

Agricultural Product Distribution Forecasting: Predicting Demand to Optimize Harvest Schedules

1. Project Overview

1.1 Rationale for the Project

- Enhancing Efficiency: Maximize resource utilization with accurate demand forecasting and optimized harvest schedules.
- Reducing Food Waste: Align harvest times with demand predictions to minimize waste.
- Improved Profit Margins: Enhance profit margins through efficient distribution, reduced wastage, and meeting customer demand.
- Competitive Advantage: Stay ahead of market trends and consistently deliver fresh produce.

1.2 Aim of the Project

- Develop a robust demand forecasting model to predict product demand based on historical data, weather conditions, and market trends.
- Identify key parameters that affect production demand.

2. Business Overview/Problem

GreenSeason Farms faces several challenges:

- Inefficient Harvesting: Underutilization of resources due to reliance on traditional seasonal patterns.
- Demand Variability: Difficulty in predicting demand fluctuations, leading to overstocking or understocking.
- Shelf Life Management: Ensuring products reach customers at peak freshness.
- Transportation Optimization: Finding optimal routes and delivery schedules to minimize transportation costs.

2.1 Tech Stack

- Pandas
- NumPy

- SciPy
- Matplotlib
- Seaborn
- Statsmodels
- Prophet
- ARIMA
- Scikit-learn

2.2 Project Scope

- Data Collection & Preparation
- Exploratory Data Analysis
- Model Evaluation
- Demand Forecasting Mode

3. Data Collection and Preparation

Import Packages

In [87]:

```
1 # Importing packages
```

```
In [44]: 1 import pandas as pd
          2 import numpy as np
          3 import matplotlib.pyplot as plt
          4 import seaborn as sns
          5 from scipy import stats
          6 from statsmodels.tsa.stattools import adfuller
          7 from statsmodels.tsa.arima.model import ARIMA
          8 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
          9 from prophet import Prophet
         10 from sklearn.ensemble import RandomForestRegressor
         11 from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error, mean_squared_error
         12
         13 import warnings
         14 warnings.filterwarnings('ignore')
         15
```

Data Collection

```
In [45]: 1 df = pd.read_csv(r"C:\Users\Amdari\OneDrive\Documents\Amdari\Agricultural Product\Dataset.csv")
```

```
In [46]: 1 df.head()
```

Out[46]:

	Date	Product	Quantity_Sold	Revenue	Temperature_Celsius	Rainfall_mm	Location	Transportation_Cost	Labor_Cost (C)
0	2015-01-01	Strawberries	690.0	70.298339	9.680289	0.077279	Field C	20	13.771809
1	2015-01-01	Apples	354.0	599.863944	9.680289	0.077279	Field A	50	9.702670
2	2015-01-01	Tomatoes	275.0	199.249505	9.680289	0.077279	Field C	20	22.356498
3	2015-01-01	Apples	1079.0	1136.111770	9.680289	0.077279	Field A	50	15.691418
4	2015-01-01	Strawberries	1008.0	111.653303	9.680289	0.077279	Field B	30	17.473353

```
In [47]: 1 df.describe()
```

Out[47]:

	Quantity_Sold	Revenue	Temperature_Celsius	Rainfall_mm	Transportation_Cost	Labor_Cost	Quality_Score	Inventory_Level
count	28382.000000	28388.000000	28377.000000	28377.000000	28388.000000	28388.000000	28388.000000	28382.000000
mean	609.197097	405.584444	20.024600	7.393332	33.465901	17.573430	0.835909	2617.164000
std	316.521726	485.206130	8.701719	4.316396	12.466448	7.232465	0.047362	1189.840000
min	41.000000	3.546003	5.011709	0.003963	20.000000	5.000083	0.769340	66.000000
25%	334.000000	94.572138	12.435564	3.621238	20.000000	11.289710	0.769340	1621.000000
50%	607.500000	256.513966	20.124216	7.358936	30.000000	17.586134	0.861953	2630.000000
75%	884.000000	513.624224	27.544882	11.062843	50.000000	23.869033	0.876056	3597.000000
max	1195.000000	3382.034941	34.985404	14.997572	50.000000	29.998989	0.876056	5172.000000

```
In [48]: 1 missing_values = df.isnull().sum()
2
3 missing_values
```

Out[48]:

Date	0
Product	0
Quantity_Sold	6
Revenue	0
Temperature_Celsius	11
Rainfall_mm	11
Location	0
Transportation_Cost	0
Labor_Cost	0
Customer	0
Quality_Score	0
Inventory_Level	6
dtype:	int64

Check for Outliers

```
In [49]: 1 cols = ('Quantity_Sold', 'Temperature_Celsius', 'Rainfall_mm', 'Inventory_Level')
        2
        3 for col in cols:
        4     z = np.abs(stats.zscore(df[col], nan_policy='omit'))
        5     outliers = df[col][z>3]
        6     print(f'{col} has {len(outliers)} outliers')
```

Quantity_Sold has 0 outliers

Temperature_Celsius has 0 outliers

Rainfall_mm has 0 outliers

Inventory_Level has 0 outliers

Check For Skewness

```
In [50]: 1 for col in cols:
        2     skew = df[col].skew()
        3     print(f'{col} has a skewness of {skew}')
```

Quantity_Sold has a skewness of 0.007049968334628765

Temperature_Celsius has a skewness of -0.005244633606917498

Rainfall_mm has a skewness of 0.017580193058957767

Inventory_Level has a skewness of -0.006928239182168006

Handle Missing Values

```
In [51]: 1 for col in cols:
        2     df[col].fillna(df[col].mean(), inplace=True)
```

Confirm that there are no more missing numbers

```
In [52]: 1 missing_values = df.isnull().sum()  
         2  
         3 missing_values
```

```
Out[52]: Date                0  
         Product            0  
         Quantity_Sold      0  
         Revenue            0  
         Temperature_Celsius 0  
         Rainfall_mm        0  
         Location           0  
         Transportation_Cost 0  
         Labor_Cost         0  
         Customer           0  
         Quality_Score      0  
         Inventory_Level    0  
         dtype: int64
```

4. Exploratory Data Analysis (EDA)

Univariate Analysis

Analysing the distribution of variables individually, for

- numerical variables,
- categorical Variables

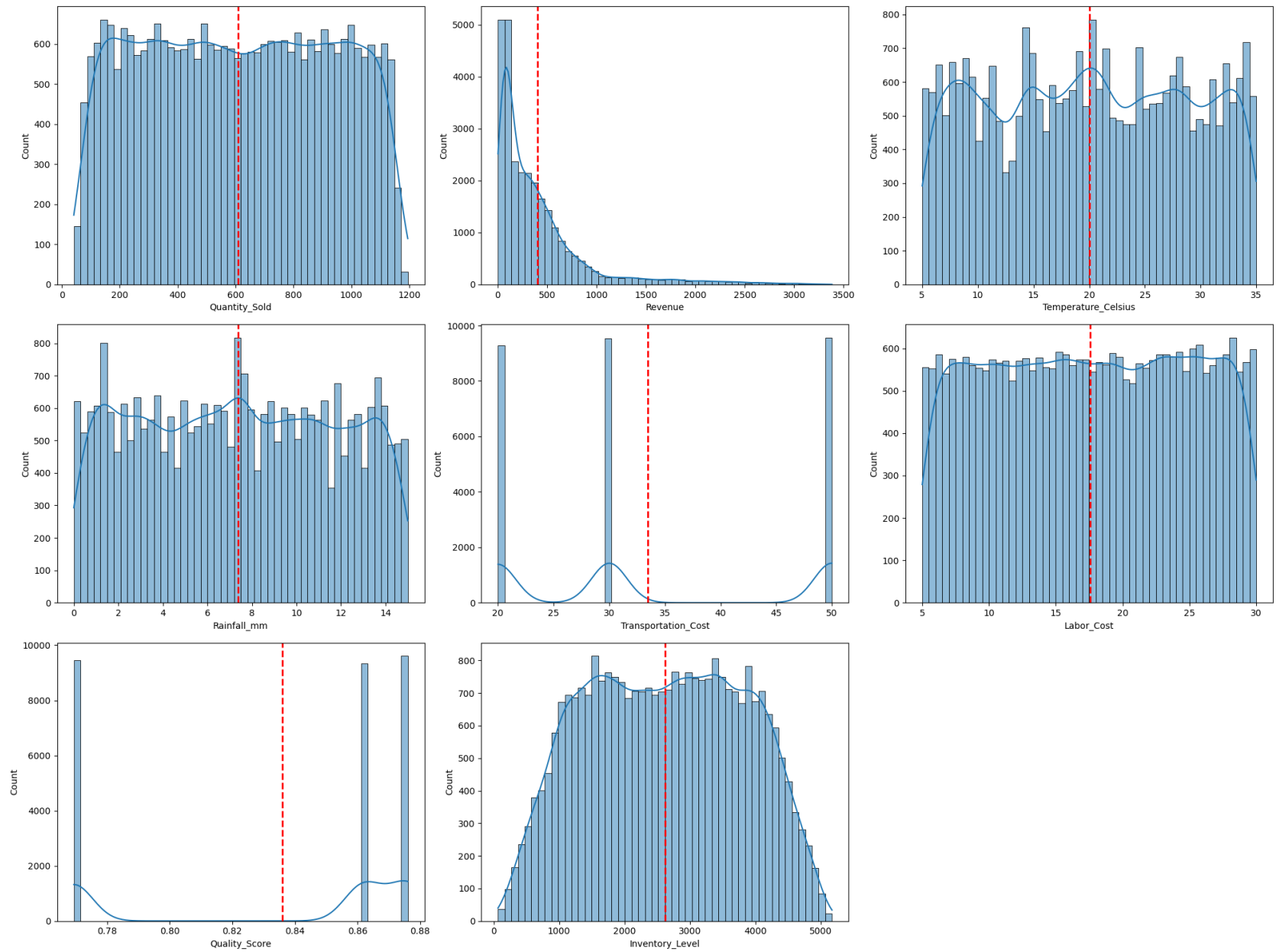
Univariate Analysis: Numerical Variables

In [53]:

1 df.info()

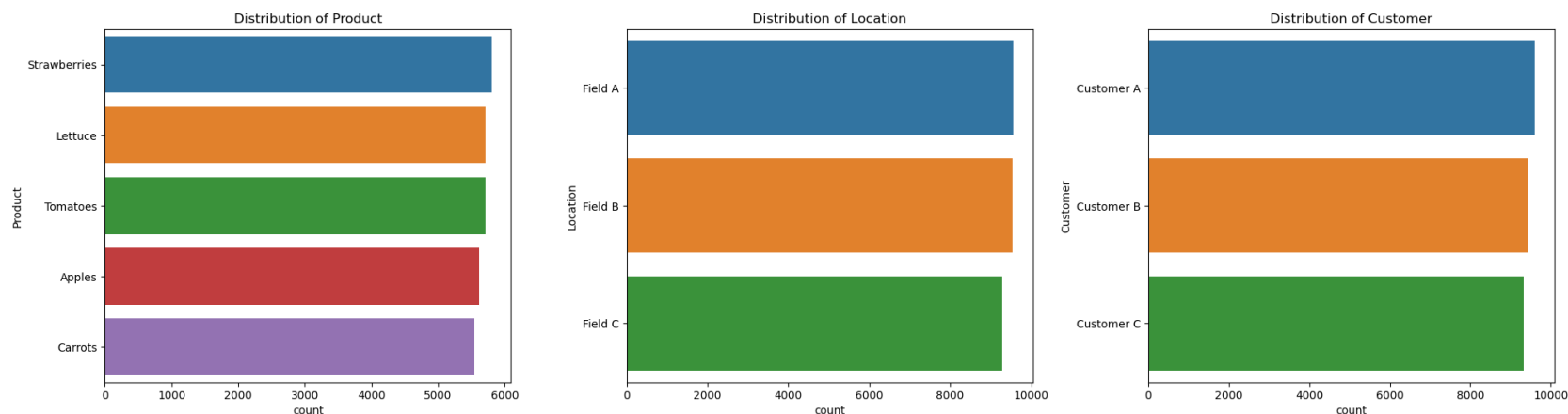
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 28388 entries, 0 to 28387
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   Date                  28388 non-null  object 
1   Product               28388 non-null  object 
2   Quantity_Sold         28388 non-null  float64
3   Revenue               28388 non-null  float64
4   Temperature_Celsius   28388 non-null  float64
5   Rainfall_mm           28388 non-null  float64
6   Location              28388 non-null  object 
7   Transportation_Cost    28388 non-null  int64  
8   Labor_Cost            28388 non-null  float64
9   Customer              28388 non-null  object 
10  Quality_Score          28388 non-null  float64
11  Inventory_Level        28388 non-null  float64
dtypes: float64(7), int64(1), object(4)
memory usage: 2.6+ MB
```

```
In [54]: 1 plt.figure(figsize=(20, 15))
          2
          3 numeric_vars = ['Quantity_Sold',
          4                     'Revenue',
          5                     'Temperature_Celsius',
          6                     'Rainfall_mm',
          7                     'Transportation_Cost',
          8                     'Labor_Cost',
          9                     'Quality_Score',
         10                     'Inventory_Level']
          11
         12 for i, var in enumerate(numeric_vars, 1):
         13     plt.subplot(3, 3, i)
         14     sns.histplot(df[var], bins=50, kde=True)
         15     plt.axvline(df[var].mean(), color='red', linestyle='dashed', linewidth =2)
         16
         17
         18 plt.tight_layout()
         19 plt.show()
```

Univariate Analysis: Categorical Variables

