

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 file_path='CommMktArrivals2012.csv'
6 # --- 1. Load and Clean Data ---
7
8 # Specify 'latin1' encoding to resolve UnicodeDecodeError
9 df = pd.read_csv(file_path, encoding='latin1')
10 # Strip whitespace from string columns to prevent matching errors
11 df = df.apply(lambda x: x.str.strip() if x.dtype == "object" else x)
12
13 # Standardize Month column for correct chronological sorting
14 month_order = ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sep', 'Oct',
15 df['Month'] = pd.Categorical(df['Month'], categories=month_order, ordered=True)
16
17 # Clean special characters from 'District Name' that cause font warnings
18 df['District Name'] = df['District Name'].str.replace('█', '', regex=False)
19
20 # Subset for volume-based analysis (Quintals)
21 df_q = df[df['Unit'] == 'Quintal'].copy()



```

```

1 # --- 2. Descriptive Statistics ---
2 print("--- Dataset Info ---")
3 info_data = []
4 for col in df.columns:
5     info_data.append({
6         'Column Name': col,
7         'Data Type': df[col].dtype,
8         'Non-Null Count': df[col].count()
9     })
10 df_info_table = pd.DataFrame(info_data)
11 display(df_info_table)

```

--- Dataset Info ---

	Column Name	Data Type	Non-Null Count	
0	District Name	object	21422	
1	Taluk Name	object	21422	
2	Market Name	object	21422	
3	Address	object	21422	
4	Telephone	object	21422	
5	Commodity	object	21422	
6	Year	int64	21422	
7	Month	category	21422	
8	Arrival	int64	21422	
9	Unit	object	21422	

Next steps:

[Generate code with df\\_info\\_table](#)

[New interactive sheet](#)

```

1 print("\n--- Summary Statistics ---")
2 desc_stats = df.describe(include='all').T
3 rows_to_drop = ['mean', 'min', '25%', '50%', '75%', 'max', 'std']
4 desc_stats = desc_stats.drop(rows_to_drop, errors='ignore') # Dropping numerical statistics
5 display(desc_stats)
6
7 # Demonstrate numpy usage: calculate mean arrival
8 mean_arrival = np.mean(df_q['Arrival'])
9 print(f"\n--- Mean Arrival Volume (Quintals) using NumPy: {mean_arrival:.2f} ---")

```

--- Summary Statistics ---

	count	unique	top	freq	mean	std	min	25%
District Name	21422	29	Hassan	1504	NaN	NaN	NaN	NaN
Taluk Name	21422	147	Mysore	669	NaN	NaN	NaN	NaN
Market Name	21422	148	MYSORE	669	NaN	NaN	NaN	NaN
Address	21422	148	NANJANGUD ROAD , BANDIPALYA MYSORE	669	NaN	NaN	NaN	NaN
Telephone	21422	148	08212482725	669	NaN	NaN	NaN	NaN
Commodity	21422	156	Maize	1021	NaN	NaN	NaN	NaN
Year	21422.0	NaN	NaN	NaN	2012.0	0.0	2012.0	2012.0
Month	21422	12	Dec	1863	NaN	NaN	NaN	NaN



Next steps:

[Generate code with desc\\_stats](#)

[New interactive sheet](#)

```
1 # --- 3. Print Data Summaries ---
2 print("\n--- Top 10 Commodities List ---")
3 display(top_10)
```

--- Top 10 Commodities List ---

	Commodity	Arrival	
0	Paddy	18554513	
1	Maize	15460281	
2	Onion	9826018	
3	Rice	9159723	
4	Potato	4467185	
5	Green Ginger	2952505	
6	Cotton	2861896	
7	Tomato	2858265	
8	Arecanut	2696704	
9	Tur	2562682	

Next steps:

[Generate code with top\\_10](#)


[New interactive sheet](#)

```

1 # 4. SUPPLY CHAIN RISK ANALYSIS (Volatility)
2 # -----
3 # We calculate the Coefficient of Variation (CV).
4 # Higher CV = More unpredictable supply = High risk for the market.
5 top_10_names = df_q.groupby('Commodity')['Arrival'].sum().nlargest(10).index
6 volatility = df_q[df_q['Commodity'].isin(top_10_names)].groupby('Commodity')['Arrival'].std().div(df_q[df_q['Commodity'].isin(top_10_names)].groupby('Commodity')['Arrival'].mean()) * 100
7 volatility['Risk_Score(CV)'] = (volatility['std'] / volatility['mean']) * 100
8
9 print("--- Supply Chain Risk Assessment ---")
10 display(volatility[['Risk_Score(CV)']].sort_values(by='Risk_Score(CV)', ascending=False))
11

```

--- Supply Chain Risk Assessment ---

	Risk_Score(CV) 
Commodity	
<b>Arecanut</b>	693.487346
<b>Rice</b>	535.230129
<b>Onion</b>	430.868343
<b>Maize</b>	415.997238
<b>Potato</b>	343.194815
<b>Tur</b>	330.170121
<b>Paddy</b>	280.496492
<b>Green Ginger</b>	268.315330
<b>Tomato</b>	255.913079
<b>Cotton</b>	206.524517

```

1 print("\n--- Raw Data Head (Top 10 Commodities) ---")
2 display(df_q[df_q['Commodity'].isin(top_10['Commodity'])].head(10))

```

--- Raw Data Head (Top 10 Commodities) ---

	District Name	Taluk Name	Market Name	Address	Telephone	Commodity	Year	Month	Arrival
7	Bagalakot	Badami	BADAMI	SECRATRY A.P.M.C.BADAMI BADAMI	220042	Maize	2012	Jan	21669
24	Bagalakot	Badami	BADAMI	SECRATRY A.P.M.C.BADAMI BADAMI	220042	Maize	2012	Feb	5184
40	Bagalakot	Badami	BADAMI	SECRATRY A.P.M.C.BADAMI BADAMI	220042	Maize	2012	Mar	10346
50	Bagalakot	Badami	BADAMI	SECRATRY A.P.M.C.BADAMI BADAMI	220042	Cotton	2012	Apr	106
57	Bagalakot	Badami	BADAMI	SECRATRY A.P.M.C.BADAMI BADAMI	220042	Maize	2012	Apr	107

```
1 from sklearn.cluster import KMeans
2 from sklearn.preprocessing import StandardScaler
3
4 # --- ADVANCED STEP: K-MEANS CLUSTERING ---
5 # Let's cluster districts by their top 3 crops (Paddy, Maize, Onion)
6 top_3_crops = ['Paddy', 'Maize', 'Onion']
7 cluster_prep = df_q[df_q['Commodity'].isin(top_3_crops)].pivot_table(
8     index='District Name', columns='Commodity', values='Arrival', aggfunc='sum', fill
9 )
10
11 # Standardize for ML
12 scaler = StandardScaler()
13 scaled_data = scaler.fit_transform(cluster_prep)
14
15 # Identify 3 clusters of Districts
16 kmeans = KMeans(n_clusters=3, random_state=42, n_init=10)
17 cluster_prep['Zone_Cluster'] = kmeans.fit_predict(scaled_data)
18
19 print("District Zones Identified:")
20 display(cluster_prep.groupby('Zone_Cluster').mean())
```

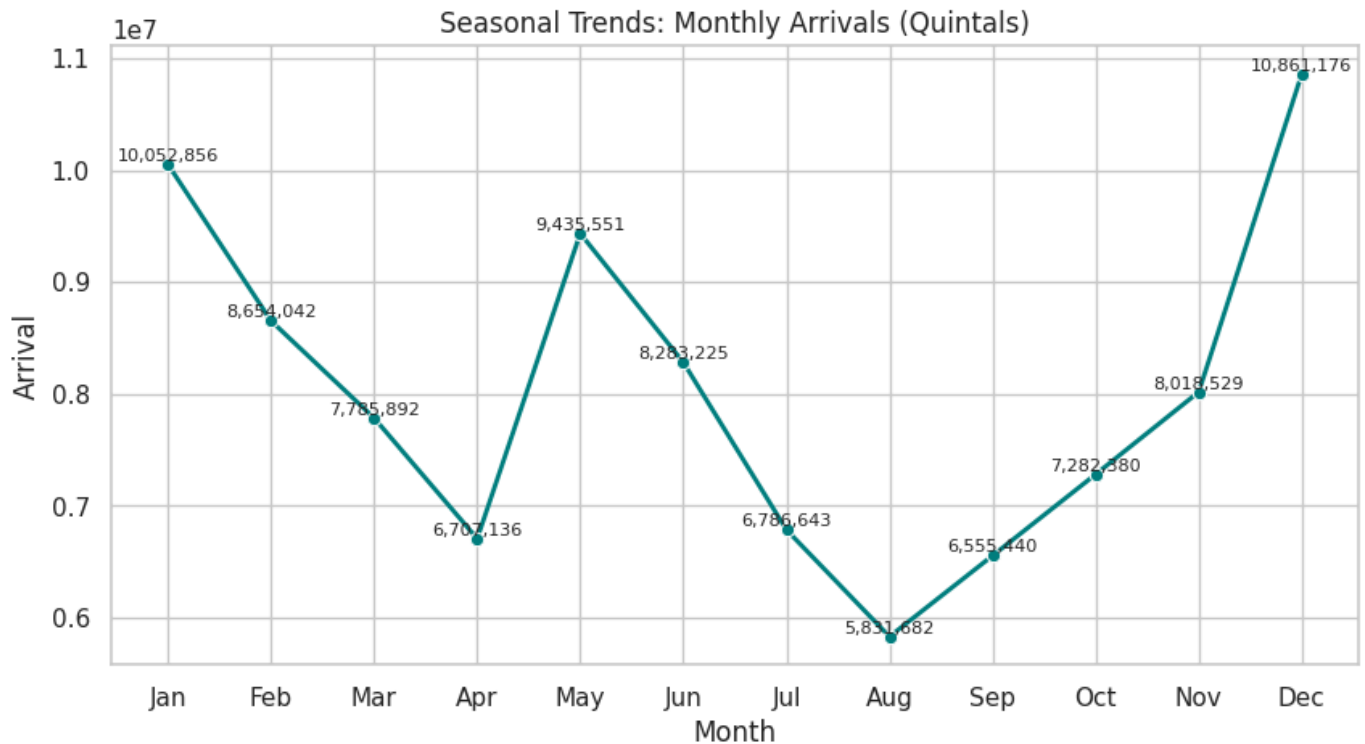
District Zones Identified:

Commodity	Maize	Onion	Paddy
Zone_Cluster			
0	7.525500e+04	6623416.0	1.300000e+02
1	2.114878e+05	128589.0	2.058246e+05
2	1.788716e+06	62274.0	2.337707e+06

```

1 # --- 5. Visualization Suite ---
2 sns.set(style="whitegrid")
3
4 # A. Seasonal Trends (Total Volume by Month)
5 monthly = df_q.groupby('Month', observed=False)['Arrival'].sum().reset_index()
6 plt.figure(figsize=(10, 5))
7 ax = sns.lineplot(data=monthly, x='Month', y='Arrival', marker='o', color='teal', li
8 plt.title('Seasonal Trends: Monthly Arrivals (Quintals)')
9 for i, point in monthly.iterrows():
10     ax.text(point['Month'], point['Arrival'], f"{point['Arrival']:,}", ha='center', v

