

In [5]: `import pandas as pd`

```
df = pd.read_csv("C:\\Users\\umaka\\Downloads\\sales_data_sample.csv", encoding=
print(df.head())
```

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	\
0	10107	30	95.70	2	2871.00	
1	10121	34	81.35	5	2765.90	
2	10134	41	94.74	2	3884.34	
3	10145	45	83.26	6	3746.70	
4	10159	49	100.00	14	5205.27	

	ORDERDATE	STATUS	QTR_ID	MONTH_ID	YEAR_ID	...	\
0	2/24/2003 0:00	Shipped	1	2	2003	...	
1	5/7/2003 0:00	Shipped	2	5	2003	...	
2	7/1/2003 0:00	Shipped	3	7	2003	...	
3	8/25/2003 0:00	Shipped	3	8	2003	...	
4	10/10/2003 0:00	Shipped	4	10	2003	...	

	ADDRESSLINE1	ADDRESSLINE2	CITY	STATE	\
0	897 Long Airport Avenue	NaN	NYC	NY	
1	59 rue de l'Abbaye	NaN	Reims	NaN	
2	27 rue du Colonel Pierre Avia	NaN	Paris	NaN	
3	78934 Hillside Dr.	NaN	Pasadena	CA	
4	7734 Strong St.	NaN	San Francisco	CA	

	POSTALCODE	COUNTRY	TERRITORY	CONTACTLASTNAME	CONTACTFIRSTNAME	DEALSIZE
0	10022	USA	NaN	Yu	Kwai	Small
1	51100	France	EMEA	Henriot	Paul	Small
2	75508	France	EMEA	Da Cunha	Daniel	Medium
3	90003	USA	NaN	Young	Julie	Medium
4	NaN	USA	NaN	Brown	Julie	Medium

[5 rows x 25 columns]

In [14]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2823 entries, 0 to 2822
Data columns (total 25 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ORDERNUMBER           2823 non-null   int64
1   QUANTITYORDERED       2823 non-null   int64
2   PRICEEACH             2823 non-null   float64
3   ORDERLINENUMBER       2823 non-null   int64
4   SALES                 2823 non-null   float64
5   ORDERDATE             2823 non-null   object
6   STATUS                2823 non-null   object
7   QTR_ID               2823 non-null   int64
8   MONTH_ID              2823 non-null   int64
9   YEAR_ID               2823 non-null   int64
10  PRODUCTLINE           2823 non-null   object
11  MSRP                  2823 non-null   int64
12  PRODUCTCODE           2823 non-null   object
13  CUSTOMERNAME          2823 non-null   object
14  PHONE                 2823 non-null   object
15  ADDRESSLINE1          2823 non-null   object
16  ADDRESSLINE2          302 non-null    object
17  CITY                  2823 non-null   object
18  STATE                 1337 non-null   object
19  POSTALCODE            2747 non-null   object
20  COUNTRY               2823 non-null   object
21  TERRITORY             1749 non-null   object
22  CONTACTLASTNAME       2823 non-null   object
23  CONTACTFIRSTNAME      2823 non-null   object
24  DEALSIZE              2823 non-null   object
dtypes: float64(2), int64(7), object(16)
memory usage: 551.5+ KB
```

In [15]: `df.describe()`

Out[15]:

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	
count	2823.000000	2823.000000	2823.000000	2823.000000	2823.0
mean	10258.725115	35.092809	83.658544	6.466171	3553.8
std	92.085478	9.741443	20.174277	4.225841	1841.8
min	10100.000000	6.000000	26.880000	1.000000	482.1
25%	10180.000000	27.000000	68.860000	3.000000	2203.4
50%	10262.000000	35.000000	95.700000	6.000000	3184.8
75%	10333.500000	43.000000	100.000000	9.000000	4508.0
max	10425.000000	97.000000	100.000000	18.000000	14082.8

In [6]: `df.shape`

Out[6]: (2823, 25)

In [13]: `#THIS WILL CHECK FOR NULL VALUES.. GIVES NUMBER OF NULL VALUES IN EACH COLUMN`
`df.isnull().sum()`

