null    float64
7   M_Yield_Ha           52 non-null    float64
8   Crop_Area_Ha         52 non-null    float64
9   S_Area_Ha            52 non-null    float64
10  M_Area_Ha            52 non-null    float64
11  S_Prod_Tot           52 non-null    float64
12  M_Prod_Tot           52 non-null    float64
dtypes: float64(7), int64(3), object(3)
memory usage: 5.4+ KB

```

```

1 # Read maize shapefile
2 maize_shp.head()

```

```

      DN
0      6
1      6
2      6
3      6
4      6

```

```

1 # Read sorghum shapefile
2 sorghum_shp.head()

```

```

      DN
0    212
1    212
2    212
3    212
4    212

```

```

1 sorghum_shp["DN"].nunique()
2

```

```
1
```

.shp files have 2 columns and the DN is the same throughout the dataframe

Feature Engineering

These new features to help assess the level of food security per subcounty. The per capita productivity of the 2 crops and the land-yield ratio

will inform on land use efficiency hence insight on how much is available to the population relatively.

```
1 # Create subcounty productivity per capita to gauge how much is available to the population on average
2 subcounty_df['S_Prod_Per_Capita'] = subcounty_df['S_Prod_Tot']/subcounty_df['POP']
3 subcounty_df['M_Prod_Per_Capita'] = subcounty_df['M_Prod_Tot']/subcounty_df['POP']
4 subcounty_df.head()
```

	OBJECTID	SUBCOUNTY_NAME	DISTRICT_NAME	POP	Area	Karamoja	S_Yield_Ha	M_Yield_Ha	Crop_Area_Ha	S_Area_Ha	M_Area_Ha
0	263	KACHERI	KOTIDO	17244	1067176155	Y	354.207411	1137.467019	7023.533691	6434.342449	5268.141111
1	264	KOTIDO	KOTIDO	52771	597575188	Y	367.890523	1162.996687	13587.990760	12455.592640	8240.141111
2	265	KOTIDO TOWN COUNCIL	KOTIDO	27389	23972401	Y	369.314177	1167.005832	1656.531855	1520.322052	8240.141111
3	266	NAKAPERIMORU	KOTIDO	38775	419111591	Y	283.324569	852.366578	7087.823334	6761.488901	4510.141111
4	267	PANYANGARA	KOTIDO	65704	880955930	Y	373.836926	1283.859882	10398.249390	10111.198130	17200.141111

```
1 #Create subcounty land-yield ratio to investigate land use efficiency
2 # Land use efficiency = (Crop Area / Total Area) * 100
3 subcounty_df['Land_Use_Efficiency'] = (subcounty_df['Crop_Area_Ha'] / subcounty_df['Area']) * 100
4 subcounty_df.head()
```

	OBJECTID	SUBCOUNTY_NAME	DISTRICT_NAME	POP	Area	Karamoja	S_Yield_Ha	M_Yield_Ha	Crop_Area_Ha	S_Area_Ha	M_Area_Ha
0	263	KACHERI	KOTIDO	17244	1067176155	Y	354.207411	1137.467019	7023.533691	6434.342449	5268.141111
1	264	KOTIDO	KOTIDO	52771	597575188	Y	367.890523	1162.996687	13587.990760	12455.592640	8240.141111
2	265	KOTIDO TOWN COUNCIL	KOTIDO	27389	23972401	Y	369.314177	1167.005832	1656.531855	1520.322052	8240.141111
3	266	NAKAPERIMORU	KOTIDO	38775	419111591	Y	283.324569	852.366578	7087.823334	6761.488901	4510.141111
4	267	PANYANGARA	KOTIDO	65704	880955930	Y	373.836926	1283.859882	10398.249390	10111.198130	17200.141111

Data Manipulation

We want to retain only what is necessary from our data set for visualization purposes.

```
1 # Add a prefix D to district_df columns to differentiate them from subcounty_df columns when doing the merge.
2 district_df.columns = ['D' + col for col in district_df.columns]
3 district_df.info()
4
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7 entries, 0 to 6
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   DOBJECTID        7 non-null      int64
1   DNAME            7 non-null      object
2   DPOP             7 non-null      int64
3   DArea            7 non-null      int64
4   DS_Yield_Ha      7 non-null      int64
5   DM_Yield_Ha      7 non-null      int64
6   DCrop_Area_Ha    7 non-null      float64
7   DS_Area_Ha       7 non-null      float64
8   DM_Area_Ha       7 non-null      float64
9   DS_Prod_Tot      7 non-null      int64
10  DM_Prod_Tot      7 non-null      int64
dtypes: float64(3), int64(7), object(1)
memory usage: 744.0+ bytes
```

```
1 # Create District Yield, District Crop Area, District Prod by summing the respective crop columns; to see the crop output at
2 district_df['District Yield'] = district_df['DS_Yield_Ha'] + district_df['DM_Yield_Ha']
3 district_df['District Crop Area'] = district_df['DS_Area_Ha'] + district_df['DM_Area_Ha']
4 district_df['District Prod'] = district_df['DS_Prod_Tot'] + district_df['DM_Prod_Tot']
5
6 district_df.head()
```

	DOBJECTID	DNAME	DPOP	DArea	DS_Yield_Ha	DM_Yield_Ha	DCrop_Area_Ha	DS_Area_Ha	DM_Area_Ha	DS_Prod_Tot	DM_Prod_Tot
0	92	ABIM	90385	2771977106	449	1040	5470.068394	3277.295971	1848.621855	1471506	1471506
1	96	AMUDAT	101790	1643582836	205	1297	5765.443719	2973.423860	2733.661014	609552	609552
2	20	KAABONG	627057	7373606003	279	945	28121.672530	20544.194960	7394.416334	5731830	5731830
3	85	KOTIDO	243157	3641539808	331	1148	53032.649450	50247.443900	1751.372284	16631904	16631904
4	5	MOROTO	127811	3570160948	128	355	5954.814048	4741.748776	1190.050606	606944	606944

```
1 # Drop the unnecessary columns e.g individual crop columns from which we just summed above as they will be redundant once we
2 columns_to_drop = ['DS_Yield_Ha', 'DM_Yield_Ha', 'DS_Area_Ha', 'DM_Area_Ha', 'DS_Prod_Tot', 'DM_Prod_Tot', 'DOBJECTID']
3
4 district_df = district_df.drop(columns=columns_to_drop)
5 district_df.head()
6
```

	DNAME	DPOP	DArea	DCrop_Area_Ha	District Yield	District Crop Area	District Prod
0	ABIM	90385	2771977106	5470.068394	1489	5125.917826	3394073
1	AMUDAT	101790	1643582836	5765.443719	1502	5707.084874	4155110
2	KAABONG	627057	7373606003	28121.672530	1224	27938.611294	12719553
3	KOTIDO	243157	3641539808	53032.649450	1479	51998.816184	18642479
4	MOROTO	127811	3570160948	5954.814048	483	5931.799382	1029412

```
1 # Merge subcounty_df to district_df to have a full enriched dataset
2 master_df = pd.merge(left = subcounty_df, right = district_df, how= "left", left_on='DISTRICT_NAME', right_on='DNAME')
```

```
1 #Drop unnecessary/duplicated columns
2 master_df.drop(columns=['DNAME', 'Karamoja'], inplace=True)
```

```
1 # Check master_df
2 master_df.head()
```

	OBJECTID	SUBCOUNTY_NAME	DISTRICT_NAME	POP	Area	S_Yield_Ha	M_Yield_Ha	Crop_Area_Ha	S_Area_Ha	M_Area_Ha	Prod
0	263	KACHERI	KOTIDO	17244	1067176155	354.207411	1137.467019	7023.533691	6434.342449	528.124229	528.124229
1	264	KOTIDO	KOTIDO	52771	597575188	367.890523	1162.996687	13587.990760	12455.592640	824.767081	824.767081
2	265	KOTIDO TOWN COUNCIL	KOTIDO	27389	23972401	369.314177	1167.005832	1656.531855	1520.322052	8.561644	8.561644
3	266	NAKAPERIMORU	KOTIDO	38775	419111591	283.324569	852.366578	7087.823334	6761.488901	45.721712	45.721712
4	267	PANYANGARA	KOTIDO	65704	880955930	373.836926	1283.859882	10398.249390	10111.198130	172.611914	172.611914

5 rows × 12 columns