```
In [55]: 1 plt.figure(figsize=(20, 10))
2
3 categorical_vars = ['Product',
4                   'Location',
5                   'Customer',]
6
7 for i, var in enumerate(categorical_vars, 1):
8     plt.subplot(2, 3, i)
9     sns.countplot(data=df, y=var, order=df[var].value_counts().index)
10    plt.title(f'Distribution of {var}')
11
12 plt.tight_layout()
13 plt.show()
```



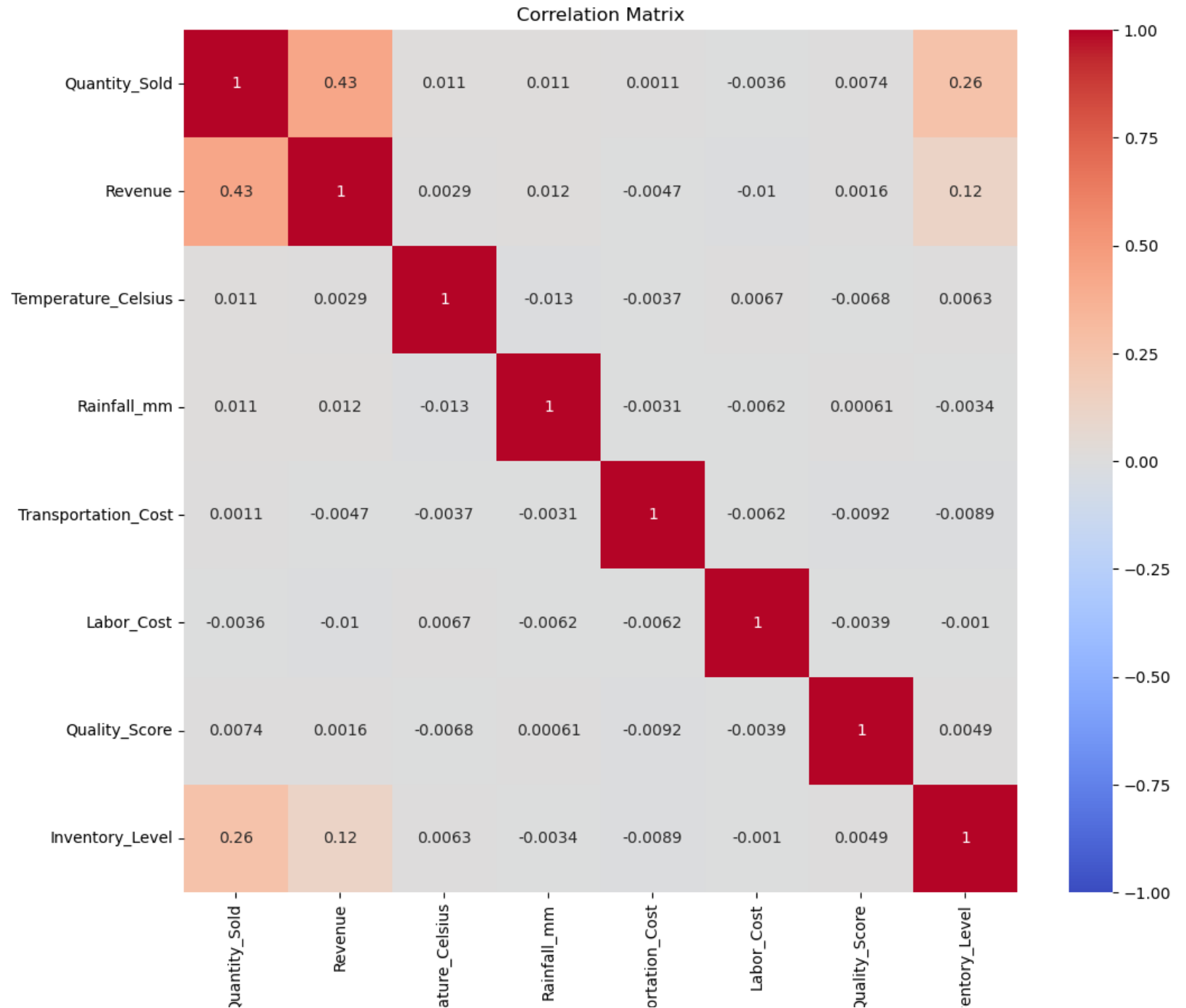
Bivariate Analysis

We look at pairs of variables, by looking at: - pairs of numerical variables - numerical and categorical variables - pairs of categorical variables

Bivariate Analysis: Numerical Variables

Starting with Correlation matrix

```
In [56]: 1 correlation_matrix = df.select_dtypes(include=np.number).corr()
          2
          3 plt.figure(figsize=(12, 10))
          4 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', vmin=-1, vmax=1)
          5 plt.title('Correlation Matrix')
          6 plt.show()
```

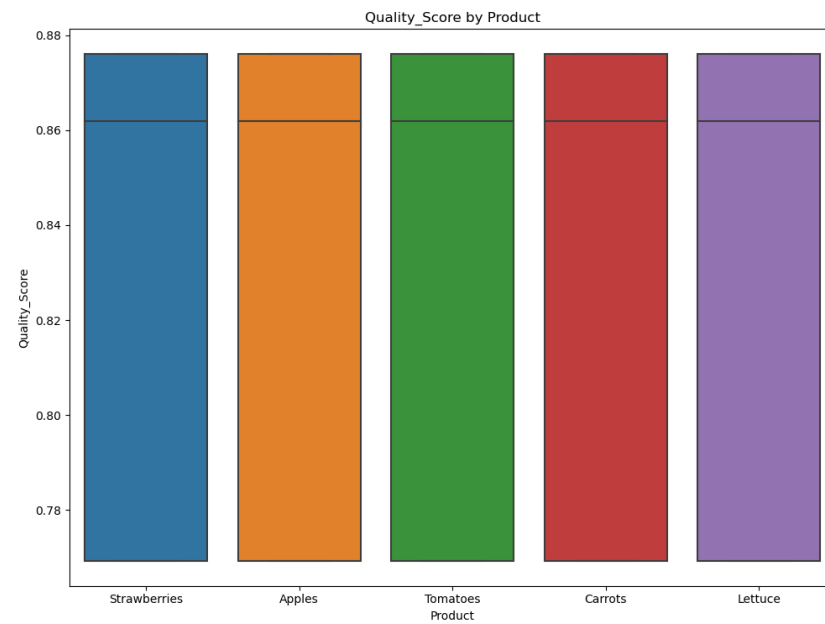
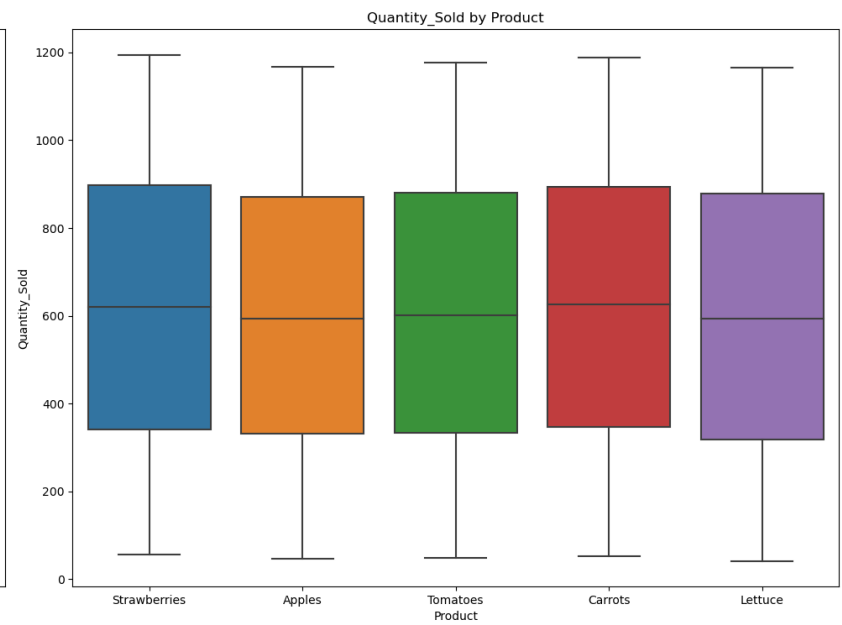
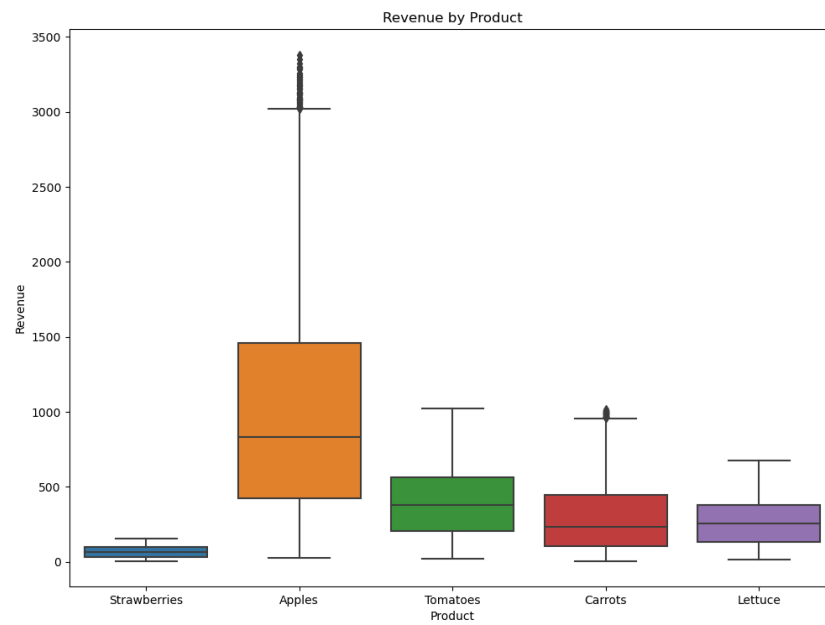


Bivariate Analysis: Numerical and Categorical Variables

Analysing 'Product' against:

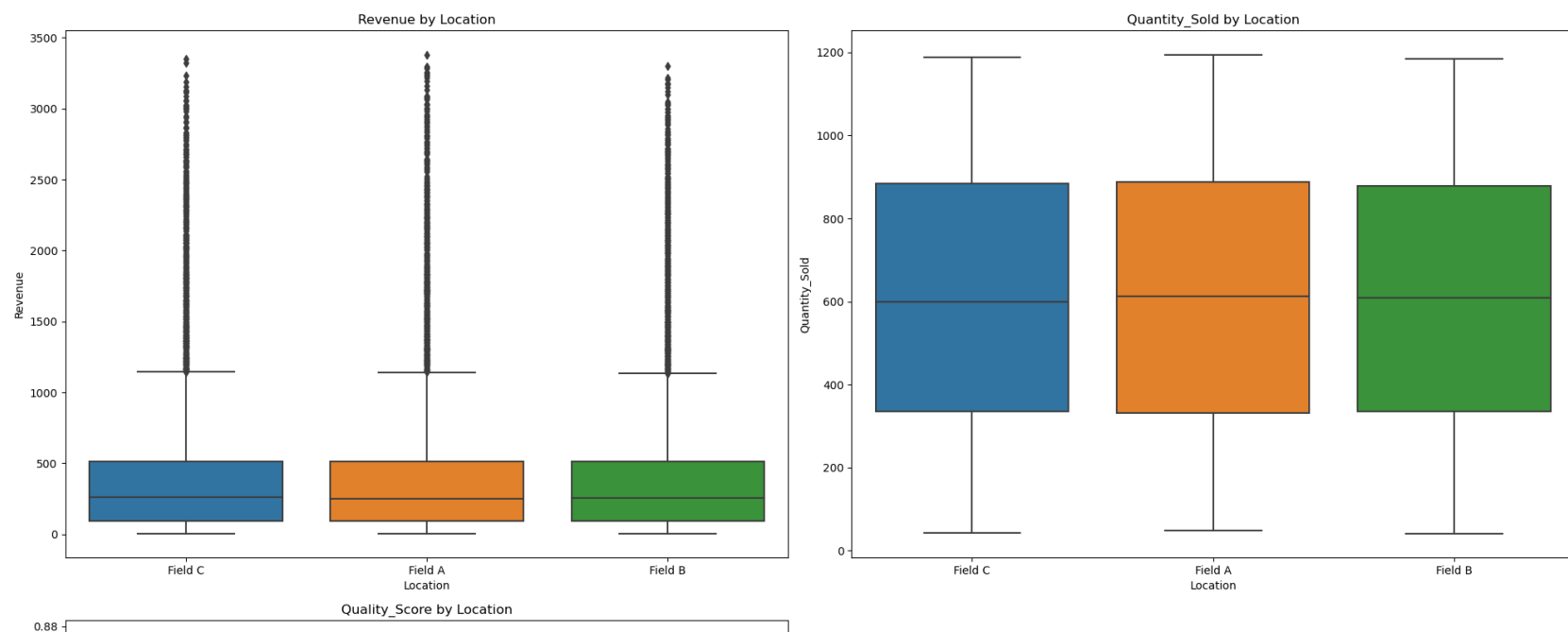
- 'Revenue'
- 'Quantity_Sold'
- 'Quality_Score'

```
In [57]: 1 plt.figure(figsize= (20, 15))
          2
          3 key_numeric_vars = ['Revenue', 'Quantity_Sold', 'Quality_Score']
          4
          5 for i,var in enumerate(key_numeric_vars,1):
          6     plt.subplot(2, 2, i)
          7     sns.boxplot(data=df, x='Product', y=var)
          8     plt.title(f'{var} by Product')
          9
         10
         11 plt.tight_layout()
         12 plt.show()
         13
```



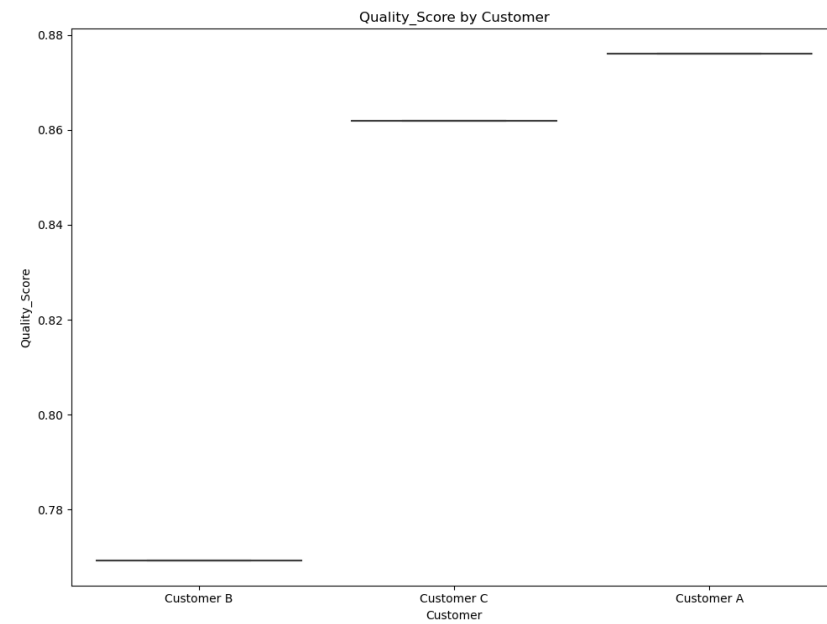
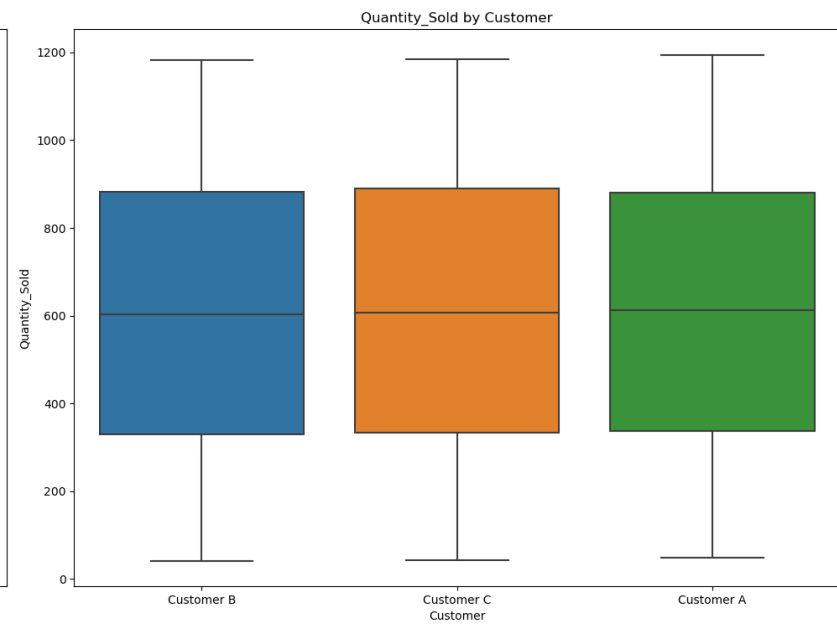
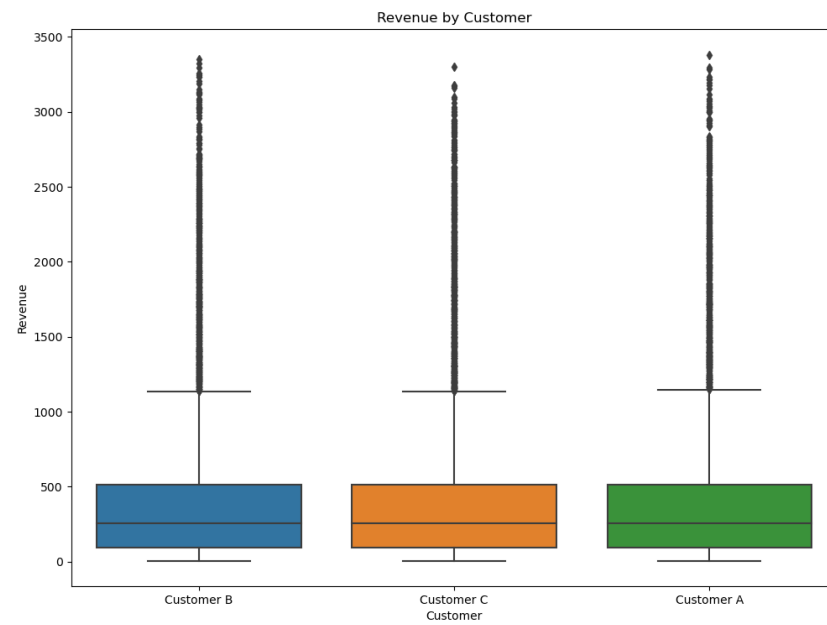
"Field" Vs the same key Variables

```
In [58]: 1 plt.figure(figsize=(20,15))
2
3 for i, var in enumerate(key_numeric_vars, 1):
4     plt.subplot(2, 2, i)
5     sns.boxplot(data=df, x='Location', y=var)
6     plt.title(f'{var} by Location')
7
8 plt.tight_layout()
9 plt.show()
10
```



Relationship between 'Customer' and the key parameters


```
In [59]: 1 plt.figure(figsize= (20, 15))
          2
          3 for i,var in enumerate(key_numeric_vars,1):
          4     plt.subplot(2, 2, i)
          5     sns.boxplot(data=df, x='Customer', y=var)
          6     plt.title(f'{var} by Customer')
          7
          8
          9 plt.tight_layout()
         10 plt.show()
         11
```



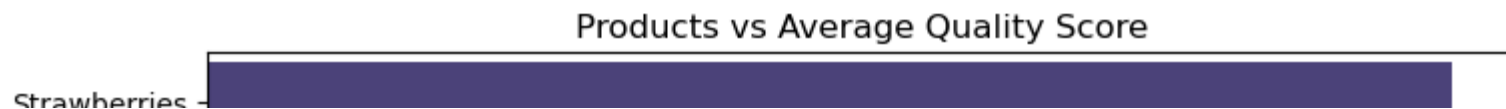
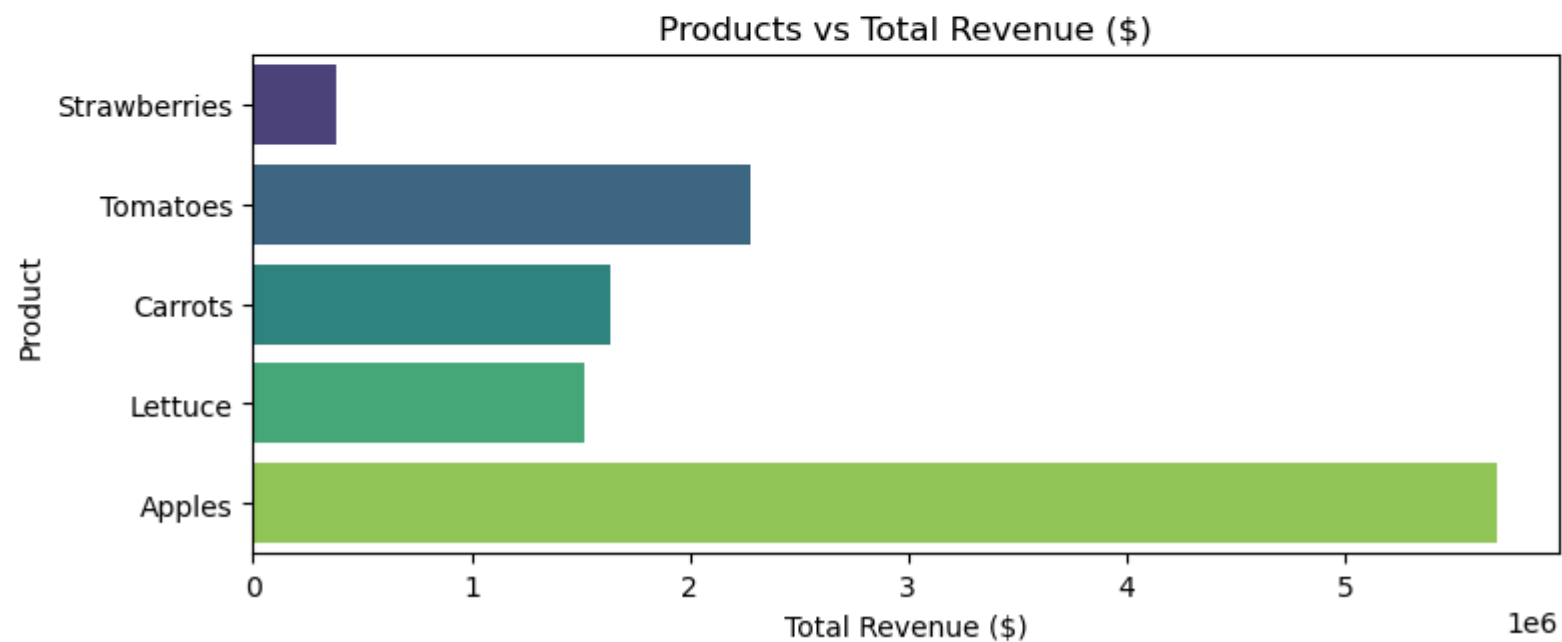
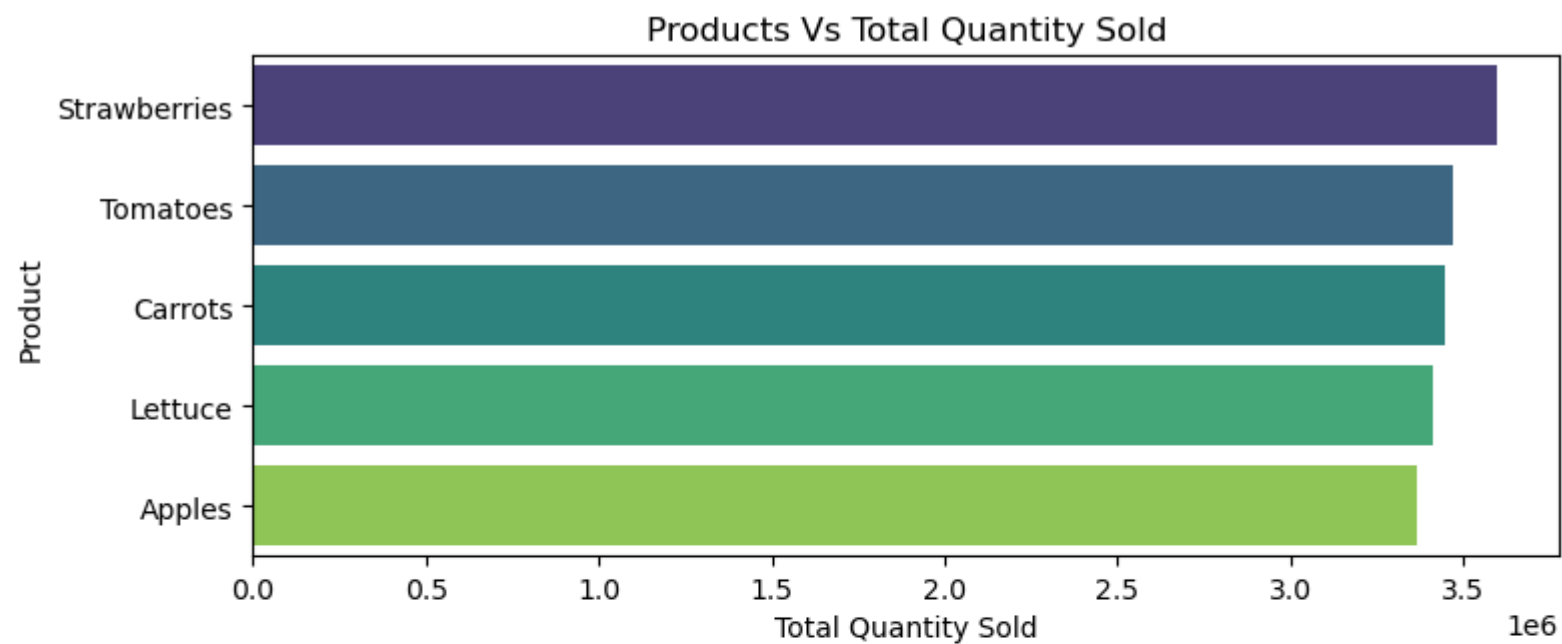
Visualise

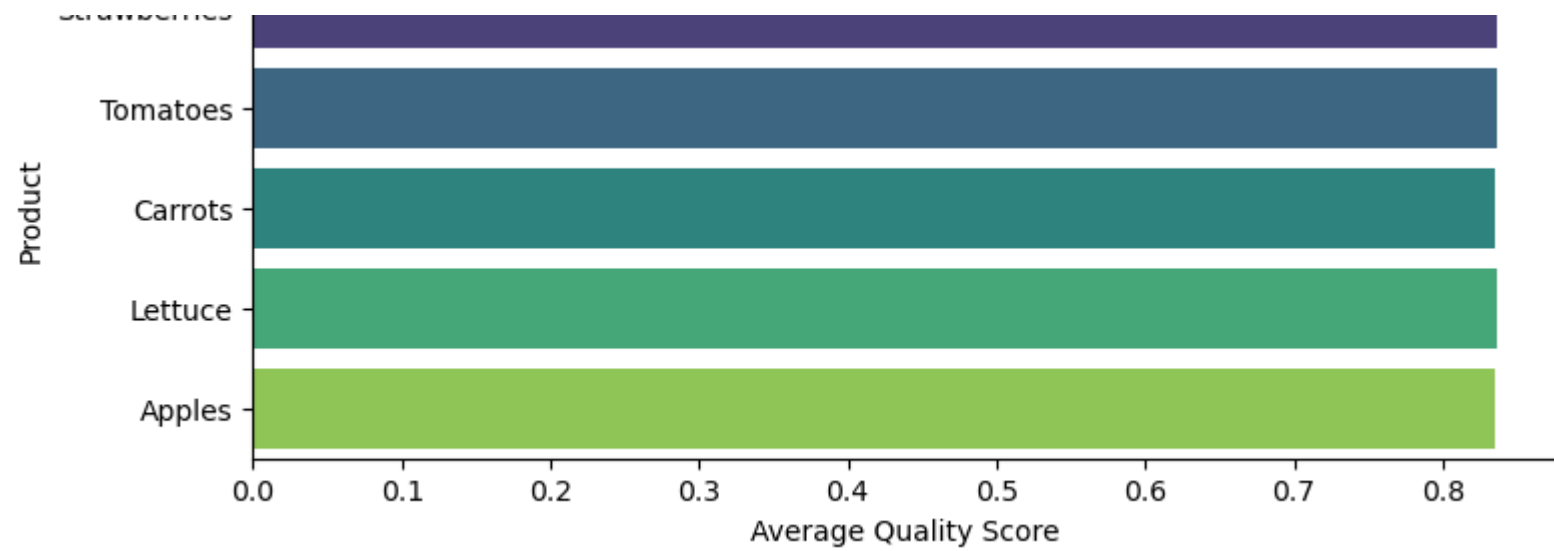
- Total 'Quantity_Sold' per product,

- Total 'Revenue' per product
- Quality_Score per product

In [60]:

```
1 # Aggregate data by product
2 product_data = df.groupby('Product').agg({'Quantity_Sold': 'sum',
3                                           'Revenue': 'sum',
4                                           'Quality_Score': 'mean'}).reset_index()
5
6 # Sort the product data by total quantity sold
7 product_data = product_data.sort_values(by='Quantity_Sold', ascending=False)
8
9 # Create subplots
10 fig, ax = plt.subplots(3, 1, figsize=(8, 10))
11
12 # Products vs Total Quantity Sold
13 sns.barplot(data=product_data, x='Quantity_Sold', y='Product', ax=ax[0], palette='viridis')
14 ax[0].set_title('Products Vs Total Quantity Sold')
15 ax[0].set_xlabel('Total Quantity Sold')
16 ax[0].set_ylabel('Product')
17
18 # Products vs Total Revenue
19 sns.barplot(data=product_data, x='Revenue', y='Product', ax=ax[1], palette='viridis')
20 ax[1].set_title('Products vs Total Revenue ($)')
21 ax[1].set_xlabel('Total Revenue ($)')
22 ax[1].set_ylabel('Product')
23
24 # Products vs Average Quality Score
25 sns.barplot(data=product_data, x='Quality_Score', y='Product', ax=ax[2], palette='viridis')
26 ax[2].set_title('Products vs Average Quality Score')
27 ax[2].set_xlabel('Average Quality Score')
28 ax[2].set_ylabel('Product')
29
30 # Adjust layout for better visualization
31 plt.tight_layout()
32 plt.show()
```



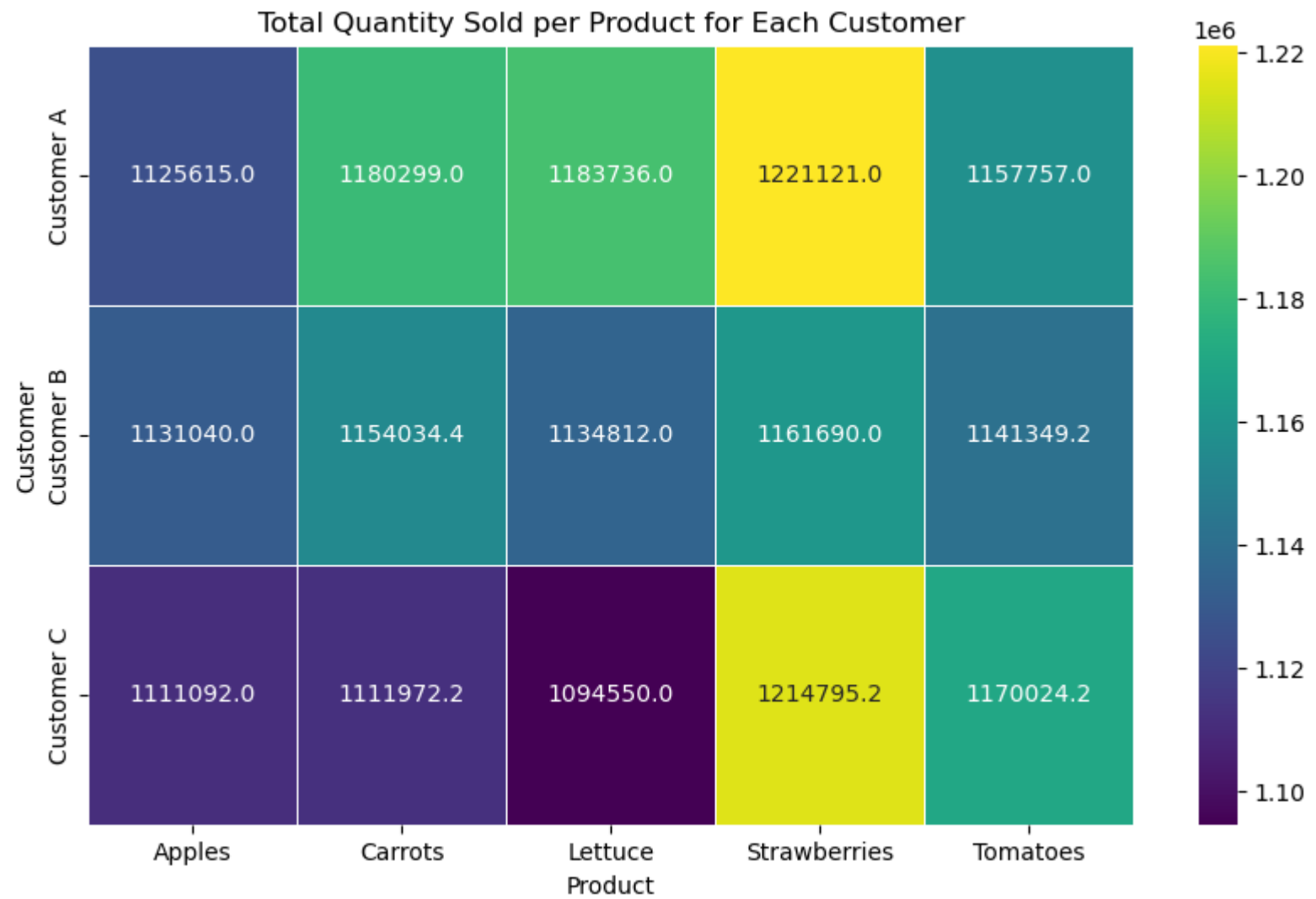


In []:

Total Quantity sold for each product per customer.

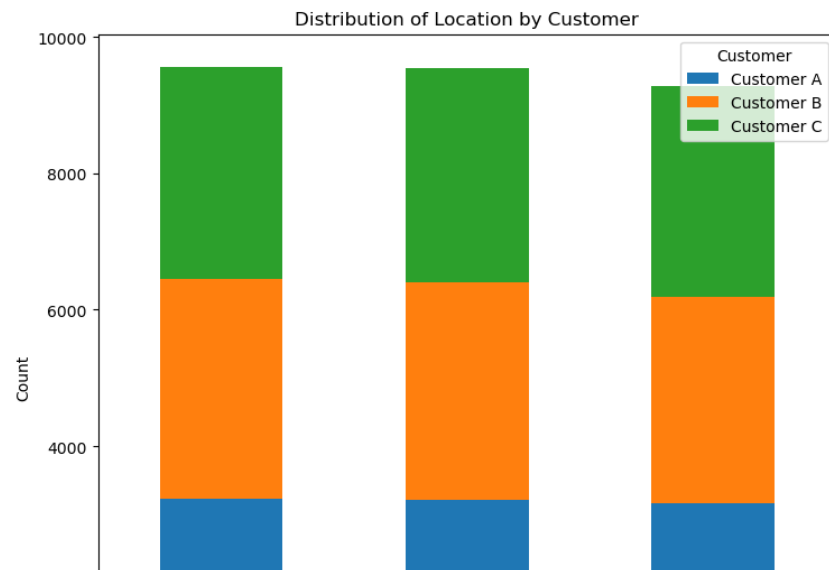
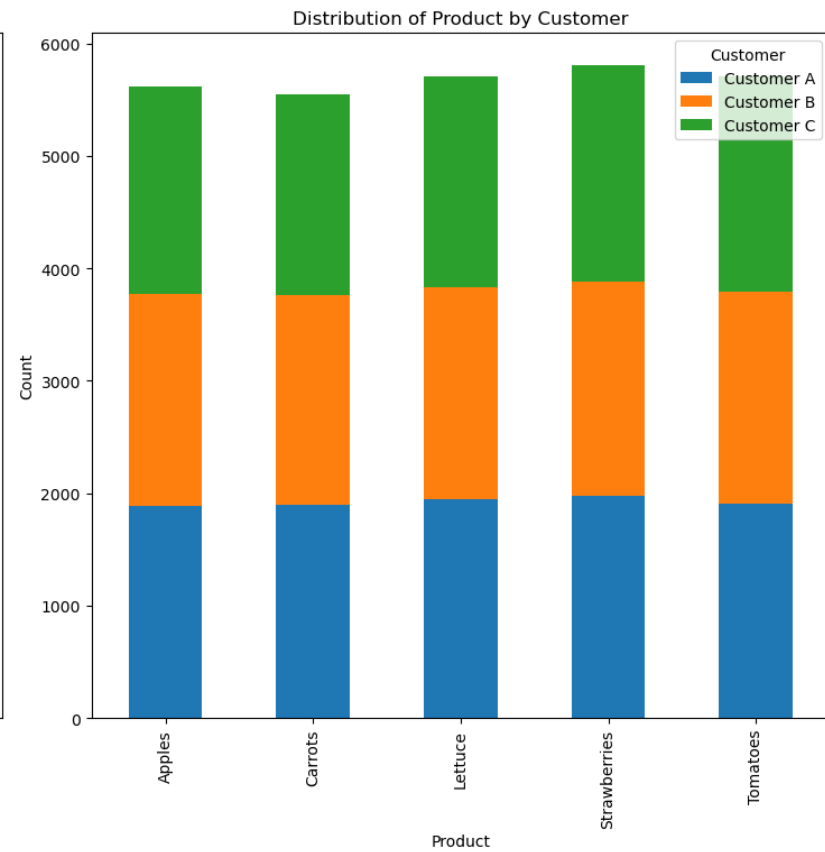
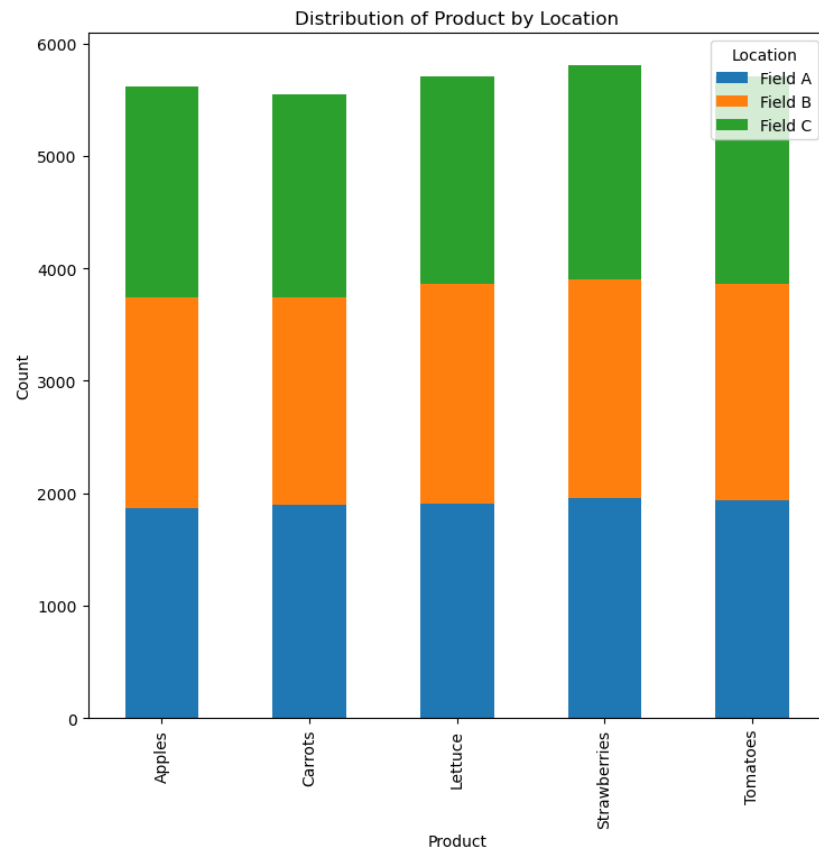
In [61]:

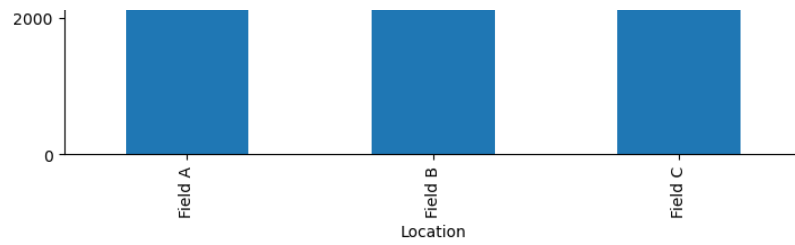
```
1 customer_data = df.groupby(['Customer', 'Product']).agg({'Quantity_Sold': 'sum',
2                                                         'Revenue': 'sum',
3                                                         'Quality_Score': 'mean'}).reset_index()
4 customer_pivot = customer_data.pivot(index='Customer', columns='Product', values='Quantity_Sold')
5
6 plt.figure(figsize=(10, 6))
7 sns.heatmap(customer_pivot, annot=True, cmap='viridis', fmt='.1f', linewidth=.5)
8 plt.title('Total Quantity Sold per Product for Each Customer')
9 plt.ylabel('Customer')
10 plt.xlabel('Product')
11 plt.show()
```



Bivariate Analysis: Two Categorical Variables


```
In [62]: 1 plt.figure(figsize=(15, 15))
          2
          3 pairs = [('Product', 'Location'),
          4               ('Product', 'Customer'),
          5               ('Location', 'Customer')]
          6 for i, (p1, p2) in enumerate(pairs):
          7     contingency_table = pd.crosstab(df[p1], df[p2])
          8
          9     ax = plt.subplot(2, 2, i + 1)
         10     contingency_table.plot(kind='bar', stacked=True, ax=ax)
         11     plt.title(f'Distribution of {p1} by {p2}')
         12     plt.ylabel('Count')
         13     plt.xlabel(p1)
         14
         15 plt.tight_layout()
         16 plt.show()
```

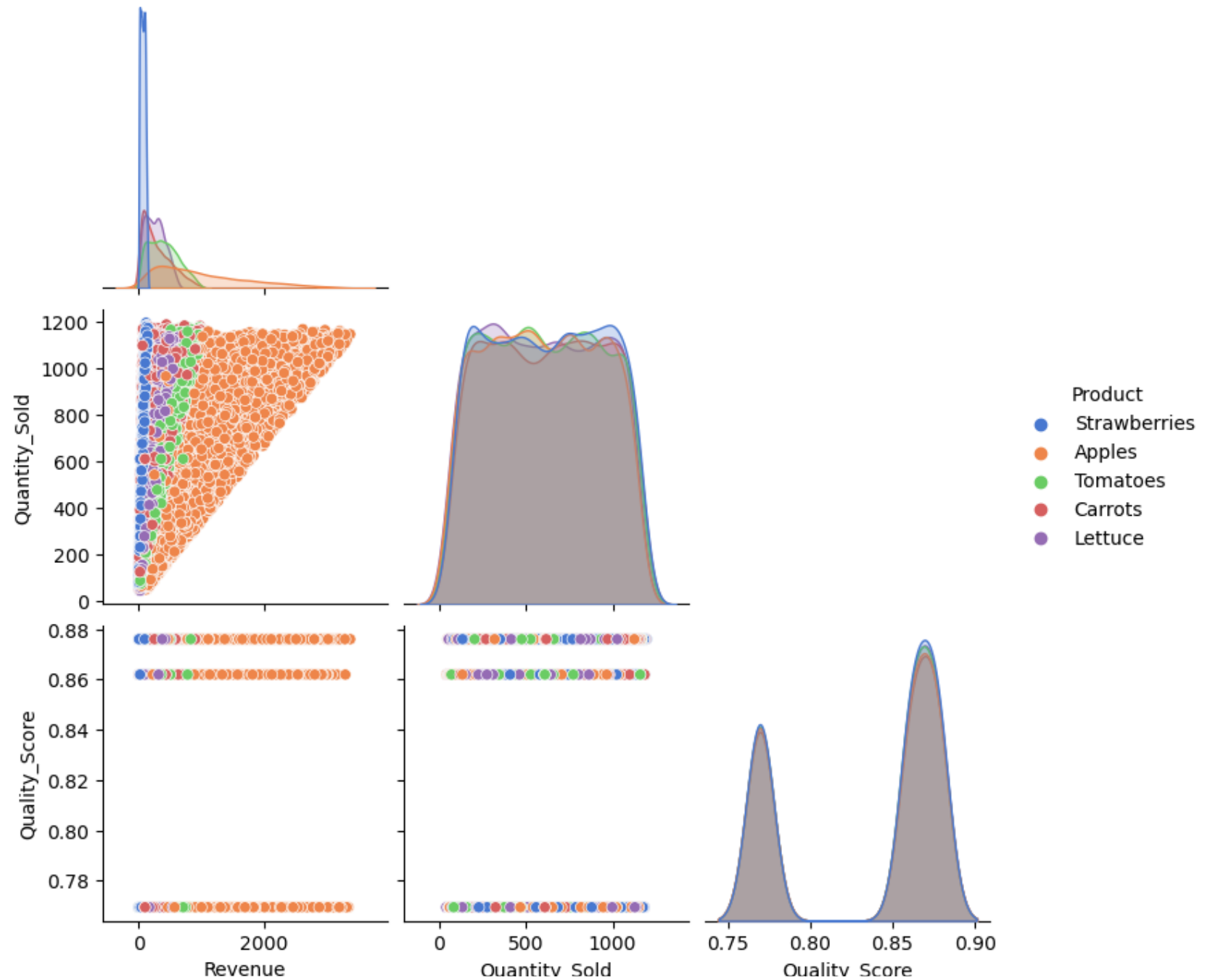




Multivariate Analysis

```
In [63]: 1 pair_plot = sns.pairplot(df, hue='Product', vars=key_numeric_vars, palette='muted', corner=True)
          2 pair_plot.fig.suptitle('Pair Plot of Key Numeric Variables by Product', y=1.02)
          3
          4 plt.show()
```

Pair Plot of Key Numeric Variables by Product

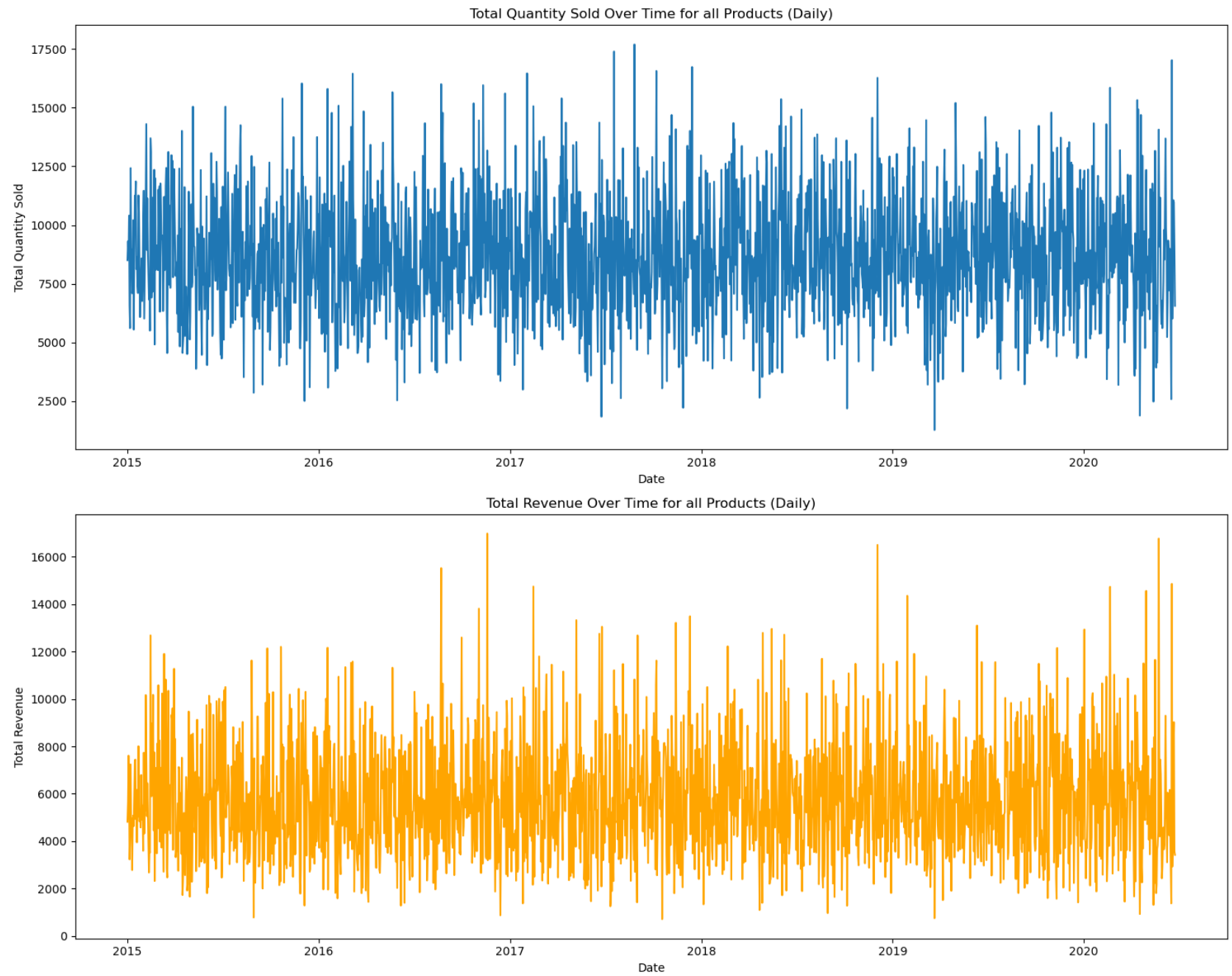


Temporal Analysis

How the Quantity_Sold and Revenue vary daily

In [64]:

```
1
2 # Convert 'Date' column to datetime
3 df['Date'] = pd.to_datetime(df['Date'])
4
5 # Group by 'Date' and aggregate 'Quantity_Sold' and 'Revenue'
6 time_based_data = df.groupby('Date').agg({'Quantity_Sold': 'sum', 'Revenue': 'sum'}).reset_index()
7
8 # Plotting the aggregated data
9 fig, ax = plt.subplots(2, 1, figsize=(15, 12))
10
11 # Quantity Sold Over Time
12 sns.lineplot(data=time_based_data, x='Date', y='Quantity_Sold', ax=ax[0])
13 ax[0].set_title('Total Quantity Sold Over Time for all Products (Daily)')
14 ax[0].set_ylabel('Total Quantity Sold')
15
16 # Revenue over Time
17 sns.lineplot(data=time_based_data, x='Date', y='Revenue', ax=ax[1], color='orange')
18 ax[1].set_title('Total Revenue Over Time for all Products (Daily)')
19 ax[1].set_ylabel('Total Revenue')
20
21 plt.tight_layout()
22 plt.show()
```



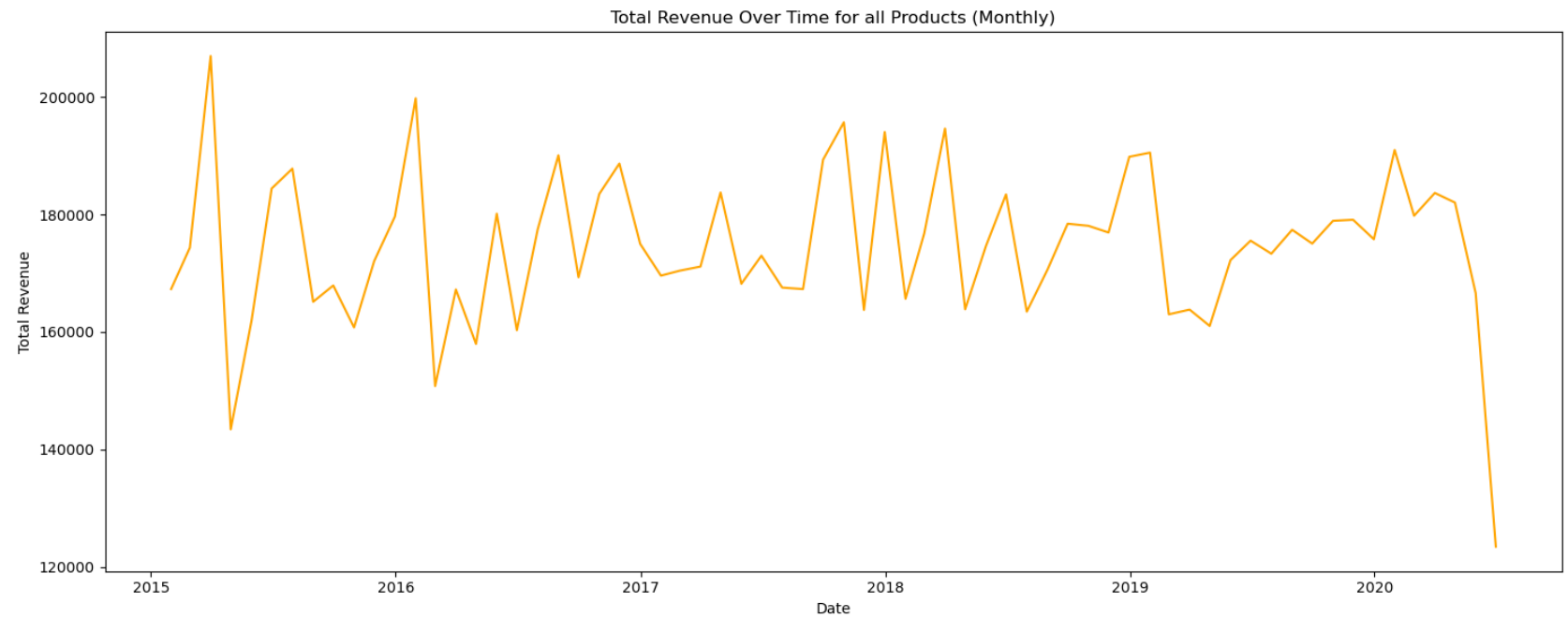
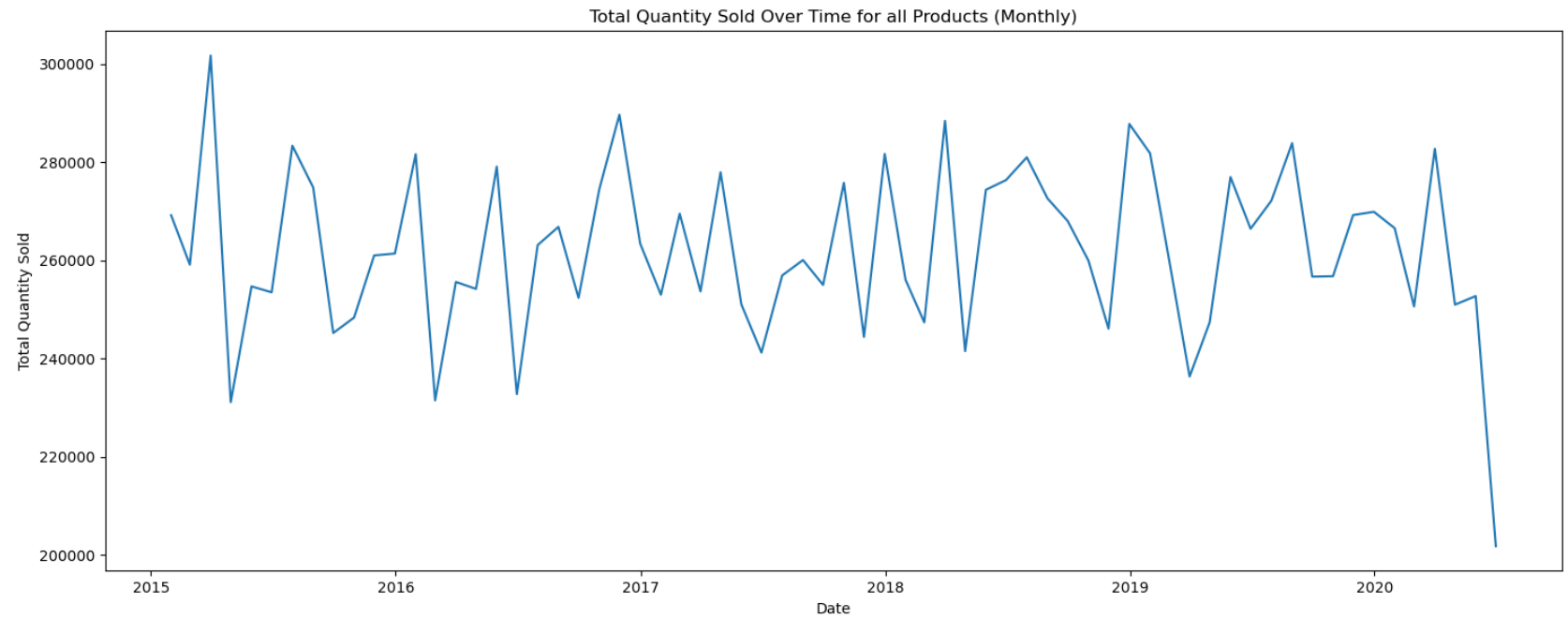
Make the Same Plot weekly

