```
1   import pandas as pd
2   from faker import Faker
3   import random
4
5   # Initialize Faker
6   fake = Faker()
7   # Load the data
8   df = pd.read_csv('fake_data.csv')
9
10  def generate_random_data(num_records):
11      random_data = []
12      for _ in range(num_records):
13          order_id = "LK" + ''.join(random.choices('0123456789', k=7))
14          email = fake.email()
15          sales = random.choice([0, random.randint(100, 1000)])
16          date = fake.date_time_this_year().strftime("%Y-%m-%d %H:%M:%S %z")
17          product_quantity = 1
18          product_name = fake.sentence(nb_words=6)  # Random product name
19          product_sku = ''.join(random.choices('0123456789ABCDEFGHIJKLMNOPQRSTUVWXYZ', k=10))
20          customer_name = fake.name()
21          customer_city = fake.city()
22          customer_zip = fake.zipcode()
23          customer_phone = ''.join(random.choices('0123456789', k=10))
24
25          random_data.append({
26              "Order id": order_id,
27              "Email": email,
28              "Sales": sales,
29              "Date": date,
30              "Product quantity": product_quantity,
31              "Product name": product_name,
32              "Product sku": product_sku,
33              "Customer Name": customer_name,
34              "Customer City": customer_city,
35              "Customer Zip": customer_zip,
36              "Customer Phone": customer_phone,
37          })
38
39      return pd.DataFrame(random_data)
40
41  # Generate 10 random records
42  random_df = generate_random_data(100000)
43
44  # Combine the sample data with the random data
45  combined_df = pd.concat([df, random_df], ignore_index=True)
46
47  # Save to CSV
48  combined_df.to_csv("combined_data.csv", index=False)
49
50  print(combined_df)
51
```

```
          Order id                      Email  Sales                 Date  \
0        LK5331309        ashley78@example.com    587  2024-07-04 12:38:16
1        LK5742967     andrewgraham@example.org      0  2024-04-14 13:57:15
2        LK1892755      ruizrobert@example.com      0  2024-05-01 15:27:08
3        LK8705549           hfrey@example.org      0  2024-01-30 17:07:39
4        LK0157039          apena@example.net    150  2024-05-13 13:35:12
...          ...                        ...    ...                  ...
199995   LK4325732       xjimenez@example.org    760  2024-01-07 07:23:24
199996   LK4923586        zgaines@example.com      0  2024-05-16 20:33:34
199997   LK9781515  kennethhernandez@example.net      0  2024-07-02 21:28:06
199998   LK1225201    anthonyperkins@example.org    867  2024-05-29 06:35:33
199999   LK0268894        steven13@example.org      0  2024-04-04 10:00:43

        Product quantity                                  Product name  \
0                      1                             Fall explain cover.
1                      1       Article education factor yeah project poor.
2                      1                         Policy identify sit age.
3                      1                         Drive recent system seat.
4                      1                            Test herself five job.
...                  ...                                           ...
199995                 1   Represent chair staff campaign something anima...
199996                 1                      Beat up door office state Mr.
199997                 1                     Guess good sort unit marriage.
199998                 1         Sea two finally wonder mention music about.
199999                 1                     Space pressure race hand perhaps.

        Product sku     Customer Name      Customer City  Customer Zip  \
0        QLZA1MEZRA        Jacob King       Jonathanside         17251
1        WD5QA9S0FX  Stephanie Shaffer         New Renee         97013
2        447Z4O8GXU        Joseph Rice        Lake Krista         75772
3        ERNY4YVY0M     Matthew Wright          Shawhaven         49529
4        9UPB6QZS16        Cory Romero     Alvaradohaven         41927
...          ...               ...               ...            ...
199995   6KCHQYDY1B     Andres Stewart    West Susanfurt         40488
199996   2AZNS8FYS6    Heather Bradley          Karinaton         89531
199997   DAQYGFAVL6     Patricia Cain   New Raymondfort         91569
199998   8AO7F93LN4     Lisa Phillips         Mosleyview         08531
```

```
199999  4IUGIAOH6R      Colleen Perez    Michelleberg        05399

        Customer Phone
0          1943237940
1          9457143601
2          6520579073
3          3483957334
4          1956670057
...               ...
199995     8435625242
199996     5061574841
199997     8950008847
199998     1509151731
199999     3473473982

[200000 rows x 11 columns]
```

```python
1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  from sklearn.linear_model import LinearRegression
4  from sklearn.tree import DecisionTreeRegressor
5  from sklearn.ensemble import RandomForestRegressor
6  import xgboost as xgb
7  from sklearn.metrics import mean_squared_error, r2_score
8  import numpy as np
9
10 # Load the data
11 df = pd.read_csv('fake_data.csv')
12
13 # Prepare the data
14 df['Date'] = pd.to_datetime(df['Date'])
15 df.set_index('Date', inplace=True)
16 df['Month'] = df.index.month
17 df['Year'] = df.index.year
18
19 # Feature and target variable
20 X = df[['Month', 'Year']]
21 y = df['Sales']
22
23 # Split the data into training and test sets
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
25
26 # Function to evaluate models
27 def evaluate_model(model, X_train, X_test, y_train, y_test):
28     model.fit(X_train, y_train)
29     y_pred = model.predict(X_test)
30     mse = mean_squared_error(y_test, y_pred)
31     r2 = r2_score(y_test, y_pred)
32     return mse, r2, y_pred
33
34 # Initialize models
35 models = {
36     'Linear Regression': LinearRegression(),
37     'Decision Tree': DecisionTreeRegressor(random_state=0),
38     'Random Forest': RandomForestRegressor(random_state=0),
39     'XGBoost': xgb.XGBRegressor(objective='reg:squarederror', random_state=0)
40 }
41
42 # Evaluate each model
43 results = {}
44 for name, model in models.items():
45     mse, r2, y_pred = evaluate_model(model, X_train, X_test, y_train, y_test)
46     results[name] = {'MSE': mse, 'R-squared': r2, 'Predictions': y_pred}
47
48 # Print results
49 for name, result in results.items():
50     print(f"{name} - MSE: {result['MSE']:.4f}, R-squared: {result['R-squared']:.4f}")
51
52 # Choose the best model based on R-squared
53 best_model_name = max(results, key=lambda k: results[k]['R-squared'])
54 best_model = models[best_model_name]
55
56 # Compare actual and predicted sales
57 y_pred_best = results[best_model_name]['Predictions']
58 comparison_df = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred_best})
59 print(comparison_df)
60
61 # Ensure indices match
62 predicted_df = pd.DataFrame(y_pred_best, index=X_test.index, columns=['Predicted_Sales'])
63 df = df.join(predicted_df)
64
65 # Predict next month's sales based on the last billing date for each customer
66 last_billing_dates = df.reset_index().groupby(['Email', 'Customer Name'])['Date'].max()
67 predictions = []
68
69 for (email, customer_name), last_date in last_billing_dates.items():
70     # Get the sales for the last billing date
71     last_sales = df.loc[df.index == last_date, 'Sales'].values[0]
72
73     next_month = (last_date.month % 12) + 1
74     next_year = last_date.year if next_month > last_date.month else last_date.year + 1
75     next_month_features = np.array([[next_month, next_year]])
76     predicted_sales = best_model.predict(next_month_features)
77
78     predictions.append((email, customer_name, last_date, last_sales, next_month, next_year, predicted_sales[0]))
79
80 predictions_df = pd.DataFrame(predictions, columns=['Email', 'Customer Name', 'Last_Billing_Date', 'Last_Sales', 'Next_Month', 'Next_Year', 'Predicted_Sale
81
82 # Save the predictions to a CSV file
83 predictions_df.to_csv('customer_next_month_predictions.csv', index=False)
84
85 print(f"Best model: {best_model_name}")
86 print(predictions_df.head())
87
```
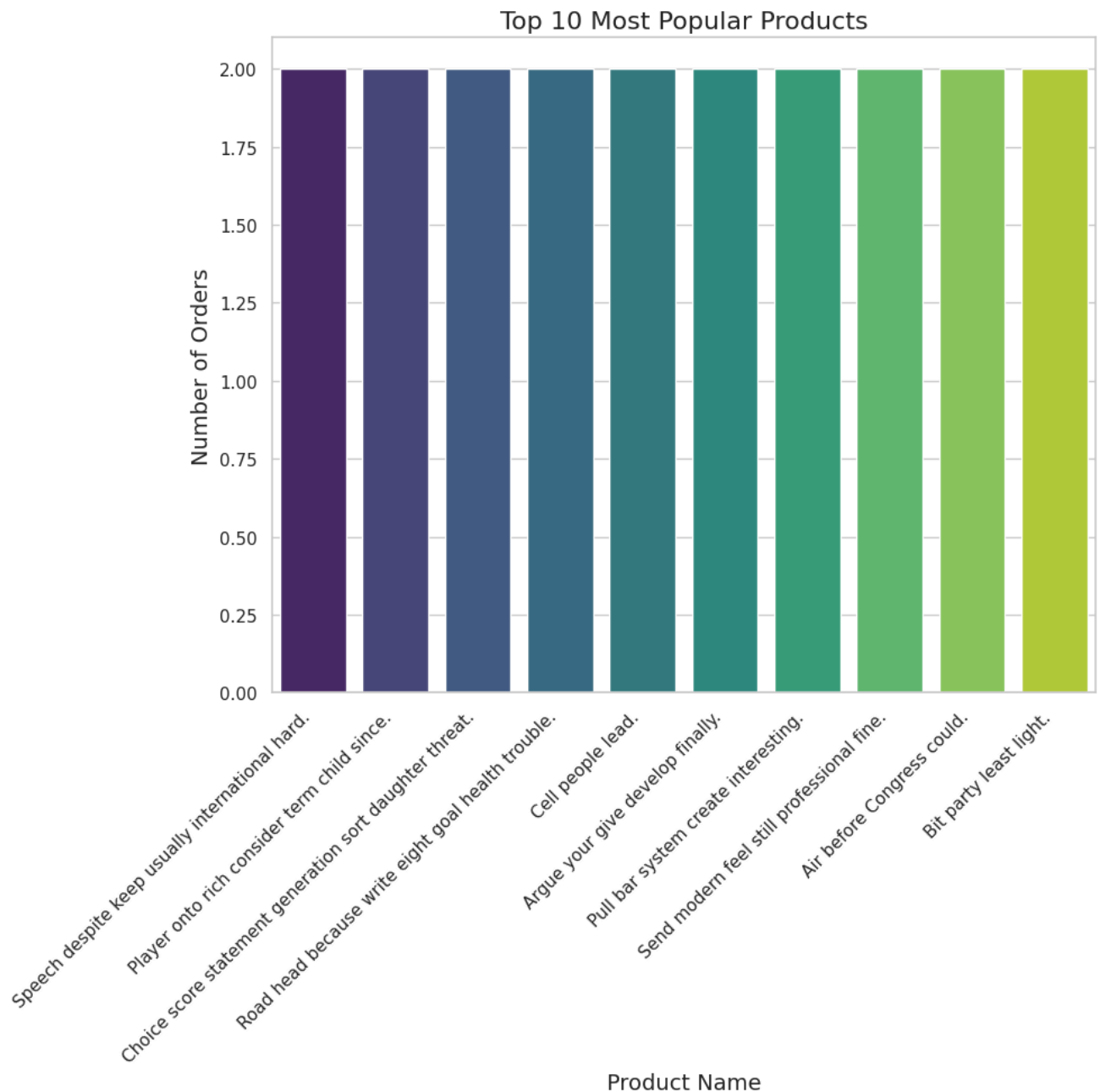
```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegress
  warnings.warn(
Best model: Linear Regression
                Email    Customer Name   Last_Billing_Date  Last_Sales  \
0    aabbott@example.net     Dawn Arroyo 2024-01-12 14:41:35         310
1    aadams@example.net  Brandon Schmidt 2024-06-10 11:37:44         618
2    aadams@example.org      Anita Soto 2024-06-24 14:55:02         608
3    aadams@example.org   Cindy Michael 2024-03-01 22:22:37         839
4    aadkins@example.com     Mark Barnes 2024-04-01 17:59:30           0

   Next_Month  Next_Year  Predicted_Sales
0           2       2024       276.441892
1           7       2024       277.664137
2           7       2024       277.664137
3           4       2024       276.930790
4           5       2024       277.175239
```

```
 1  import pandas as pd
 2  import matplotlib.pyplot as plt
 3  import seaborn as sns
 4
 5
 6  # Set the style for seaborn plots
 7  sns.set(style='whitegrid')
 8
 9  # 1. Product Popularity Analysis
10  plt.figure(figsize=(10, 10))
11  top_products = df['Product name'].value_counts().head(10)
12  sns.barplot(x=top_products.index, y=top_products.values, palette='viridis')
13  plt.title('Top 10 Most Popular Products', fontsize=16)
14  plt.xlabel('Product Name', fontsize=14)
15  plt.ylabel('Number of Orders', fontsize=14)
16  plt.xticks(rotation=45, ha='right')
17  plt.tight_layout()
18  plt.show()
19
20
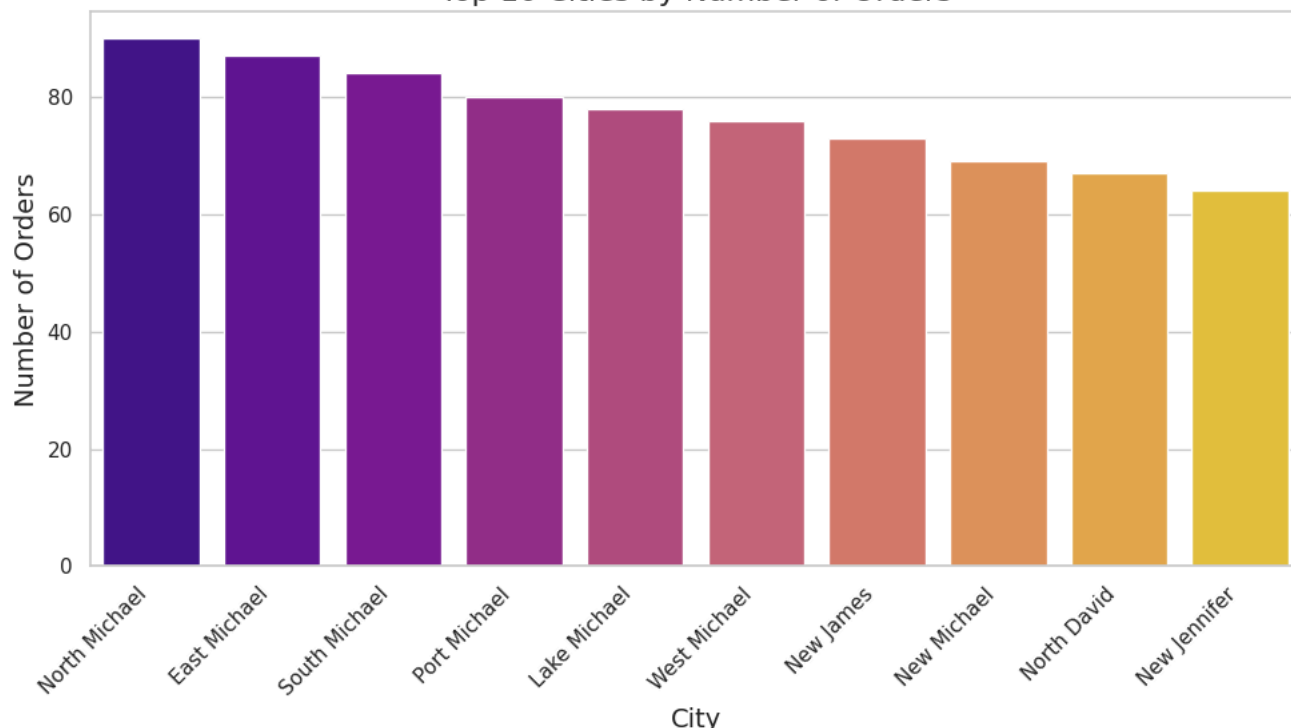```

