

Hand written Notes

-On-

Exploratory Data Analysis (Python)

Tools used -

- 1. NumPy**
- 2. Pandas**
- 3. Matplotlib**
- 4. Seaborn**
- 5. Urge to learn**

Divyanshu Vyas

Exploratory Data Analysis with Python

➔ List of Tasks to be covered:

- ① Initial Data Exploration
- ② Univariate Analysis
- ③ Bivariate Analysis

- ④ Dealing with duplicate rows and missing values.
- ⑤ Correlation Analysis.

* → libraries

- ➔ import numpy as np , pandas as pd → Data Analysis
- ➔ import matplotlib.pyplot as plt ; ➔ import seaborn as sns
- ➔ import Calmap :- Calendar & heatmap library
- ➔ from pandas-profiling import ProfileReport
↳ makes EDA structured.

* Dataset used → supermarket-sales data.

Invoice ID	Branch	City	Customers Type	Gender	Product line
unit Price	Qty	Tax	...	17 columns	

* To view the datatypes → df.dtypes →

col-name	datatype
⋮	⋮

 series.

* `df['date'] = pd.to_datetime(df['date'])`

* To set date ('Date') column as the index →

→ `df.set_index('Date', inplace=True)`

Summary statistics

`df.describe()`

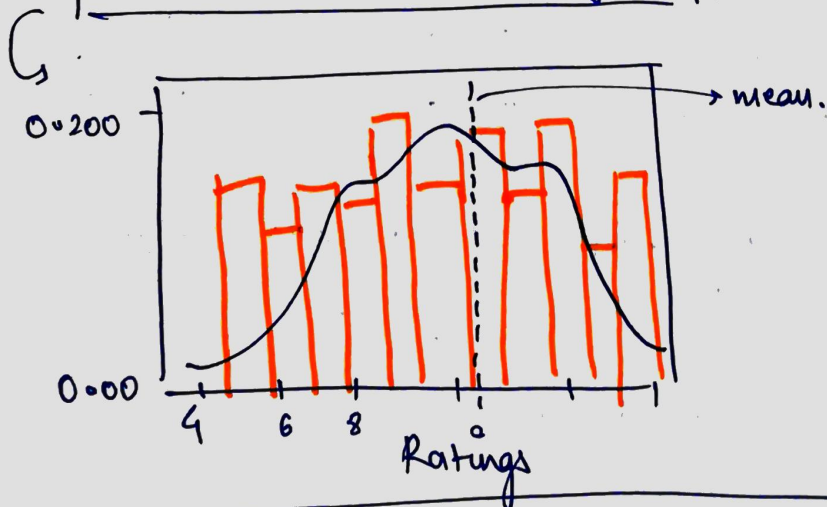
	← col-names →			
Count	✓	✓	✓	✓
mean	✓	✓	✓	✓
std	✓	✓	✓	✓
min	✓	✓	✓	✓
25%	✓	✓	✓	✓
50%	✓	✓	✓	✓
75%	✓	✓	✓	✓
max	✓	✓	✓	✓

Univariate Analysis

"Ask questions & Answer them using your numpy / pandas / matplotlib / seaborn skills."

Q1 → what does customer ratings distribⁿ look like?
Is it skewed?

Ans → `sns.distplot(df['Rating'])`

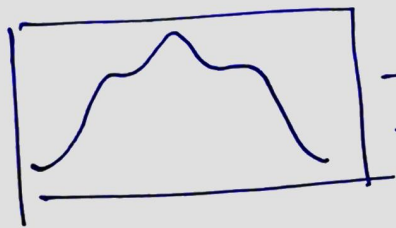


`plt.axvline(x=np.mean(df['Rating']), ls='--')`

Drawing v-line for 25th & 75th percentile.

for 25th percentile →

plt.axvline(x=np.percentile(df['Rating'], 25), ls='--')
75 for 75th %ile.



→ looks like a Uniform Distribⁿ.
→ No skew in left or Right.

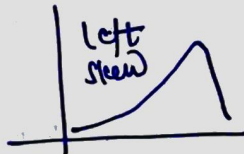
Q2:- How to see 'visually' the distribution of all the Numerical variables?

Ans → df.hist() or df.hist(figsize=(10,10))

Note →



Right skew



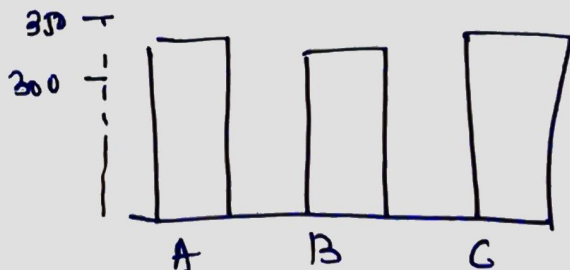
Left skew

(Trick.
whenever you can
write legend, that's
the skew)

Note → Highly correlated stuff has almost identical correlⁿ.

Q3:- The column 'Branch' has 3 categories in it → 'A', 'B', 'C'.
What's the count of each?

Ans → sns.countplot(df['Branch'])



Not much
difference.

NOTE → To get the exact count that the countplot displays

```
df['Branch'].value_counts()
```

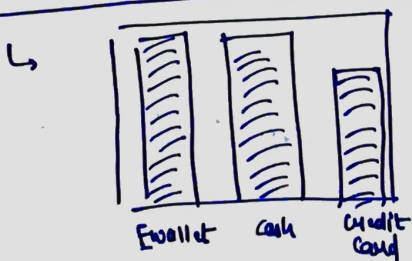
{
A 342
B 333
C 328
Name: Branch, dtype: int64}

→ We can use countplot and value_counts for any categorical variable.

Example - Perform Count Analysis for 'Payment' column
(this column has payment method categories)

∴

```
sns.countplot(df['Payments'])  
df['payments'].value_counts()
```



Hence credit card is the least popular kind of payment method.

BIVARIATE ANALYSIS

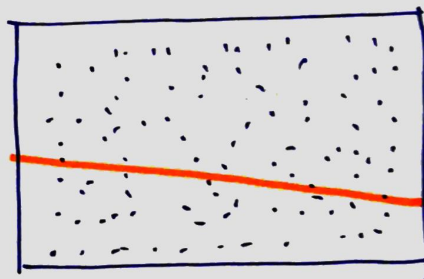
→ To check if there's a relationship b/w various columns (features)

Q:- Is there a relship b/w Gross income & customer Ratings?

`sns.scatterplot(dt['Ratings'], dt['gross income'])`

gross income

`sns.regplot(dt['Ratings'], dt['gross income'])`



Ratings →

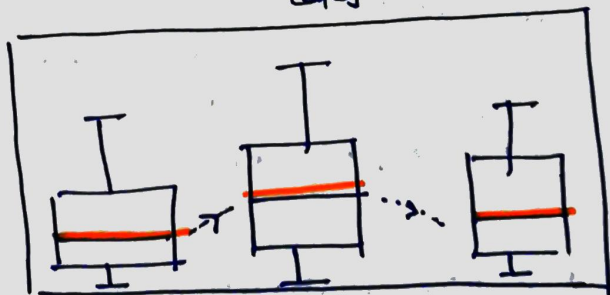
To see the baseline (orange)

Flat-Trendline → No relship.

Q:- To see the Relnshp b/w Y (Gross Income) with a categorical variable X (Branch)

⇒ Boxplot.

→ `sns.boxplot(x = dt['Branch'], y = dt['Gross Income'])`
Categorical.



A

C

B

Branch.

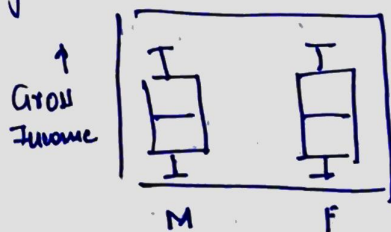
C > B > A

median Gross income

Similar Example → To see how gross income varies with gender.

sns.boxplot (x=df['Gender'], y=df['Gross Income'])

NOK → For a vertical box plot, make sure you put the reference categorical variable on the x-axis.



Q:- Time Trend Analysis → Is there any trend that the gross income (which is the earning of supermarket) (from the people) follows?

NOK → ① First convert date column into DateTime Object.

② Cannot move directly. First aggregate the income for every day since each date appears multiple times (As expected).

↓
GROUP-BY
or df['Date'].groupby('Date').sum()

→ df.groupby(df.index).mean()

↳ make a new dataframe & perform the Analysis.

DEALING with Duplicate Rows & Missing Values.

Q1:- To find number of duplicate Rows →

↳ `df.duplicated().sum()`

Returns T/F
(Boolean)

sums the No. of True's up.

Q2:- How to view the duplicated Rows?

↳ `df[df.duplicated()]`

Q3:- How to drop duplicate Rows - Permanently?

Ans - `df.drop_duplicates(inplace=True)`

↳ After this now if you run the `duplicated().sum()`
↳ Ans should be 0.

Missing Values

Q1:- How to find out the No. of missing values in each column.

Soln - `df.isna().sum()`

Series (output)

col-1	5
col-2	0
col-3	0
col-4	3
⋮	⋮

Q2:- How to visualize the No. of Missing values?

↳ `sns.heatmap(df.isnull())`

There are various ways to deal with Missing data.

2 Basic ways → ① Fill with 0
② Fill with mean(s).

Q3:- How to Fill the missing values with mean?

↳ `df.fillna(df.mean(), inplace=True)`

↳ But this only fills up the numeric columns.

↳ Categorical missing values still remain empty.

Q4:- How to Fill categorical missing values?
(missing/nans in categorical columns)

↳ Fill in with Mode.

NOTE → To get the mode for each column,
we do →

`df.mode().iloc[0]`

↳ only first row gives mode.

◦ `df.fillna(df.mode().iloc[0], inplace=True)`

↳ Finally No missing values.

Note - Whatever we've done, the EDA till now can be nicely done in one command by the Pandas Profiler. (Recommended for small datasets).

→ `data = pd.read_csv('xyz.csv')`

→ `profile = ProfileReport(data)`

→ `profile`

↳

Detailed EDA report

CORRELATION ANALYSIS.

Method 1 → `np.corrcoef(df['A'], df['B'])`

Method 2 → `df.corr()` → Correlation matrix.

↳ `np.round(df.corr(), 2)`

↳ rounded to 2 decimal places.

Method 3 → `sns.heatmap(df.corr())`

↳ `sns.heatmap(df.corr(), annot=True)`