```
1 master_df.describe()
```

	OBJECTID	POP	Area	S_Yield_Ha	M_Yield_Ha	Crop_Area_Ha	S_Area_Ha	M_Area_Ha	S_Prod_Tot	M_Prod_Tot
count	52.000000	52.000000	5.200000e+01	52.000000	52.000000	52.000000	52.000000	52.000000	5.200000e+01	5.200000e+01
mean	787.865385	28934.692308	5.331913e+08	274.165405	940.259552	2839.646974	2253.143395	536.300569	6.557443e+05	5.57443e+05
std	280.101314	20865.122974	4.913308e+08	118.569907	321.641901	3110.505917	2954.355858	724.092288	9.915839e+05	7.915839e+05
min	263.000000	1418.000000	2.121209e+06	108.156411	0.000000	0.171390	0.130941	0.000000	1.728126e+01	0.000000
25%	597.750000	16558.500000	1.568923e+08	173.034066	743.075879	964.876031	405.394759	79.821743	1.210555e+05	6.000000e+04
50%	810.500000	23053.500000	3.848356e+08	277.255206	1016.684002	1654.265138	1231.824455	326.479336	2.543687e+05	2.843687e+05
75%	982.250000	39461.000000	7.749029e+08	368.246437	1203.548665	3267.564651	2429.985069	740.296675	6.040942e+05	8.104942e+05
max	1320.000000	100919.000000	2.069555e+09	560.313070	1396.991494	13587.990760	12964.499730	3840.698081	4.582294e+06	4.582294e+06

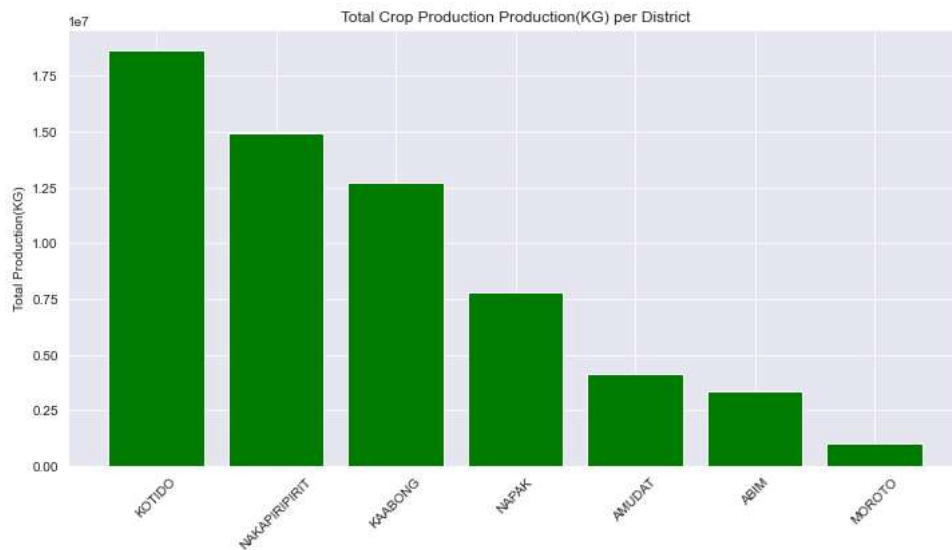
Summary statistics for 52 subcounties showing wide variation in population (1,418 to 100,919) and agricultural production. Mean sorghum yield is 274 kg/ha versus 940 kg/ha for maize, with average total crop area of 2,840 hectares per sub-county. District-level aggregates indicate mean production of 9.9 million tons across an average crop area of 22,675 hectares.

```
1 master_df.info()
2 # Since our left table was subcounties, we have 52 rows of data with the enjoined district data.

<class 'pandas.core.frame.DataFrame'>
Int64Index: 52 entries, 0 to 51
Data columns (total 21 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   OBJECTID                             52 non-null     int64
1   SUBCOUNTY_NAME                      52 non-null     object
2   DISTRICT_NAME                        52 non-null     object
3   POP                                  52 non-null     int64
4   Area                                 52 non-null     int64
5   S_Yield_Ha                           52 non-null     float64
6   M_Yield_Ha                           52 non-null     float64
7   Crop_Area_Ha                         52 non-null     float64
8   S_Area_Ha                            52 non-null     float64
9   M_Area_Ha                            52 non-null     float64
10  S_Prod_Tot                           52 non-null     float64
11  M_Prod_Tot                           52 non-null     float64
12  S_Prod_Per_Capita                    52 non-null     float64
13  M_Prod_Per_Capita                    52 non-null     float64
14  Land_Use_Efficiency                  52 non-null     float64
15  DPOP                                 52 non-null     int64
16  DArea                                52 non-null     int64
17  DCrop_Area_Ha                        52 non-null     float64
18  District Yield                       52 non-null     int64
19  District Crop Area                    52 non-null     float64
20  District Prod                         52 non-null     int64
dtypes: float64(12), int64(7), object(2)
memory usage: 8.9+ KB
```

✖ Exploratory Data Analysis

```
1 #Sort the district productivity from highest to lowest to see the output across.
2 district_df_sorted = master_df.sort_values(by="District Prod", ascending=False)
3
4 plt.figure(figsize=(12,6))
5
6
7 plt.bar(district_df_sorted["DISTRICT_NAME"], district_df_sorted["District Prod"], label="Maize Productivity", color="green")
8
9 #Details
10 plt.xticks(rotation=45)
11 plt.ylabel("Total Production(KG)")
12 plt.title("Total Crop Production Production(KG) per District")
13 # plt.legend()
14 plt.show()
```



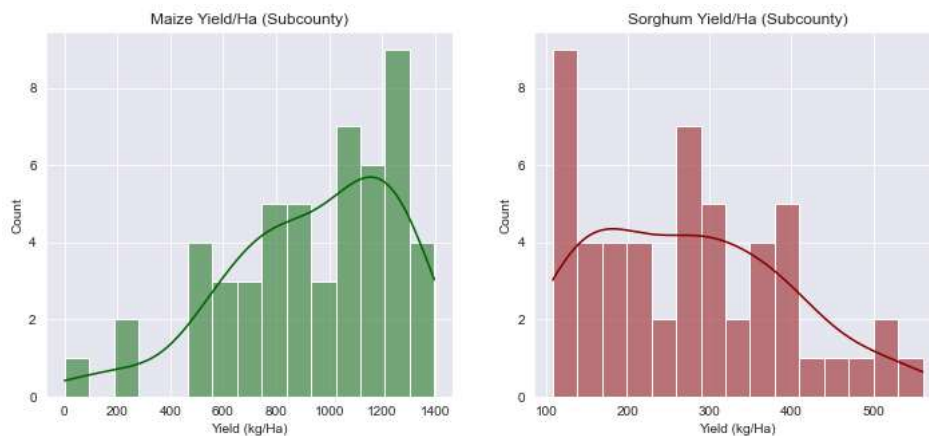
Observation:

Kotido has the highest crop production at 18642479 Kgs while Moroto is performing dismally.

```

1 #Plot a histogram for each crop to investigate the distribution
2 fig, axes = plt.subplots(1,2, figsize=(12,5))
3
4 sns.set_style("darkgrid")
5 sns.histplot(master_df["M_Yield_Ha"], bins=15, kde=True, ax=axes[0], color="darkgreen")
6 axes[0].set_title("Maize Yield/Ha (Subcounty)")
7 axes[0].set_xlabel("Yield (kg/Ha)")
8
9 sns.histplot(master_df["S_Yield_Ha"], bins=15, kde=True, ax=axes[1], color="darkred")
10 axes[1].set_title("Sorghum Yield/Ha (Subcounty)")
11 axes[1].set_xlabel("Yield (kg/Ha)")
12
13 plt.show()

```



Observation:

Maize yield has a positive distribution with more subcounties producing between 1000-1200Kgs while sorghum has a negative distribution with more subcounties producing between 100-300Kgs.

This implies maize production is well established in most areas.

```

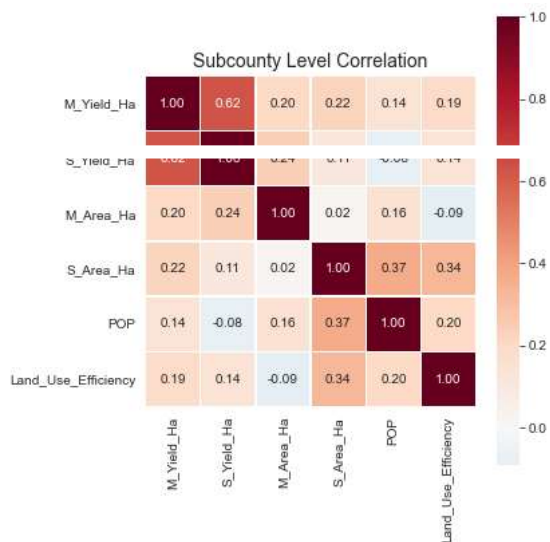
1 # Perform a correlation analysis to analyze the variables' relationship
2
3 plt.figure(figsize=(6,6))
4
5 # Include efficiency and per capita metrics in correlation
6 subcounty_corr = master_df[["M_Yield_Ha", "S_Yield_Ha", "M_Area_Ha", "S_Area_Ha", "POP",

```

```

7         "Land_Use_Efficiency"]].corr()
8
9     sns.heatmap(subcounty_corr,
10                 annot=True,
11                 cmap="RdBu_r",
12                 center=0,
13                 fmt='.2f',
14                 square=True,
15                 linewidths=0.5)
16
17 #Details
18 plt.title("Subcounty Level Correlation" , fontsize=14)
19 plt.tight_layout()
20 plt.show()

```



Observations

- Maize and sorghum yields have a relatively strong correlation, implying both cereals tend to perform well in the given conditions.
- Population has a weak positive correlation with land use. Hence, a higher population does not imply better land use.
- Maize and sorghum areas have a negligible relationship; there is no relation on the area on which they are cultivated.

Export dataframes as csv

```

1 # Download the district data
2 district_df.to_csv('TABLES/karamoja_district.csv', index=False)
3
4 # Download the subcounty data
5 subcounty_df.to_csv('TABLES/karamoja_subcounty.csv', index=False)
6
7 # Download master_df
8 master_df.to_csv('TABLES/karamoja_master.csv', index=False)
9
10

```