```
<ipython-input-22-14fc856ef185>:12: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set

  sns.barplot(x=top_products.index, y=top_products.values, palette='viridis')
```



Top 10 Most Popular Products

```
 1 # 2. Customer Segmentation by Location
 2 plt.figure(figsize=(10, 6))
 3 top_cities = df['Customer City'].value_counts().head(10)
 4 sns.barplot(x=top_cities.index, y=top_cities.values, palette='plasma')
 5 plt.title('Top 10 Cities by Number of Orders', fontsize=16)
 6 plt.xlabel('City', fontsize=14)
 7 plt.ylabel('Number of Orders', fontsize=14)
 8 plt.xticks(rotation=45, ha='right')
 9 plt.tight_layout()
10 plt.show()
11
12
13
14
```

```
<ipython-input-23-2006098d045a>:4: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `l

  sns.barplot(x=top_cities.index, y=top_cities.values, palette='plasma')
```
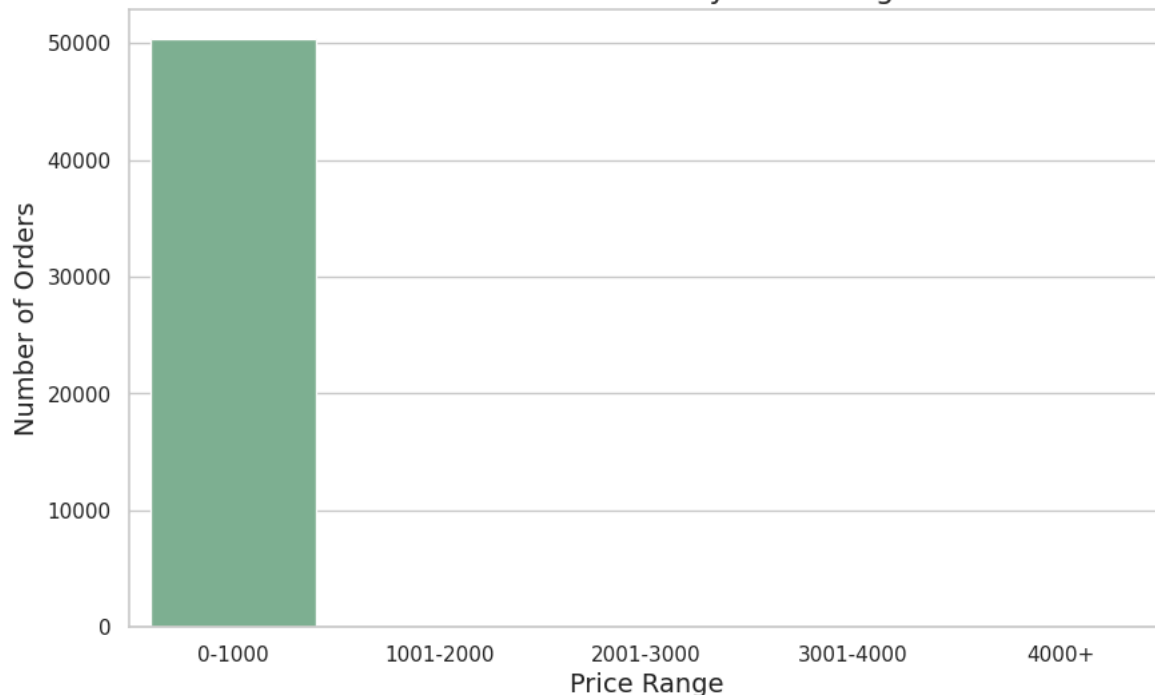
## Top 10 Cities by Number of Orders



```
 1  # 4. Price Point Analysis
 2  df['Sales'] = pd.to_numeric(df['Sales'], errors='coerce')
 3  price_ranges = pd.cut(df['Sales'], bins=[0, 1000, 2000, 3000, 4000, float('inf')],
 4                        labels=['0-1000', '1001-2000', '2001-3000', '3001-4000', '4000+'])
 5  price_range_counts = price_ranges.value_counts()
 6  plt.figure(figsize=(10, 6))
 7  sns.barplot(x=price_range_counts.index, y=price_range_counts.values, palette='crest')
 8  plt.title('Distribution of Sales by Price Range', fontsize=16)
 9  plt.xlabel('Price Range', fontsize=14)
10  plt.ylabel('Number of Orders', fontsize=14)
11  plt.show()
12
```

```
<ipython-input-25-f7feba57bb0b>:7: FutureWarning:

    Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `l

        sns.barplot(x=price_range_counts.index, y=price_range_counts.values, palette='crest')
```
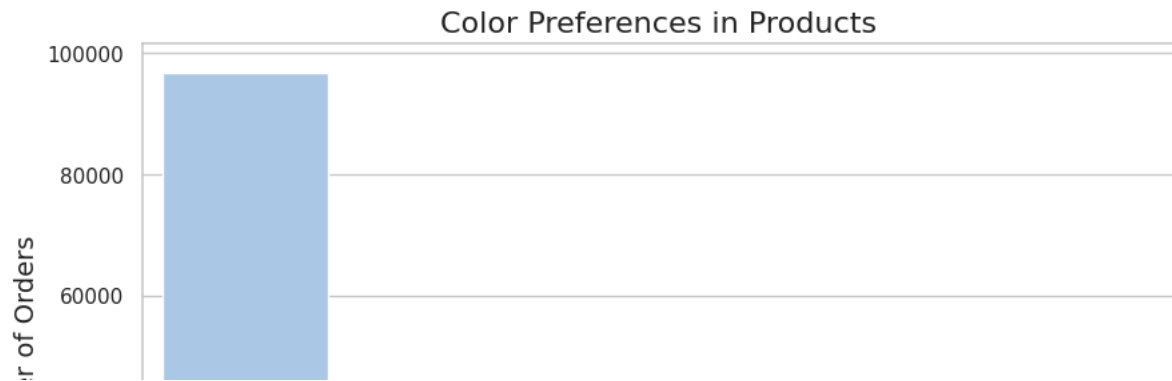


Distribution of Sales by Price Range

```
 1
 2   # 5. Time-based Analysis
 3   # df['Date'] = pd.to_datetime(df['Date'])
 4   # orders_by_hour = df['Date'].dt.hour.value_counts().sort_index()
 5   # plt.figure(figsize=(14, 7))
 6   # sns.lineplot(x=orders_by_hour.index, y=orders_by_hour.values, marker='o', color='royalblue')
 7   # plt.title('Number of Orders by Hour of Day', fontsize=16)
 8   # plt.xlabel('Hour of Day', fontsize=14)
 9   # plt.ylabel('Number of Orders', fontsize=14)
10   # plt.grid(True)
11   # plt.show()
12
13   # 6. Color Preference Analysis
14   def extract_color(product_name):
15       colors = ['Blue', 'Red', 'Pink', 'Green', 'Grey', 'White', 'Ivory', 'Navy', 'Rose', 'Lavender', 'Mustard']
16       for color in colors:
17           if color.lower() in product_name.lower():
18               return color
19       return 'Other'
20
21   df['Color'] = df['Product name'].apply(extract_color)
22   color_counts = df['Color'].value_counts()
23   plt.figure(figsize=(10, 6))
24   sns.barplot(x=color_counts.index, y=color_counts.values, palette='pastel')
25   plt.title('Color Preferences in Products', fontsize=16)
26   plt.xlabel('Color', fontsize=14)
27   plt.ylabel('Number of Orders', fontsize=14)
28   plt.xticks(rotation=45, ha='right')
29   plt.show()
30
31
32
```

```
<ipython-input-26-eb013bfabf40>:23: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `l

  sns.barplot(x=color_counts.index, y=color_counts.values, palette='pastel')
```

## Color Preferences in Products



```
1 # 3. Size Analysis
2 plt.figure(figsize=(8, 6))
3 size_counts = df['Product name'].str.extract(r'(\d+X|[XLS])').value_counts()
4 size_counts.plot(kind='pie', autopct='%1.1f%%', colors=sns.color_palette('Set2'))
5 plt.title('Distribution of Product Sizes', fontsize=12)
6 plt.ylabel('')
7 plt.show()
```

### Distribution of Product Sizes