

```
In [3]: import pandas as pd

In [5]: Amazon_sales_data = pd.read_csv(r"C:\Users\steph\OneDrive\Documents\Desktop\Amazon sales data\amazon_sales_data 2025.csv")

In [7]: Amazon_sales_data.head()

Out[7]:
```

	Order ID	Date	Product	Category	Price	Quantity	Total Sales	Customer Name	Customer Location	Payment Method	Status
0	ORD0001	14-03-2025	Running Shoes	Footwear	60	3	180	Emma Clark	New York	Debit Card	Cancelled
1	ORD0002	20-03-2025	Headphones	Electronics	100	4	400	Emily Johnson	San Francisco	Debit Card	Pending
2	ORD0003	15-02-2025	Running Shoes	Footwear	60	2	120	John Doe	Denver	Amazon Pay	Cancelled
3	ORD0004	19-02-2025	Running Shoes	Footwear	60	3	180	Olivia Wilson	Dallas	Credit Card	Pending
4	ORD0005	10-03-2025	Smartwatch	Electronics	150	3	450	Emma Clark	New York	Debit Card	Pending

```
In [9]: Amazon_sales_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 258 entries, 0 to 249
Data columns (total 11 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Order ID              258 non-null    object
 1   Date                  258 non-null    object
 2   Product               258 non-null    object
 3   Category              258 non-null    object
 4   Price                 258 non-null    int64
 5   Quantity              258 non-null    int64
 6   Total Sales           258 non-null    int64
 7   Customer Name         258 non-null    object
 8   Customer Location     258 non-null    object
 9   Payment Method        258 non-null    object
10   Status                258 non-null    object
dtypes: int64(3), object(8)
memory usage: 21.6+ KB

In [11]: # Checking the correct Date format

In [13]: Amazon_sales_data["Date"] = pd.to_datetime(Amazon_sales_data["Date"],format="%d-%b-%Y")

In [15]: # Changing the column name

In [17]: Amazon_sales_data.rename(columns={Amazon_sales_data.columns[1]:"Order_Date"},inplace=True)

In [19]: Amazon_sales_data["Order_Date"] = pd.to_datetime(Amazon_sales_data["Order_Date"],errors = "coerce")

In [21]: Amazon_sales_data["Order_Date"].head()

Out[21]:
```

	Order_Date
0	2025-03-14
1	2025-03-20
2	2025-02-15
3	2025-02-19
4	2025-03-10

```
Name: Order_Date, dtype: datetime64[ns]

for checking that product name and all in correct space or extra space

if its had extra space we need to remove them

In [27]: Amazon_sales_data.apply(lambda x:x.str.strip() if x.dtype=="object" else x)

Out[27]:
```

	Order ID	Order_Date	Product	Category	Price	Quantity	Total Sales	Customer Name	Customer Location	Payment Method	Status
0	ORD0001	2025-03-14	Running Shoes	Footwear	60	3	180	Emma Clark	New York	Debit Card	Cancelled
1	ORD0002	2025-03-20	Headphones	Electronics	100	4	400	Emily Johnson	San Francisco	Debit Card	Pending
2	ORD0003	2025-02-15	Running Shoes	Footwear	60	2	120	John Doe	Denver	Amazon Pay	Cancelled
3	ORD0004	2025-02-19	Running Shoes	Footwear	60	3	180	Olivia Wilson	Dallas	Credit Card	Pending
4	ORD0005	2025-03-10	Smartwatch	Electronics	150	3	450	Emma Clark	New York	Debit Card	Pending
...	...	...	...	...	...	...	...	...	...	...	...
245	ORD0246	2025-03-17	T-Shirt	Clothing	20	2	40	Daniel Harris	Miami	Debit Card	Cancelled
246	ORD0247	2025-03-30	Jeans	Clothing	40	1	40	Sophia Miller	Dallas	Debit Card	Cancelled
247	ORD0248	2025-03-05	T-Shirt	Clothing	20	2	40	Chris White	Denver	Debit Card	Cancelled
248	ORD0249	2025-03-08	Smartwatch	Electronics	150	3	450	Emily Johnson	New York	Debit Card	Cancelled
249	ORD0250	2025-02-19	Smartphone	Electronics	500	4	2000	Emily Johnson	Seattle	Amazon Pay	Completed