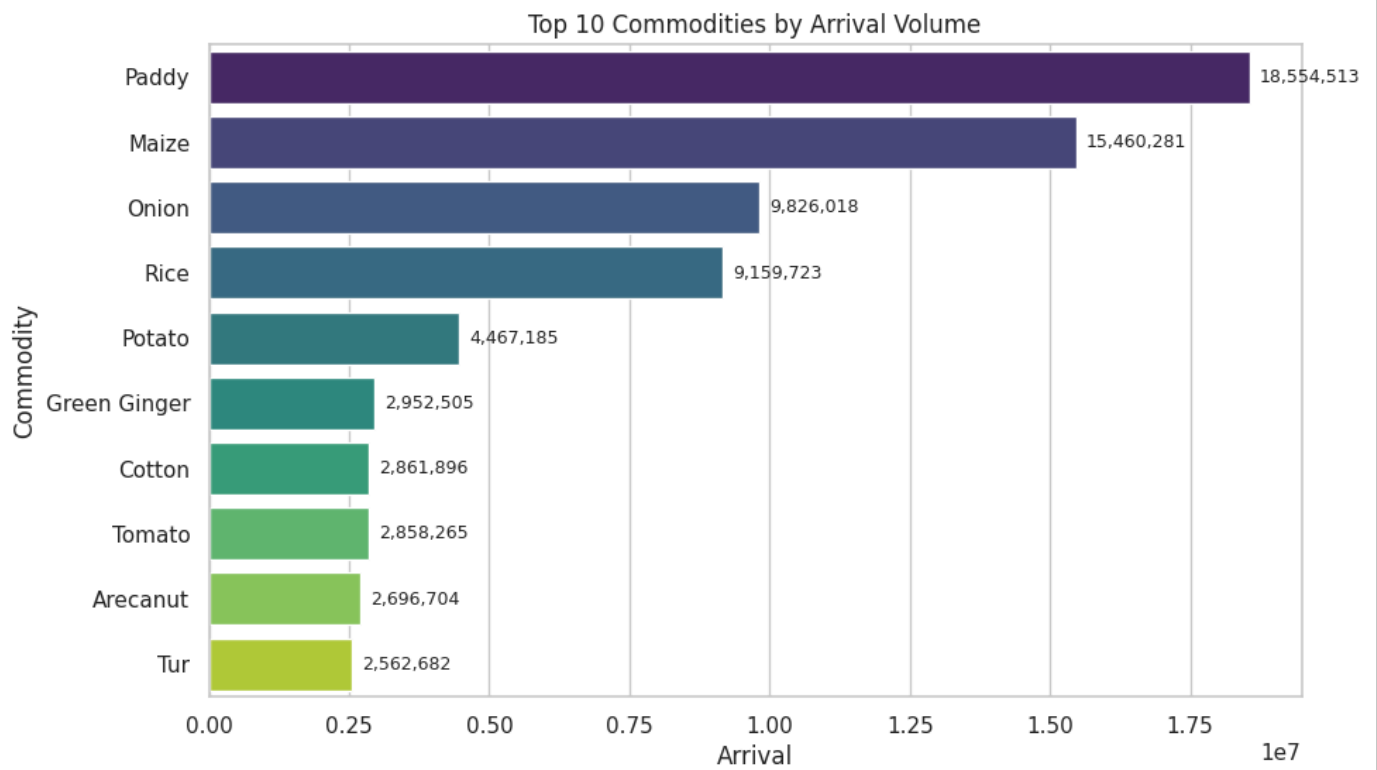
```



```

1 # B. Top 10 Commodities
2 top_10 = df_q.groupby('Commodity')['Arrival'].sum().sort_values(ascending=False).head(10)
3 plt.figure(figsize=(10, 6))
4 ax = sns.barplot(data=top_10, x='Arrival', y='Commodity', palette='viridis', hue='Commodity')
5 plt.title('Top 10 Commodities by Arrival Volume')
6 for p in ax.patches:
7     ax.annotate(f'{p.get_width():.0f}', (p.get_width(), p.get_y() + p.get_height() / 2),
8               ha='left', va='center', xytext=(5, 0), textcoords='offset points', fontweight='bold')