```
In [65]: 1 time_based_data = df.groupby('Date').agg({'Quantity_Sold': 'sum',
2                                             'Revenue': 'sum'}).reset_index()
3
4 time_based_data = time_based_data.set_index('Date').resample('W').sum().reset_index()
5
6 fig, ax = plt.subplots(2, 1, figsize=(15,12))
7
8 #Quantity Sold over Time
9 sns.lineplot(data=time_based_data, x='Date', y='Quantity_Sold', ax=ax[0])
10 ax[0].set_title('Total Quantity Sold Over Time for all Products (Weekly)')
11 ax[0].set_ylabel('Total Quantity Sold')
12
13
14 # Revenue over Time
15 sns.lineplot(data=time_based_data, x='Date', y='Revenue', ax=ax[1], color='orange')
16 ax[1].set_title('Total Revenue Over Time for all Products (Weekly)')
17 ax[1].set_ylabel('Total Revenue')
18
19 plt.tight_layout()
20 plt.show()
21
```



Same plot Monthly

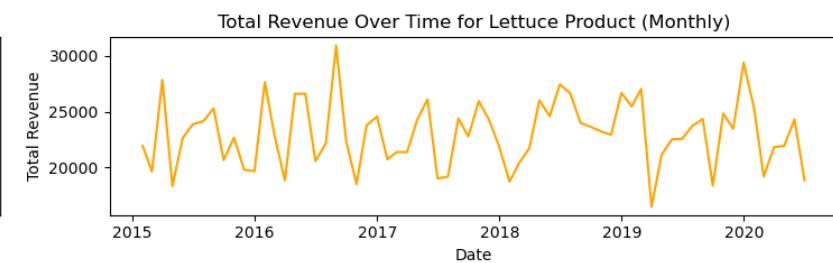
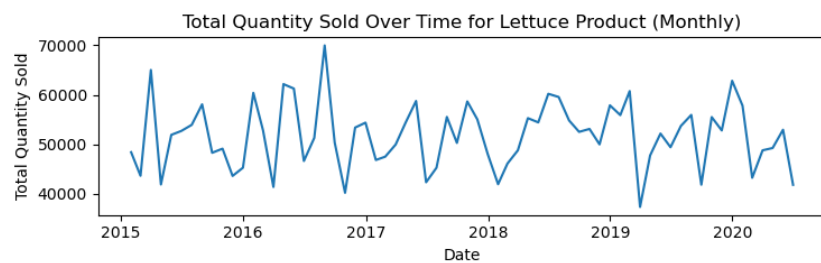
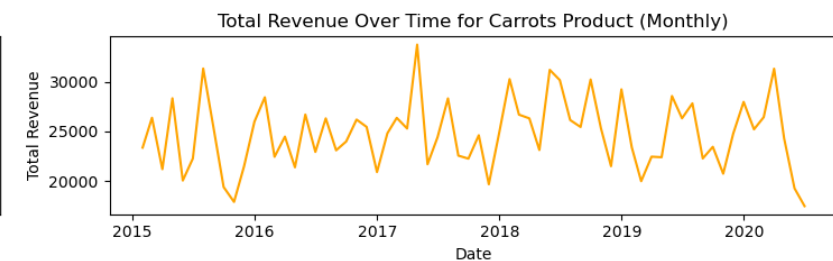
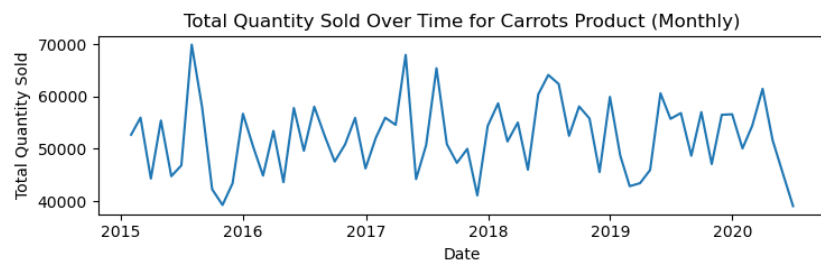
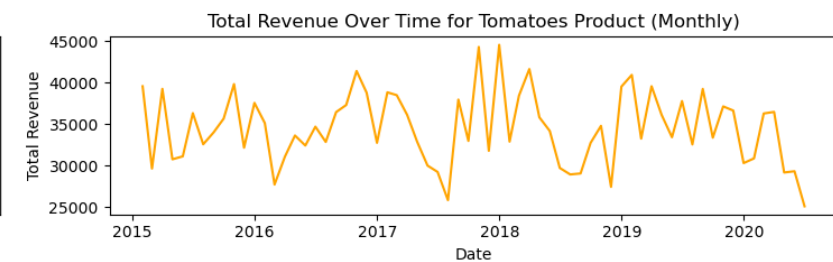
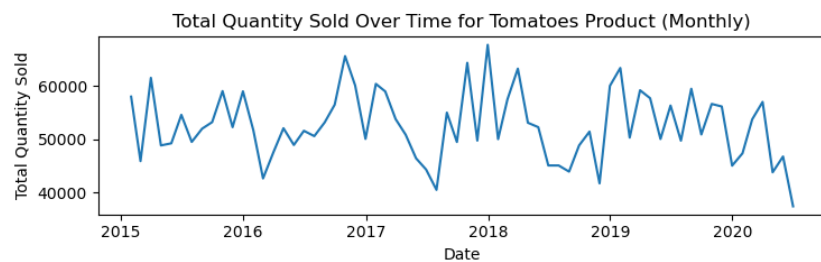
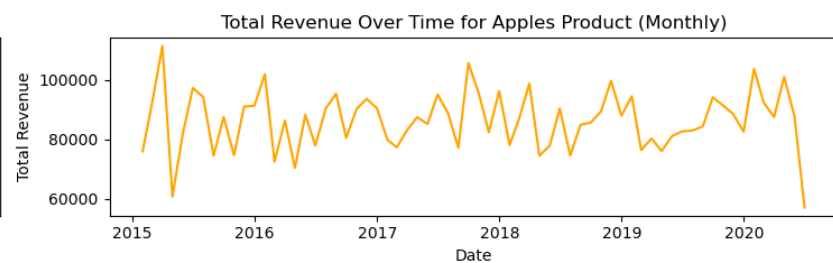
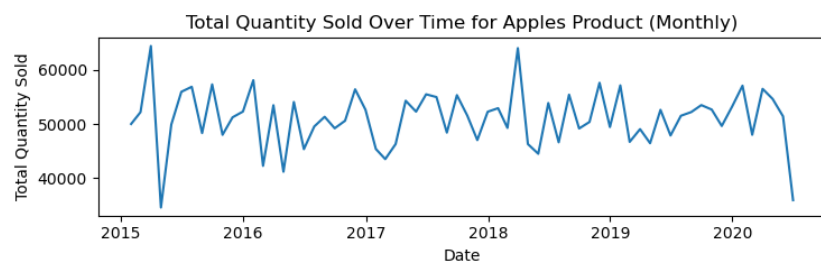
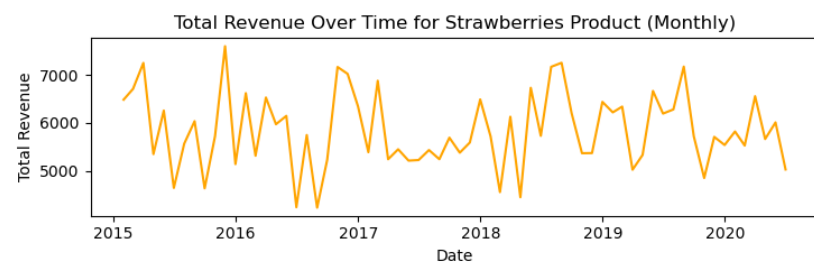
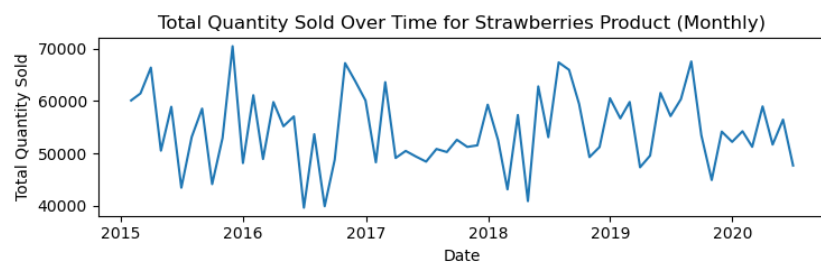
```
In [66]: 1 time_based_data = df.groupby('Date').agg({'Quantity_Sold': 'sum',
2                                               'Revenue': 'sum'}).reset_index()
3
4 time_based_data = time_based_data.set_index('Date').resample('M').sum().reset_index()
5
6 fig, ax = plt.subplots(2, 1, figsize=(15,12))
7
8 #Quantity Sold over Time
9 sns.lineplot(data=time_based_data, x='Date', y='Quantity_Sold', ax=ax[0])
10 ax[0].set_title('Total Quantity Sold Over Time for all Products (Monthly)')
11 ax[0].set_ylabel('Total Quantity Sold')
12
13
14 # Revenue over Time
15 sns.lineplot(data=time_based_data, x='Date', y='Revenue', ax=ax[1], color='orange')
16 ax[1].set_title('Total Revenue Over Time for all Products (Monthly)')
17 ax[1].set_ylabel('Total Revenue')
18
19 plt.tight_layout()
20 plt.show()
21
```



PLot qunatity sold and revenue per product

In [67]:

```
1 products = df['Product'].unique()
2
3
4 fig, ax = plt.subplots(len(products), 2, figsize=(15, 12))
5
6 for i, product in enumerate(products):
7     time_based_data = df[df['Product'] == product]
8
9     time_based_data = time_based_data.groupby('Date').agg({'Quantity_Sold': 'sum', 'Revenue': 'sum'}).res
10
11     time_based_data = time_based_data.set_index('Date').resample('M').sum().reset_index()
12
13     # Quantity Sold over Time
14     sns.lineplot(data=time_based_data, x='Date', y='Quantity_Sold', ax=ax[i, 0])
15     ax[i, 0].set_title(f'Total Quantity Sold Over Time for {product} Product (Monthly)')
16     ax[i, 0].set_ylabel('Total Quantity Sold')
17
18     # Revenue over Time
19     sns.lineplot(data=time_based_data, x='Date', y='Revenue', ax=ax[i, 1], color='orange')
20     ax[i, 1].set_title(f'Total Revenue Over Time for {product} Product (Monthly)')
21     ax[i, 1].set_ylabel('Total Revenue')
22
23 plt.tight_layout()
24 plt.show()
```

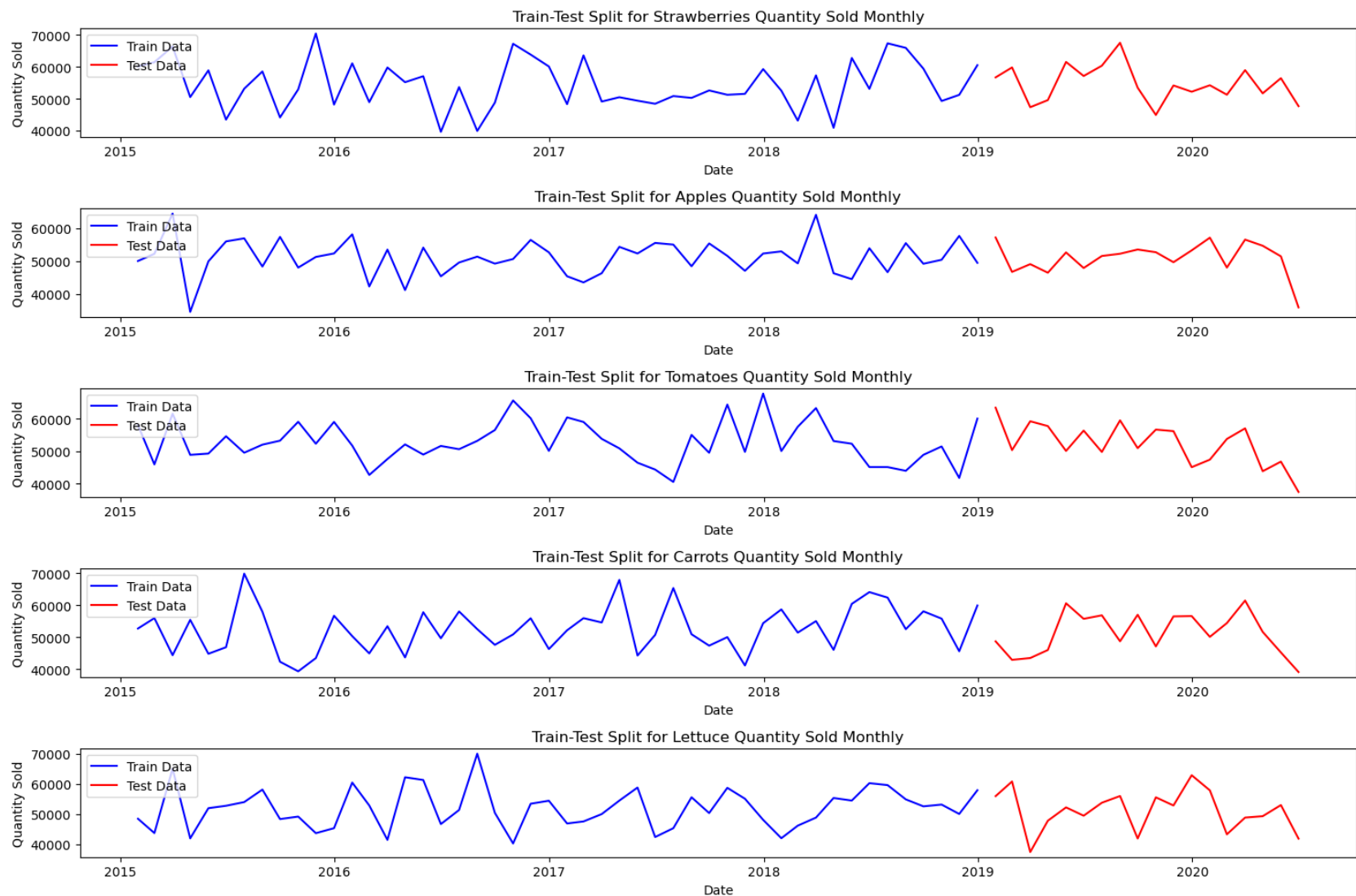
Feature Engineering

Splitting Dataset into train And test sets

```
In [68]: 1 train_data = dict()
          2 test_data = dict()
          3
          4 for product in products:
          5     df_product = df[df['Product'] == product]
          6
          7     df_product = df_product.groupby('Date').agg({'Quantity_Sold': 'sum',
          8                                             'Revenue': 'sum'}).reset_index()
          9
          10    df_product = df_product.set_index('Date').resample('M').sum().reset_index()
          11
          12    train_data[product] = df_product[df_product['Date'].dt.year < 2019].reset_index()
          13    test_data[product] = df_product[df_product['Date'].dt.year >= 2019].reset_index()
          14
```

Visualize the split

```
In [69]: 1 fig, ax = plt.subplots(len(products), 1, figsize=(15, 10))
2
3 for i, product in enumerate(products):
4     ax[i].plot(train_data[product]['Date'],
5               train_data[product]['Quantity_Sold'],
6               color='blue')
7
8     ax[i].plot(test_data[product]['Date'],
9               test_data[product]['Quantity_Sold'],
10              color='red')
11     ax[i].set_title(f'Train-Test Split for {product} Quantity Sold Monthly')
12     ax[i].set_xlabel('Date')
13     ax[i].set_ylabel('Quantity Sold')
14     ax[i].legend(['Train Data', 'Test Data'], loc='upper left')
15
16 plt.tight_layout()
17 plt.show()
18
```



5. Modelling

The models that will be used include: - ARIMA - Prophet, - RandomForestRegressor.

For evaluation, we will be using:

- Mean Absolute Error (MAE),
- Root Mean Square Error (RMSE).

Forecasting with ARIMA

Stationarity test

```
In [70]: 1 def adf_test(series):
2         result = adfuller(series)
3         dfoutput = pd.Series(result[0:4], index=['Test Statistic',
4         'p-value',
5         '#lags used',
6         'number of observations'])
7
8         for key, value in result[4].items():
9             dfoutput[f'Critical Value ({key})'] = value
10
11         return dfoutput
12
```

Performing the Sationarity test.

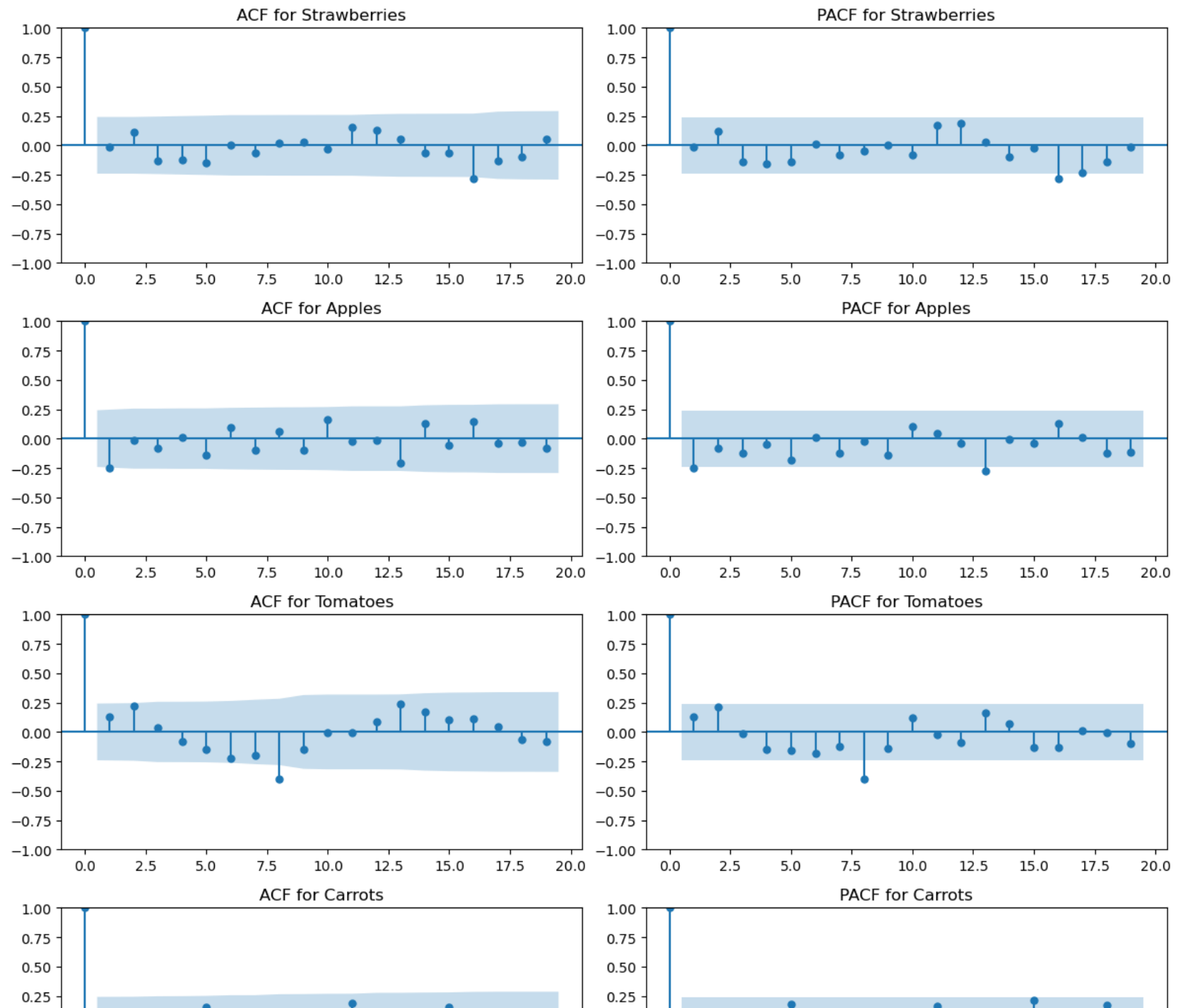
```
In [71]: 1 stationarity_test = dict()
2
3 for product in products:
4     product_data = pd.concat([train_data[product], test_data[product]])
5     stationarity_test[product] = adf_test(product_data['Quantity_Sold'])
6
7 pd.DataFrame(stationarity_test)
```

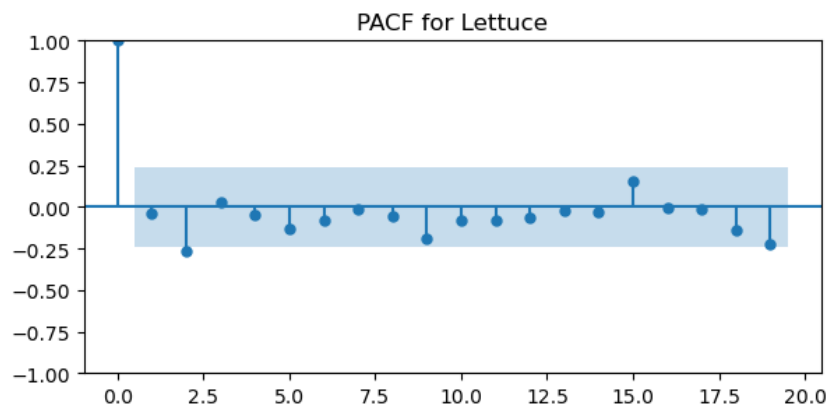
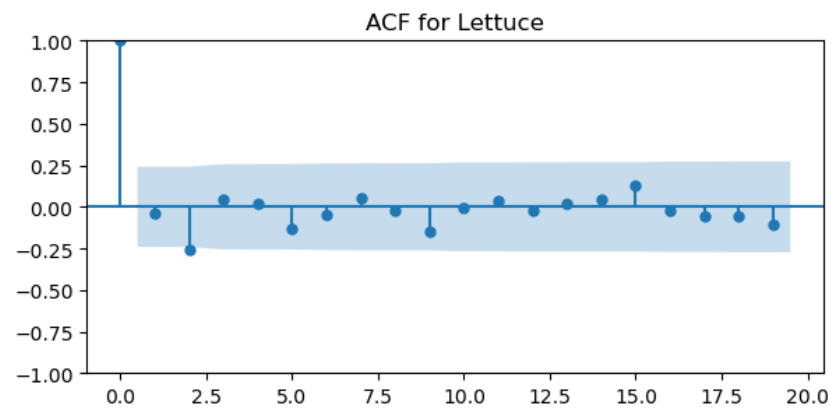
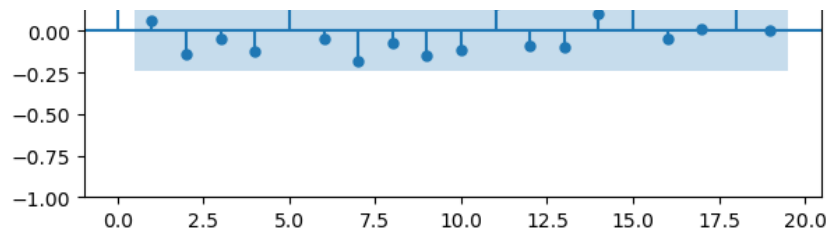
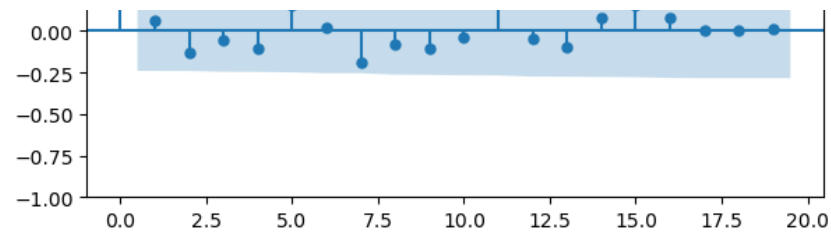
Out[71]:

	Strawberries	Apples	Tomatoes	Carrots	Lettuce
Test Statistic	-5.058980	-9.876209e+00	-4.368359	-7.209263e+00	-7.433840e+00
p-value	0.000017	3.876435e-17	0.000338	2.257203e-10	6.258368e-11
#lags used	2.000000	0.000000e+00	8.000000	0.000000e+00	1.000000e+00
number of observations	63.000000	6.500000e+01	57.000000	6.500000e+01	6.400000e+01
Critical Value (1%)	-3.538695	-3.535217e+00	-3.550670	-3.535217e+00	-3.536928e+00
Critical Value (5%)	-2.908645	-2.907154e+00	-2.913766	-2.907154e+00	-2.907887e+00
Critical Value (10%)	-2.591897	-2.591103e+00	-2.594624	-2.591103e+00	-2.591493e+00

Plot ACF and PACF

```
In [72]: 1 fig, ax = plt.subplots(len(products), 2, figsize=(12, 15))
          2
          3 for i, product in enumerate(products):
          4     product_data = pd.concat([train_data[product], test_data[product]])
          5
          6     plot_acf(product_data['Quantity_Sold'], ax=ax[i][0])
          7     ax[i][0].set_title(f'ACF for {product}')
          8
          9     plot_pacf(product_data['Quantity_Sold'], ax=ax[i][1])
         10     ax[i][1].set_title(f'PACF for {product}')
         11
         12 plt.tight_layout()
         13 plt.show()
```





Build the models, and evaluate them

```
In [73]: 1 mae, rmse, mape = 0, 0, 0
2         arima_mape_scores = dict()
3
4         arima_forecasts = dict()
5
6         for product in products:
7             train = train_data[product]['Quantity_Sold'].values
8             test = test_data[product]['Quantity_Sold'].values
9
10            model = ARIMA(train, order=(1, 0, 1))
11            model_fit = model.fit()
12
13            forecast = model_fit.forecast(steps=len(test))
14            arima_forecasts[product] = forecast
15
16            product_mae = mean_absolute_error(test, forecast)
17            product_rmse = mean_squared_error(test, forecast, squared=False)
18            product_mape = mean_absolute_percentage_error(test, forecast) * 100
19
20            mae += product_mae
21            rmse += product_rmse
22            mape += product_mape
23
24            arima_mape_scores[product] = product_mape
25
26        n_products = len(products)
27
28        mae /= n_products
29        rmse /= n_products
30        mape /= n_products
31
32        print('MAE:', mae)
33        print('RMSE:', rmse)
34        print('MAPE:', mape)
```

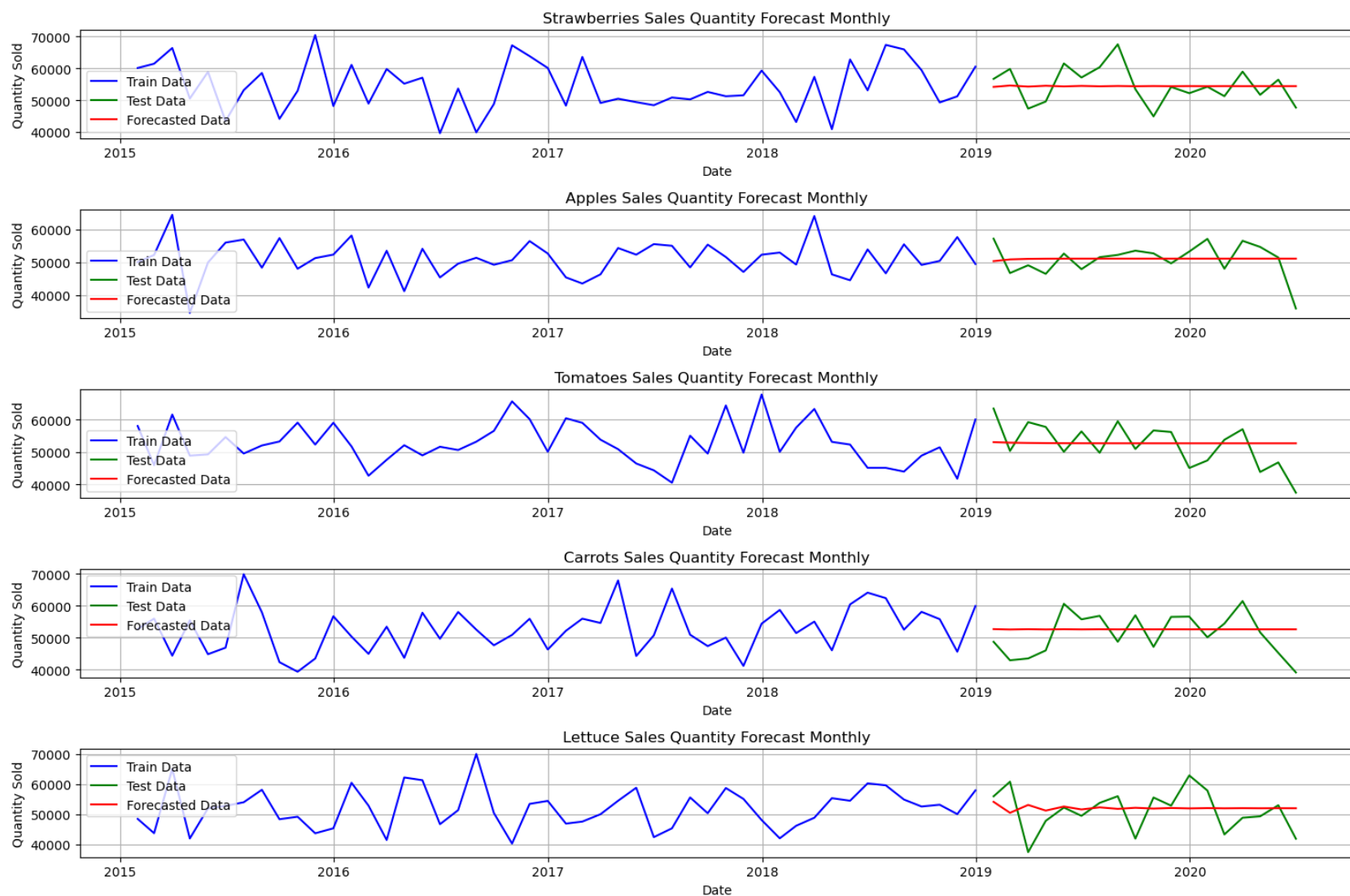
MAE: 4915.419962117331

RMSE: 6059.511862209205

MAPE: 9.971574059237227

Plot the train, test and forecast data.

```
In [74]: 1 import matplotlib.pyplot as plt
2
3 plt.figure(figsize=(15, 10))
4
5 for i, product in enumerate(products):
6     plt.subplot(len(products), 1, i + 1)
7
8     plt.plot(train_data[product]['Date'],
9              train_data[product]['Quantity_Sold'],
10             label='Actual Quantity Sold (Train Data)',
11             color='blue')
12     plt.plot(test_data[product]['Date'],
13             test_data[product]['Quantity_Sold'],
14             label='Actual Quantity Sold (Test Data)',
15             color='green')
16     plt.plot(test_data[product]['Date'],
17             arima_forecasts[product],
18             label='Forecasted Quantity to be Sold',
19             color='red',
20             linestyle='-')
21     plt.title(f'{product} Sales Quantity Forecast Monthly')
22     plt.xlabel('Date')
23     plt.ylabel('Quantity Sold')
24     plt.legend(['Train Data', 'Test Data', 'Forecasted Data'])
25     plt.grid(True)
26
27 plt.tight_layout()
28 plt.show()
29
```



Forecasting With Prophet

With prophet, Train and Evaluate the Models

In [82]:

```
1 mae, rmse, mape = 0, 0, 0
2
3
4 prophet_models = dict()
5 prophet_forecasts = dict()
6 prophet_mape_scores = dict()
7
8 for product in products:
9     train = train_data[product][['Date', 'Quantity_Sold']]
10    test = test_data[product][['Date', 'Quantity_Sold']]
11
12    train.columns = ['ds', 'y']
13    test.columns = ['ds', 'y']
14
15    model_prophet = Prophet(yearly_seasonality=True, daily_seasonality=False)
16    model_prophet.fit(train)
17    prophet_models[product] = model_prophet
18
19    future_dates = model_prophet.make_future_dataframe(periods=len(test), freq='M')
20
21    all_forecast = model_prophet.predict(future_dates)
22    prophet_forecasts[product] = all_forecast
23
24    test_forecast = all_forecast[-len(test):]
25
26    product_mae = mean_absolute_error(test['y'], test_forecast['yhat'])
27    product_rmse = mean_squared_error(test['y'], test_forecast['yhat'], squared=False)
28    product_mape = mean_absolute_percentage_error(test['y'], test_forecast['yhat']) * 100
29
30    prophet_mape_scores[product] = product_mape
31
32    mae += product_mae
33    rmse += product_rmse
34    mape += product_mape
35
36 n_products = len(products)
37
38 mae /= n_products
39 rmse /= n_products
40 mape /= n_products
41
```

```
42 print('\nMAE:', mae)
43 print('RMSE:', rmse)
44 print('MAPE:', mape)
```

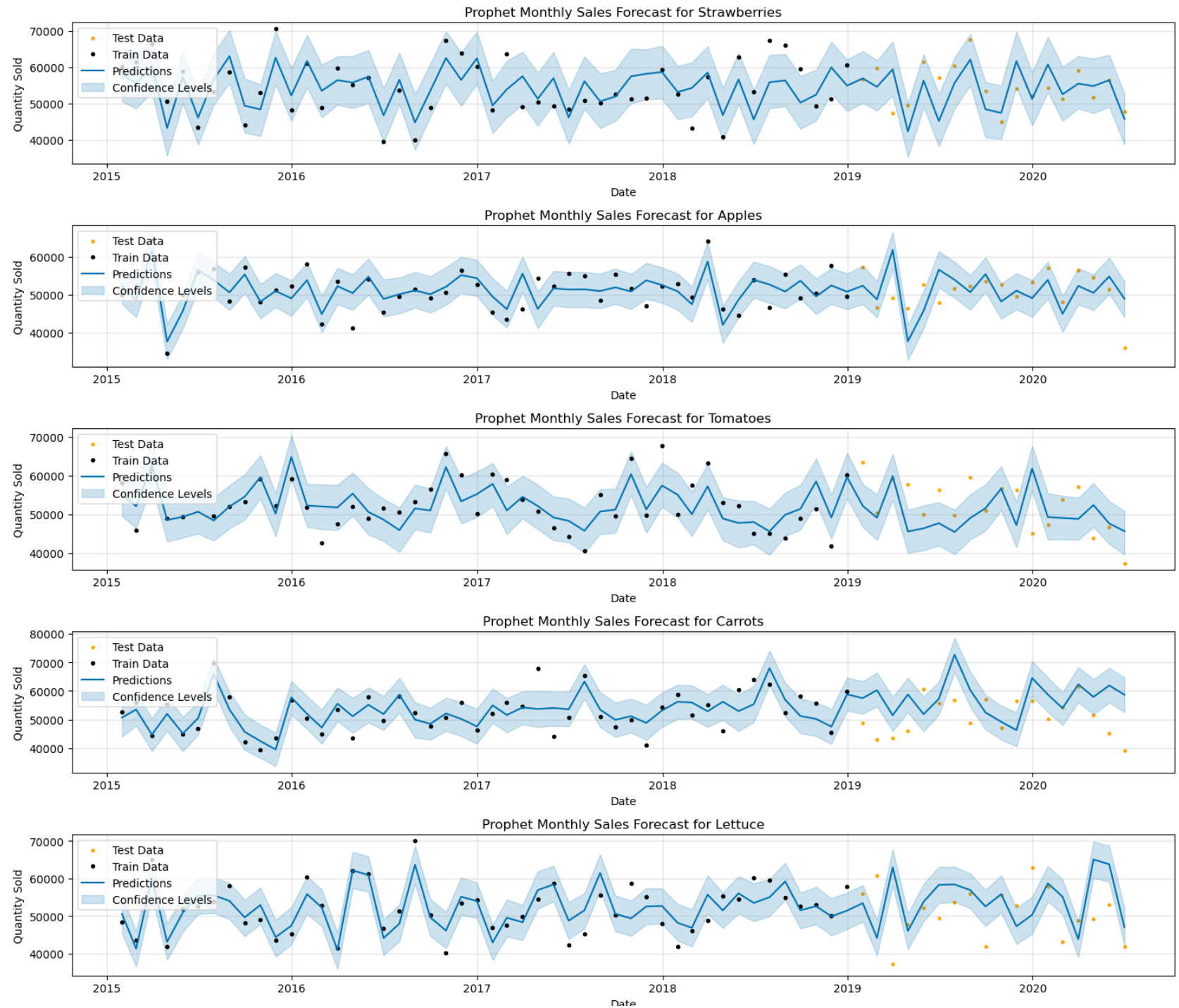
```
10:38:28 - cmdstanpy - INFO - Chain [1] start processing
10:38:29 - cmdstanpy - INFO - Chain [1] done processing
10:38:29 - cmdstanpy - INFO - Chain [1] start processing
10:38:29 - cmdstanpy - INFO - Chain [1] done processing
10:38:30 - cmdstanpy - INFO - Chain [1] start processing
10:38:30 - cmdstanpy - INFO - Chain [1] done processing
10:38:30 - cmdstanpy - INFO - Chain [1] start processing
10:38:30 - cmdstanpy - INFO - Chain [1] done processing
10:38:31 - cmdstanpy - INFO - Chain [1] start processing
10:38:31 - cmdstanpy - INFO - Chain [1] done processing
```

```
MAE: 6557.378688869758
RMSE: 8136.562065696507
MAPE: 13.280854877120092
```

Visualization Of the Result

In [83]:

```
1
2 # Create subplots for each product
3 fig, axes = plt.subplots(len(products), 1, figsize=(15, 13))
4
5 if len(products) == 1:
6     axes = [axes] # Ensure axes is always a list
7
8 for i, product in enumerate(products):
9     axes[i].scatter(test_data[product]['Date'],
10                    test_data[product]['Quantity_Sold'],
11                    color='orange',
12                    s=6)
13
14     category_plot = prophet_models[product].plot(prophet_forecasts[product],
15                                                    ax=axes[i],
16                                                    figsize=(15, 7))
17
18     axes[i].set_title(f'Prophet Monthly Sales Forecast for {product}')
19     axes[i].set_xlabel('Date')
20     axes[i].set_ylabel('Quantity Sold')
21
22     axes[i].legend(['Test Data', 'Train Data', 'Predictions', 'Confidence Levels'], loc='upper left')
23
24 plt.tight_layout()
25 plt.show()
```



Forecasting With RandomForestRegressor

Train And Evaluate the models

In [84]:

```
1
2 # Initialize variables for evaluation metrics
3 rf_mae, rf_rmse, rf_mape = 0, 0, 0
4
5 # Dictionaries to store models and forecasts
6 rf_mape_scores = dict()
7 rf_X_train_dates = dict()
8 rf_X_test_dates = dict()
9 rf_y_train = dict()
10 rf_y_test = dict()
11 rf_forecasts = dict()
12
13 for product in products:
14     df_product = df[df['Product'] == product]
15
16     df_product = pd.get_dummies(df_product, drop_first=True)
17
18     df_product = df_product.set_index('Date').resample('M').sum().reset_index()
19
20     df_product['Year'] = df_product['Date'].dt.year
21     df_product['Month'] = df_product['Date'].dt.month
22     df_product['Day'] = df_product['Date'].dt.day
23     df_product['WeekOfYear'] = df_product['Date'].dt.isocalendar().week
24
25     # Prepare the data
26     train = df_product[df_product['Year'] < 2019].reset_index()
27     test = df_product[df_product['Year'] >= 2019].reset_index()
28
29     rf_X_train_dates[product] = train['Date']
30     rf_X_test_dates[product] = test['Date']
31
32     train = train.drop(columns=['Date'])
33     test = test.drop(columns=['Date'])
34
35     X_train = train.drop(columns=['Quantity_Sold', 'Revenue'])
36     y_train = train['Quantity_Sold']
37
38     X_test = test.drop(columns=['Quantity_Sold', 'Revenue'])
39     y_test = test['Quantity_Sold']
40
41     rf_y_train[product] = y_train
```

```
42 rf_y_test[product] = y_test
43
44 # Train the RandomForest model
45 model = RandomForestRegressor(n_estimators=100, random_state=42)
46 model.fit(X_train, y_train)
47
48 forecast = model.predict(X_test)
49 rf_forecasts[product] = forecast
50
51 product_mae = mean_absolute_error(y_test, forecast)
52 product_rmse = mean_squared_error(y_test, forecast, squared=False)
53 product_mape = mean_absolute_percentage_error(y_test, forecast) * 100
54
55 rf_mae += product_mae
56 rf_rmse += product_rmse
57 rf_mape += product_mape
58
59 rf_mape_scores[product] = product_mape
60
61 # Calculate average metrics
62 n_products = len(products)
63 rf_mae /= n_products
64 rf_rmse /= n_products
65 rf_mape /= n_products
66
67 print('Random Forest MAE:', rf_mae)
68 print('Random Forest RMSE:', rf_rmse)
69 print('Random Forest MAPE:', rf_mape)
70
```

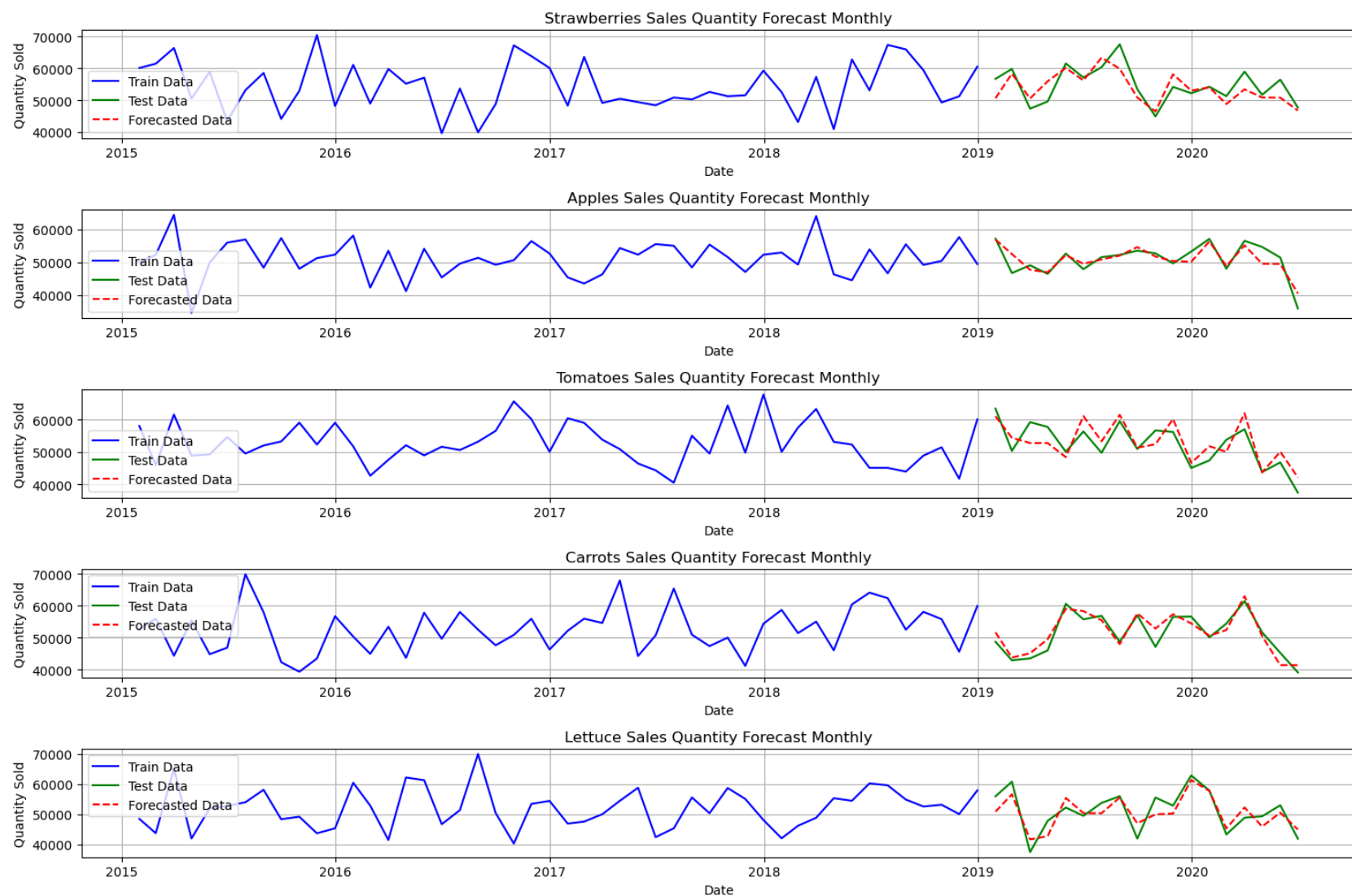
Random Forest MAE: 2661.934684588824

Random Forest RMSE: 3177.6926944638376

Random Forest MAPE: 5.203538526320775

In [85]:

```
1 plt.figure(figsize=(15, 10))
2
3
4 for i, product in enumerate(products):
5     plt.subplot(n_products, 1, i + 1)
6
7     if product in rf_X_train_dates and product in rf_y_train:
8         # Plot actual vs forecast
9         plt.plot(rf_X_train_dates[product], rf_y_train[product], color='blue', label='Actual Quantity Sold')
10        plt.plot(rf_X_test_dates[product], rf_y_test[product], color='green', label='Actual Quantity Sold')
11        plt.plot(rf_X_test_dates[product], rf_forecasts[product], color='red', linestyle='dashed', label='Forecasted Quantity')
12
13        plt.title(f'{product} Sales Quantity Forecast Monthly')
14        plt.xlabel('Date')
15        plt.ylabel('Quantity Sold')
16        plt.legend(['Train Data', 'Test Data', 'Forecasted Data'])
17        plt.grid(True)
18    else:
19        print(f"Data for '{product}' not found in rf_X_train_dates or rf_y_train.")
20
21 plt.tight_layout()
22 plt.show()
```



Conclusion

```
In [86]: 1 arima_mape_scores['--Average--'] = np.mean([v for k, v in arima_mape_scores.items()])
2 prophet_mape_scores['--Average--'] = np.mean([v for k, v in prophet_mape_scores.items()])
3 rf_mape_scores['--Average--'] = np.mean([v for k, v in rf_mape_scores.items()])
4
5
6 pd.DataFrame({'ARIMA': arima_mape_scores,
7               'Prophet': prophet_mape_scores,
8               'RandomForestRegressor': rf_mape_scores})
```

Out[86]:

	ARIMA	Prophet	RandomForestRegressor
Strawberries	8.252509	8.665478	5.482484
Apples	7.672138	10.521305	3.627103
Tomatoes	11.009154	12.012712	6.526076
Carrots	11.643805	18.719716	4.069554
Lettuce	11.280264	16.485064	6.312475
--Average--	9.971574	13.280855	5.203539

Findings

1. Demand Patterns:

- Seasonality significantly influences demand, with certain months showing peaks.

2. Key Influencers:

- Weather conditions, particularly temperature and rainfall, have a strong correlation with demand.

3. Product Variability:

- Different products exhibit varied demand patterns and seasonal peaks.

Recommendations

1. Adjust Harvest Schedules:

- Use demand forecasts to adjust harvest schedules, ensuring alignment with market demand.

2. Diversify Product Offering:

- Consider diversifying the product range to stabilize revenue throughout the year.

3. Enhance Data Collection:

- Improve data collection methods for more accurate and comprehensive forecasting.

4. Leverage the Best Performing Model:

- The RandomForestRegressor consistently outperformed ARIMA and Prophet across all products. Thus, it is recommended to primarily use the RandomForestRegressor for demand forecasting.

5. Tailor Strategies to Product-Specific Insights:

- Strawberries:
 - With the lowest error rate, focus on refining the model to maintain accuracy.
 - Implement strategies to optimize strawberry production based on accurate forecasts.
- Apples:
 - Continue using RandomForestRegressor due to its high accuracy.
 - Consider exploring additional features to improve the model further.
- Tomatoes, Carrots, and Lettuce:
 - Despite higher error rates compared to strawberries and apples, the RandomForestRegressor still outperforms other models.
 - Focus on data quality and potentially additional features specific to these crops to reduce errors.

6. Model Improvement and Maintenance:

- Regularly monitor the performance of the RandomForestRegressor to ensure its predictive accuracy remains high.
- Periodically retrain the model with new data to adapt to changing market conditions and trends.

7. Data Collection Enhancements:

- Improve the quality and granularity of data collected, particularly for products with higher prediction errors, to enhance model accuracy.
- Consider additional data points such as soil conditions, pest infestations, and more detailed weather patterns.

Future Work

1. Incorporate More Variables:

- Include additional factors like market trends and competitor data for more robust forecasting.

2. Continuous Model Improvement:

- Regularly update and refine models to adapt to changing market conditions.

3. Real-time Forecasting:

- Implement real-time data processing and forecasting to enable more responsive decision-making.

In []:

1