250 rows x 11 columns

```
In [29]: Amazon_sales_data.columns = Amazon_sales_data.columns.str.strip().str.replace(" ", "_") # for column name

In [31]: Amazon_sales_data.columns

Out[31]: Index(['Order_ID', 'Order_Date', 'Product', 'Category', 'Price', 'Quantity', 'Total_Sales', 'Customer_Name', 'Customer_Location', 'Payment_Method', 'Status'],
        dtype='object')

In [33]: Amazon_sales_data.dropna(subset=["Order_Date", "Price", "Quantity", "Total_Sales"], inplace=True)

In [35]: Amazon_sales_data.to_csv("Clean_data_csv",index=False) # this data for PowerBI Analysis

In [37]: # Exploratory Data Analysis (EDA)

In [39]: import pandas as pd

In [41]: import matplotlib.pyplot as plt

In [42]: import seaborn as sns

In [43]: # best performing product categories

In [44]: best_product_categories = Amazon_sales_data.groupby("Category")[["Total_Sales"]].sum().sort_values(by = "Total_Sales",ascending=False)

In [45]: best_product_categories

Out[45]:
```

	Total_Sales
Category	
Electronics	129950
Home Appliances	105000
Footwear	4320
Clothing	3540
Books	1035

```
In [47]: best_product_categories.plot(kind = "bar",color = "Red",figsize=(8,5))

plt.xlabel("Category")
plt.ylabel("Total_Sales")
plt.title("Best Product Categories by Total Sales")
plt.tight_layout()
plt.show()
```

```
In [48]: # Sales Trend Over Time

In [49]: daily_sales = Amazon_sales_data.groupby("Order_Date")["Total_Sales"].sum().sort_index()

In [50]: daily_sales

Out[50]:
```

Order_Date	Total_Sales
2025-02-02	3680
2025-02-03	3360
2025-02-04	6815
2025-02-05	5480
2025-02-06	11480
2025-02-07	2520
2025-02-08	1640
2025-02-09	3550
2025-02-10	8865
2025-02-11	3550
2025-02-12	3560
2025-02-13	4860
2025-02-14	1815
2025-02-15	510
2025-02-16	9540
2025-02-17	1885
2025-02-18	8810
2025-02-19	2195
2025-02-20	6730
2025-02-21	4680
2025-02-22	1880
2025-02-23	2570
2025-02-24	6980
2025-02-25	6380
2025-02-26	1580
2025-02-27	380
2025-02-28	7210
2025-03-01	1880
2025-03-02	4575
2025-03-03	130
2025-03-04	5980
2025-03-05	2980
2025-03-06	9720
2025-03-07	8660
2025-03-08	3125
2025-03-09	280
2025-03-10	3965
2025-03-11	1240
2025-03-12	710
2025-03-13	4165
2025-03-14	4350
2025-03-15	7190
2025-03-16	2735
2025-03-17	540
2025-03-18	3780
2025-03-19	4115
2025-03-20	3630
2025-03-21	1560
2025-03-22	950
2025-03-23	3980
2025-03-24	9520
2025-03-25	6815
2025-03-26	2970
2025-03-27	15
2025-03-28	2480
2025-03-29	2480
2025-03-30	300
2025-03-31	7465
2025-04-01	6680
2025-04-02	320
2025-04-03	3180

Name: Total\_Sales, dtype: int64

```
In [51]: daily_sales.plot(kind = "line",color = "Green",figsize = (8,5))

plt.xlabel("Order_Date")
plt.ylabel("Total_Sales")
plt.title("Sales Trend Over Time")
plt.tight_layout()
plt.show()
```