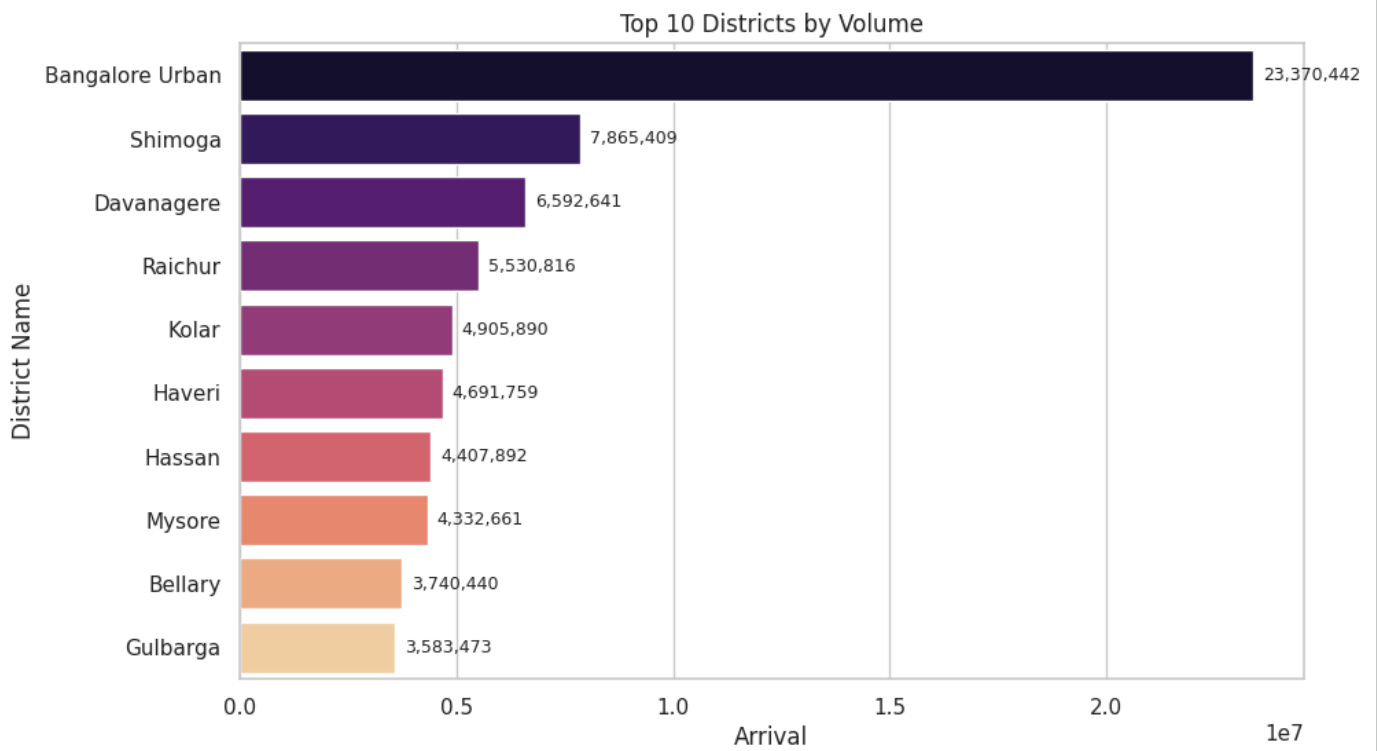
```



```

1 # C. Top 10 Districts
2 districts = df_q.groupby('District Name')['Arrival'].sum().sort_values(ascending=False)
3 plt.figure(figsize=(10, 6))
4 ax = sns.barplot(data=districts, x='Arrival', y='District Name', palette='magma', hue='District Name')
5 plt.title('Top 10 Districts by Volume')
6 for p in ax.patches:
7     ax.annotate(f'{p.get_width():.0f}', (p.get_width(), p.get_y() + p.get_height()),
8             ha='left', va='center', xytext=(5, 0), textcoords='offset points', fontweight='bold')

```



```

1 # D. Professional Seasonality Heatmap
2 plt.figure(figsize=(12, 6))
3 top_5 = df_q.groupby('Commodity')['Arrival'].sum().nlargest(5).index
4 pivot = df_q[df_q['Commodity'].isin(top_5)].pivot_table(index='Commodity', columns='I
5 sns.heatmap(pivot.div(pivot.max(axis=1), axis=0), cmap="YlGnBu", annot=True)
6 plt.title("Supply Chain Pressure: Peak Arrival Heatmap")

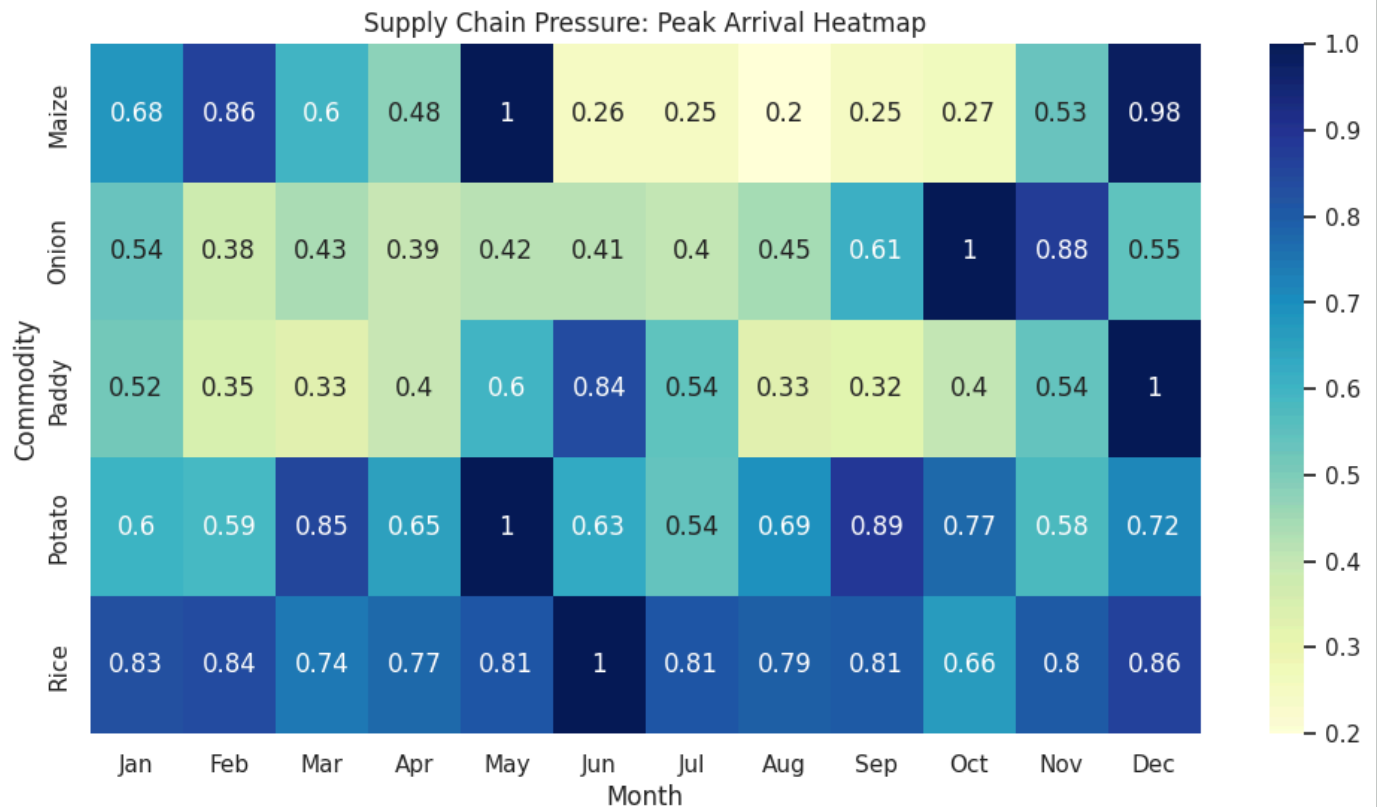
```



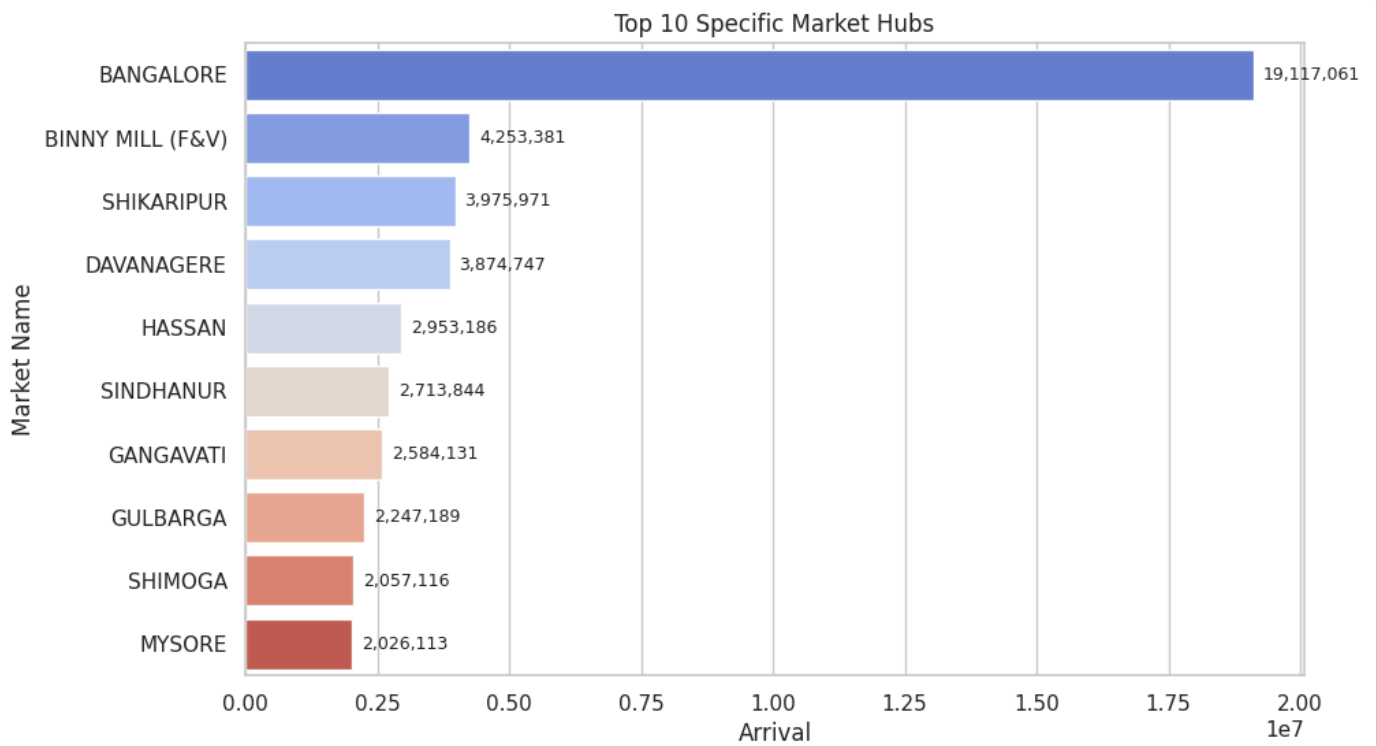
/tmp/ipython-input-3960951349.py:4: FutureWarning:

The default value of observed=False is deprecated and will change to observed=True in a fu

Text(0.5, 1.0, 'Supply Chain Pressure: Peak Arrival Heatmap')



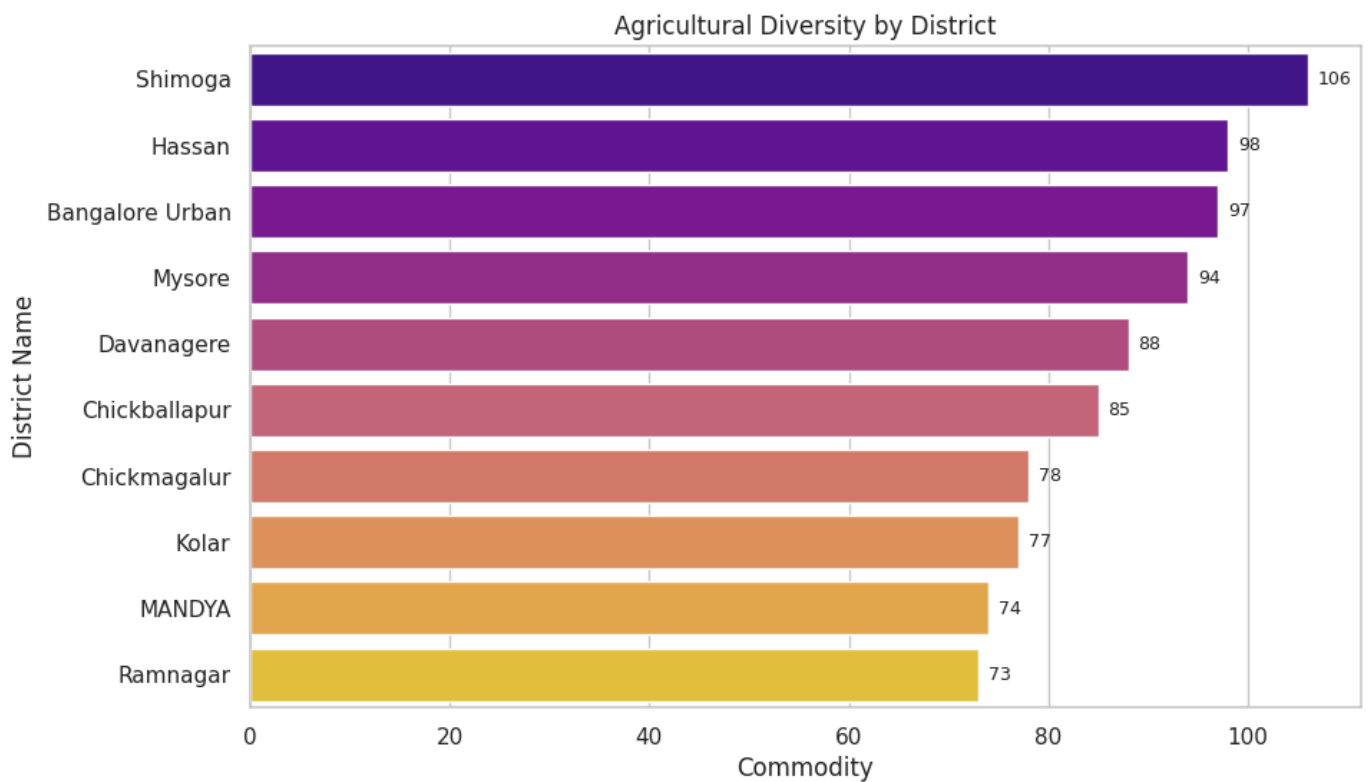
```
1 # E. Top Market Hubs
2 markets = df_q.groupby('Market Name')['Arrival'].sum().sort_values(ascending=False).l
3 plt.figure(figsize=(10, 6))
4 ax = sns.barplot(data=markets, x='Arrival', y='Market Name', palette='coolwarm', hue:
5 plt.title('Top 10 Specific Market Hubs')
6 for p in ax.patches:
7     ax.annotate(f'{p.get_width():.0f}', (p.get_width(), p.get_y() + p.get_height() ,
8         ha='left', va='center', xytext=(5, 0), textcoords='offset points', fo
9
```



```

1 # F. Regional Diversity (Unique Commodities per District)
2 diversity = df.groupby('District Name')['Commodity'].nunique().sort_values(ascending=
3 plt.figure(figsize=(10, 6))
4 ax = sns.barplot(data=diversity, x='Commodity', y='District Name', palette='plasma',
5 plt.title('Agricultural Diversity by District')
6 for p in ax.patches:
7     ax.annotate(f'{p.get_width():.0f}', (p.get_width(), p.get_y() + p.get_height() /
8         ha='left', va='center', xytext=(5, 0), textcoords='offset points', fr
9

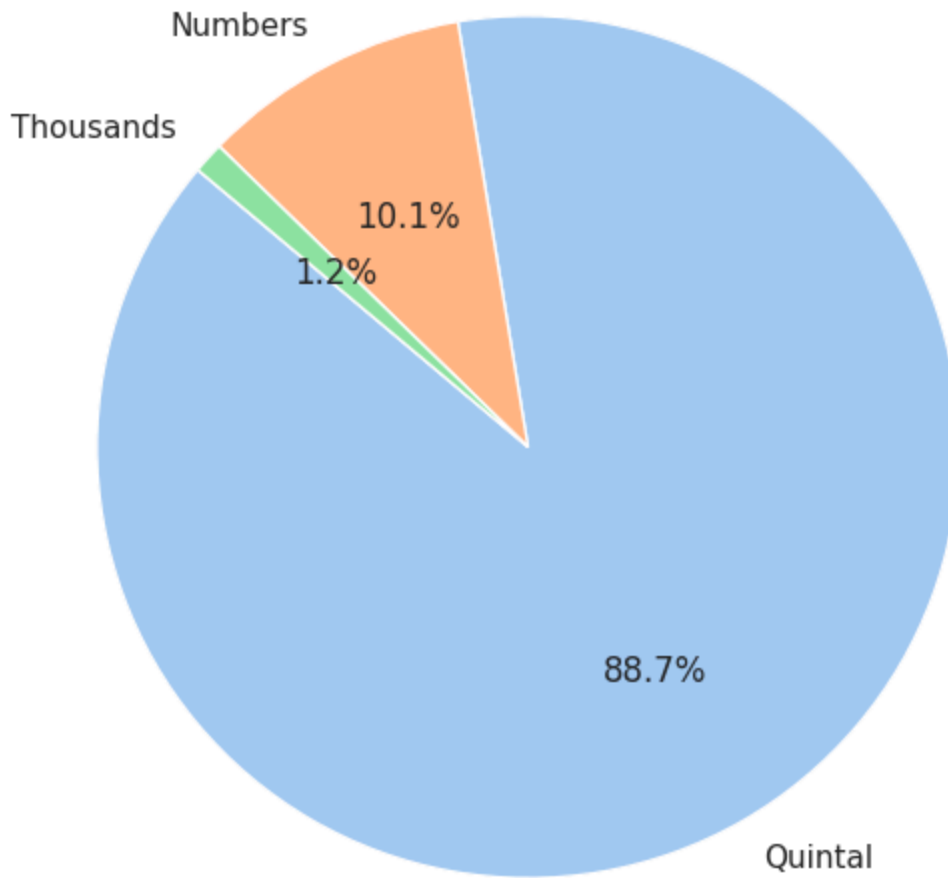
```



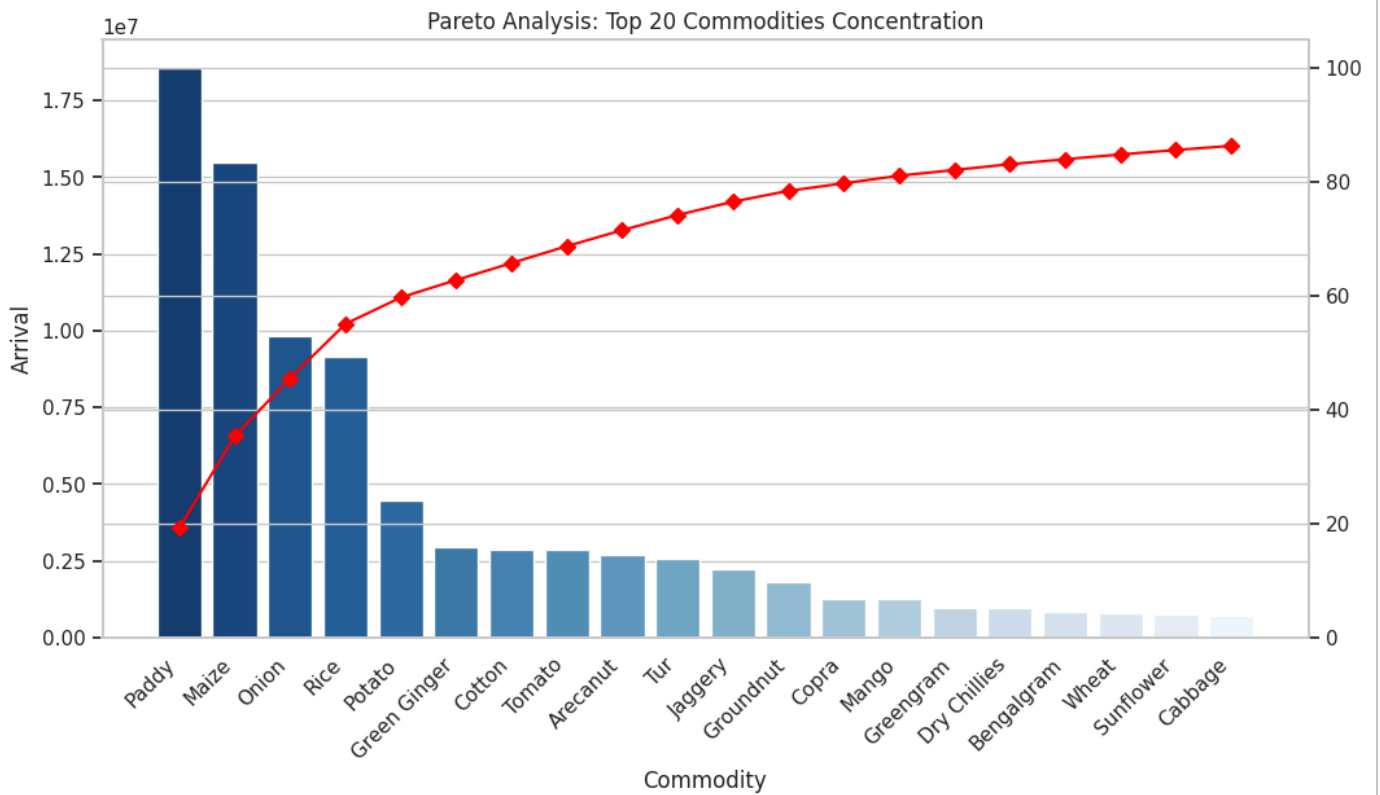
```
1 # G. Unit Composition
2 plt.figure(figsize=(7, 7))
3 df['Unit'].value_counts().plot.pie(autopct='%1.1f%%', startangle=140, colors=sns.col
4 plt.title('Data Distribution by Unit')
5 plt.ylabel('')
6
7
8
```

Text(0, 0.5, '')

Data Distribution by Unit



```
1 # H. Pareto Analysis (Volume Concentration)
2 comm_totals = df_q.groupby('Commodity')['Arrival'].sum().sort_values(ascending=False)
3 comm_totals['Cumulative %'] = 100 * comm_totals['Arrival'].cumsum() / comm_totals['Arrival'].sum()
4 fig, ax1 = plt.subplots(figsize=(12, 6))
5 sns.barplot(data=comm_totals.head(20), x='Commodity', y='Arrival', ax=ax1, palette='magma')
6 plt.xticks(rotation=45, ha='right')
7 ax2 = ax1.twinx()
8 ax2.plot(comm_totals.head(20)['Commodity'], comm_totals.head(20)['Cumulative %'], color='red')
9 ax2.set_ylim(0, 105)
10 plt.title('Pareto Analysis: Top 20 Commodities Concentration')
11 plt.savefig('8_pareto_analysis.png')
12
```

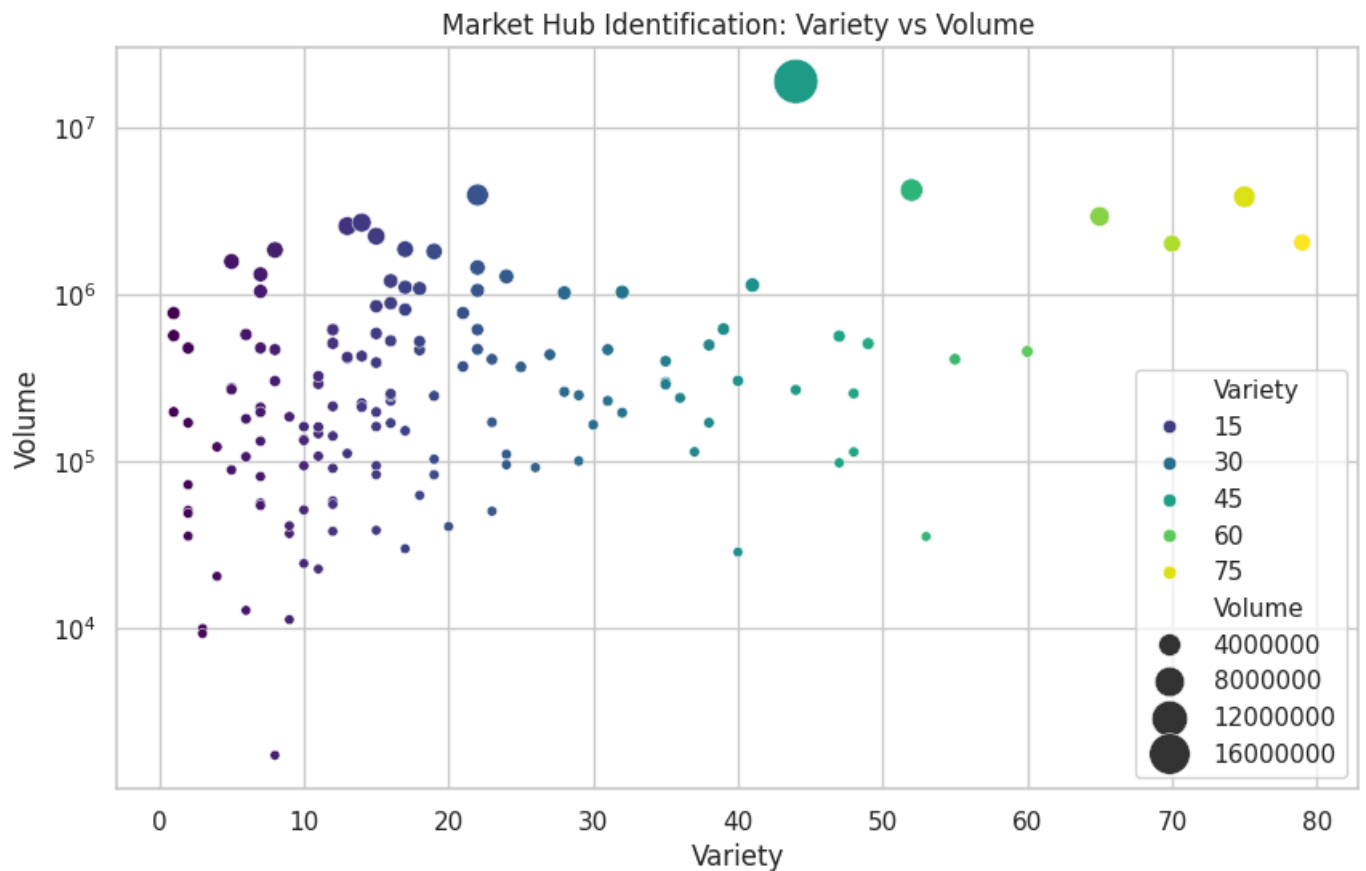


```

1 # I. Market Diversity Plot
2 market_stats = df_q.groupby(['District Name', 'Market Name']).agg(Variety=('Commodity
3 plt.figure(figsize=(10, 6))
4 sns.scatterplot(data=market_stats, x='Variety', y='Volume', hue='Variety', palette='
5 plt.yscale('log')
6 plt.title("Market Hub Identification: Variety vs Volume")

```

Text(0.5, 1.0, 'Market Hub Identification: Variety vs Volume')



```
1 # J MARKET DIVERSITY INDEX ---
2 # This identifies which markets are the backbone of food security
3 market_diversity = df_q.groupby(['District Name', 'Market Name']).agg(
4     Commodity_Count=('Commodity', 'nunique'),
5     Total_Volume=('Arrival', 'sum')
6 ).reset_index()
7
8 # Categorize Markets
9 def categorize_market(row):
10     if row['Commodity_Count'] > 50: return 'Mega Hub'
11     if row['Commodity_Count'] > 20: return 'Regional Center'
12     return 'Specialized/Local'
13
14 market_diversity['Market_Type'] = market_diversity.apply(categorize_market, axis=1)
15
16 plt.figure(figsize=(12, 6))
17 sns.scatterplot(data=market_diversity, x='Commodity_Count', y='Total_Volume',
18                 hue='Market_Type', size='Total_Volume', sizes=(50, 500), alpha=0.7)
19 plt.yscale('log')
20 plt.title('Market Classification: Diversity vs. Volume', fontsize=15)
21 plt.xlabel('Number of Unique Commodities', fontsize=12)
```

```
22 plt.ylabel('Total volume (Log Scale)', fontsize=12)
23 plt.show()
```