```
Out[13]: ORDERNUMBER      0
          QUANTITYORDERED  0
          PRICEEACH        0
          ORDERLINENUMBER  0
          SALES             0
          ORDERDATE        0
          STATUS           0
          QTR_ID           0
          MONTH_ID        0
          YEAR_ID          0
          PRODUCTLINE      0
          MSRP             0
          PRODUCTCODE      0
          CUSTOMERNAME     0
          PHONE            0
          ADDRESSLINE1     0
          ADDRESSLINE2     2521
          CITY             0
          STATE            1486
          POSTALCODE       76
          COUNTRY          0
          TERRITORY        1074
          CONTACTLASTNAME  0
          CONTACTFIRSTNAME 0
          DEALSIZE         0
          dtype: int64
```

```
In [8]: #THIS WILL REMOVE ROWS WHICH CONTAINS TOTAL NULL VALUES
        df.dropna(how="all")
```

Out[8]:

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES
0	10107	30	95.70	2	2871.00
1	10121	34	81.35	5	2765.90
2	10134	41	94.74	2	3884.34
3	10145	45	83.26	6	3746.70
4	10159	49	100.00	14	5205.27
...
2818	10350	20	100.00	15	2244.40
2819	10373	29	100.00	1	3978.51
2820	10386	43	100.00	4	5417.57
2821	10397	34	62.24	1	2116.16
2822	10414	47	65.52	9	3079.44

2823 rows × 25 columns



```
In [11]: #TO REMOVE DUPLICATES
df.drop_duplicates()
```

Out[11]:

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES
0	10107	30	95.70	2	2871.00
1	10121	34	81.35	5	2765.90
2	10134	41	94.74	2	3884.34
3	10145	45	83.26	6	3746.70
4	10159	49	100.00	14	5205.27
...
2818	10350	20	100.00	15	2244.40
2819	10373	29	100.00	1	3978.51
2820	10386	43	100.00	4	5417.57
2821	10397	34	62.24	1	2116.16
2822	10414	47	65.52	9	3079.44

2823 rows × 25 columns



```
In [29]: #HANDLING MISSING VALUES
df['ADDRESSLINE2']=df['ADDRESSLINE2'].fillna('null')
df['STATE']=df['STATE'].fillna('null')
df['TERRITORY']=df['TERRITORY'].fillna('null')
df['POSTALCODE']=df['POSTALCODE'].fillna('null')
```

```
In [30]: df.isnull().sum()
```

```
Out[30]: ORDERNUMBER      0
          QUANTITYORDERED  0
          PRICEEACH        0
          ORDERLINENUMBER  0
          SALES             0
          ORDERDATE        0
          STATUS           0
          QTR_ID           0
          MONTH_ID         0
          YEAR_ID          0
          PRODUCTLINE      0
          MSRP             0
          PRODUCTCODE      0
          CUSTOMERNAME     0
          PHONE            0
          ADDRESSLINE1     0
          ADDRESSLINE2     0
          CITY             0
          STATE            0
          POSTALCODE       0
          COUNTRY          0
          TERRITORY        0
          CONTACTLASTNAME  0
          CONTACTFIRSTNAME 0
          DEALSIZE         0
          dtype: int64
```

```
In [31]: #FIXING INCONSISTENCIES
df['COUNTRY'] = df['COUNTRY'].str.upper()
```

```
In [32]: #CONVERTING DATE COLUMNS TO PROPER FORMAT
df['ORDERDATE'] = pd.to_datetime(df['ORDERDATE'])
```

```
In [33]: df
```

Out[33]:

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES
0	10107	30	95.70	2	2871.00
1	10121	34	81.35	5	2765.90
2	10134	41	94.74	2	3884.34
3	10145	45	83.26	6	3746.70
4	10159	49	100.00	14	5205.27
...
2818	10350	20	100.00	15	2244.40
2819	10373	29	100.00	1	3978.51
2820	10386	43	100.00	4	5417.57
2821	10397	34	62.24	1	2116.16
2822	10414	47	65.52	9	3079.44

2823 rows × 25 columns



```
In [34]: #CLEANED DATA SET
df.to_csv("cleaned_sales_data.csv", index=False)
```

In []: