

Data Analysis with Python

- 1) Import libraries
 - import pandas as pd
 - import numpy as np
 - import matplotlib.pyplot as plt
 - import seaborn as sns
- 2) Import Datasets
 - 1) df = pd.read_csv("pathName" / "fileName")
 - gives no. of rows & columns
 - 2) df.shape
 - top values
 - 3) df.head()
 - bottom values
 - 4) df.tail()
- 3) Data cleaning
 - Info for Data cleaning
 - 1) df.info()
 - Null Values T/F
 - 2) df.isnull() or
 - Sum of all Null Values df.isnull().sum().sum()
 - Removes NaN (missing)
 - Deletes columns with 0 values in df.info()
 - 3) df.drop(['colName'], axis=1, inplace=True)
 - removes Null values
 - 4) df.dropna()
 - Removes entire row
 - 5) df2 = df.dropna(how='any')
 - df2 = df.dropna(how='all')
 - If any Nan values
 - If all Nan values
 - Remove duplicates
 - sum-duplicate values
 - 6) df['colName'] = df['colName'].duplicated().sum()
 - removes duplicates
 - 7) df.drop_duplicates()
 - fill Nan (missing)
 - fill Nan
 - 8) df = df.fillna(value=0)
 - fill Nan of sp. col
 - 9) df['colName'].fillna('unknown', inplace=True)
 - fill Nan with Prev. val
 - 10) df2 = df.fillna(method='pad')
 - fill Nan with forw. val
 - 11) df2 = df.fillna(method='bfill')
 - fill Nan in diff. col
 - 12) df2 = df.fillna({'colName1': 'value1', 'colName2': 'value2'})
 - fill Nan with Mean
 - 13) df2 = df.fillna(value=df['colName'].mean())
 - replaces multiple values
 - 14) df.loc[rowNo, 'colName'] = 'new value'
 - removes multiple values
 - 15) for x in df.index:
 - If df.loc[x, 'ColName'] = 'old value':
df.loc[x, 'ColName'] = 'New value' or
df.drop[x, inplace=True]
 - Replace values
 - replace any/Nan value
 - 16) df2 = df.replace(to_replace=np.nan, value='0')
or any value
 - Converts value from numeric to string
 - 17) def conv(values):
 - if values == 0:
return 'no'
 - else:
return 'yes'
 - df['colName'] - df['colName'].apply(conv)
 - Renames col Name
 - 18) df.rename(columns={'oldValue': 'newValue'})
 - Descriue col
 - 19) df.describe() or
 - 20) df[['colName']].describe()

Change Data Type

- Find the datatypes 21) df.dtypes
- changes Datatype(date) 22) df['Date'] = pd.to_datetime(df['Date'])
- changes Date format to YYYY-mm-DD 23) df['Date'] = pd.to_datetime(df['Date']).strftime('%Y-%m-%d')
- change 11 (Other) 24) df['ColName'] = df['ColName'].astype('float')

Save file

- Save cleaned file 25)

```
import os
directory_path = 'C:/Users/shagu/Documents/Data Analyst/SQL'

csv_path = os.path.join(directory_path, 'practice1.csv')
df.to_csv(csv_path, index=False)
print(f"Table data saved to {csv_path}")

Table data saved to C:/Users/shagu/Documents/Data Analyst/SQL\practice1.csv
```

- To Know current CSV Path 26) os.getcwd()
- Find list of ColNames 27) df.columns.tolist()
- Merge 2 tables 28) pd.merge(df1, df2, on='ColName', how='inner')
using ID

Missing Values

- 1) find % of Nan values

df.isnull().mean() * 100

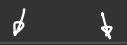
- 2) If below 5%

MCAR



i) CCA

ii) Simple Imputation
(univariate)



Numerical

Categorical

- i) Mean
- ii) Median
- iii) Arbitrary

- i) Mode
- ii) Arbitrary

If above 5%

MCAR / MAR



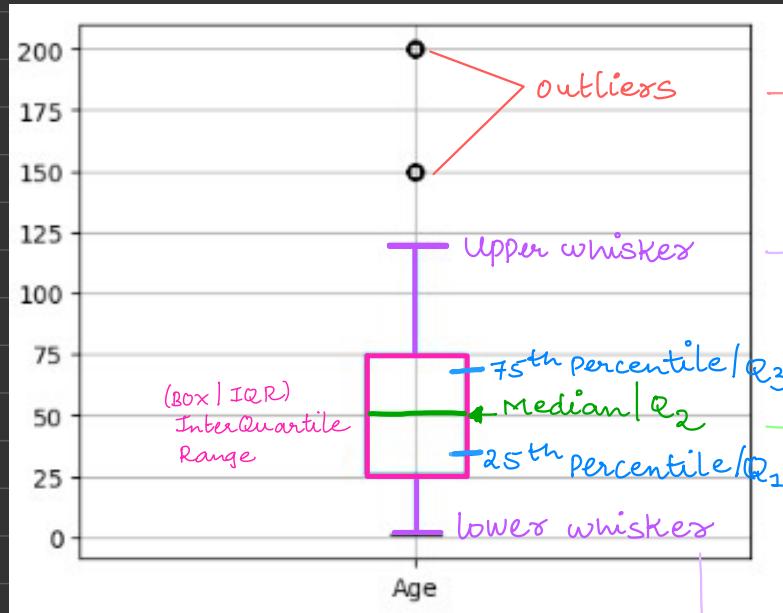
KNN

MNAR



Iterative

Box plot → detects outliers



→ Data points that fall outside 1.5 times the IQR
Age = 150 & 200

→ Starts at Q₃ & extends upto largest non-outlier

the value below which 75% data lies. Age = 75

→ Middle value of dataset. Age = 50

the value below which 25% data lies. Age = 25

starts at Q₁ & extends upto smallest non-outlier

outliers

IQR - a rectangular box

whiskers

represents middle 50% of data

Vertical lines extending from the box & Not considered outliers.

plotted individually as dots or circles

4) EDA
(Exploratory Data Analysis)

• Group By
• Data Binning

• Pearson Correlation

① Group by

gb = df.groupby(['ColName1', 'ColName2'])['ColName3'].count()
.aggregate('count').reset_index(name = 'NewColumnName') Non-Null rows only

gb = df.groupby(['ColName1', 'ColName2']).size()
.reset_index(name = 'NewColumnName') - counts all rows

Categorical Values

gb = df.groupby('ColName1', as_index = False)[['ColName2']].sum()
.sort_values(by = 'ColName2', ascending = False)

Numerical Values

gb[.sort_values(by = 'Count', ascending = False)] → For sorting in descending

② Data Binning

bins = np.linspace(min(df['ColName']), max(df['ColName']), 4) creates equal space bins
labels = ("low", "moderate", "high") → 1 no. fewer than bin → 1 more than
df['New_Col'] = pd.cut(df['ColName'], bins = bins, labels = labels, no. of labels include_lowest = True)

③ Pearson Coefficient

Converts DF to Float

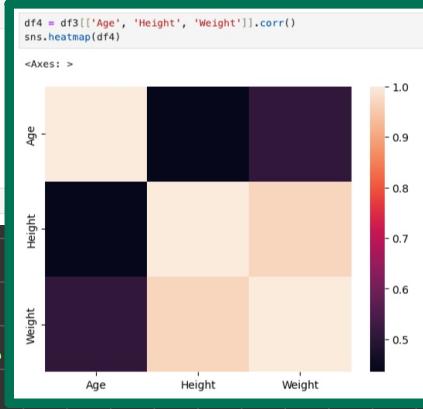
- 1) df.corr()
- 2) df[['ColName', 'ColName2']].corr()
- 3) df.apply(pd.to_numeric, errors='coerce')
- 4) sns.heatmap(df)

```
# Convert entire DataFrame to float, non-convertible values become NaN
df4 = df3.apply(pd.to_numeric, errors='coerce').corr()
df4
```

	name	Age	Height	Weight
name	NaN	NaN	NaN	NaN
Age	NaN	1.000000	0.435143	0.509119
Height	NaN	0.435143	1.000000	0.967105
Weight	NaN	0.509119	0.967105	1.000000

```
sns.heatmap(df4)
```

more +ve value & light colour → more +ve relationship



bins = np.linspace(min(state['Orders']), max(state['Orders']), 4)
 labels = ('Min', 'Moderate', 'Max') → 4-1 = 3 labels
 state['Status'] = pd.cut(state['Orders'], bins=bins, labels=labels, include_lowest=True)

	State	Orders	Status
14	Uttar Pradesh	4807	Max
10	Maharashtra	3810	Max
7	Karnataka	3240	Moderate
2	Delhi	2740	Moderate
9	Madhya Pradesh	2252	Moderate
0	Andhra Pradesh	2051	Moderate
5	Himachal Pradesh	1568	Min
8	Kerala	1137	Min
4	Haryana	1109	Min
3	Gujarat	1066	Min
1	Bihar	1062	Min
6	Jharkhand	953	Min
15	Uttarakhand	824	Min
12	Rajasthan	555	Min
11	Punjab	495	Min
13	Telangana	312	Min

If NOT: then manually enter bin edges (1 more than labels)

- close to +1 = +ve relationship x↑, y↑
- close to -1 = -ve relationship x↑, y↓
- close to 0 = No relation

Libraries

- Seaborn
- 1) sns.load_dataset(" ")
 - 2) sns.countplot(data=df, x=' ')
 - 3) sns.barplot(data=df, x=' ', y=' ')
 - 4) sns.lineplot()
 - 5) sns.displot()
 - 6) sns.Scatterplot()
 - 7) sns.histplot()

Parameters

- hue = Category
- palette = Colour
- bins = binning
- KDE = For density

Labels ax = sns.barplot()
 the For bars in ax.containers:
 bars ax.bar_label(bars)

figure 1) sns.set({'figure.figsize':(5,5)})
 size 2) plt.figure(figsize=5,5)

marker = '*'
 linestyle = '--'
 linewidth = 2
 colors = 'g'

plt.legend(title=' ',
 bbox_to_anchor=(1.05, 1),
 loc='upper left')

- Matplotlib
- 1) plt.plot(x=' ', y=' ')
 - 2) plt.title(' ')
 - 3) plt.xlabel(' ')
 - 4) plt.ylabel(' ')
 - 5) plt.grid()
 - 6) plt.grid(axis='x') · only on x-axis
 - 7) plt.show()
 - 8) plt.Legend(title=' ') or plt.legend()
 - line 9) plt.plot()
 - bar 10) plt.bar()
 - horizontal bar 11) plt.barch()
 - pie 12) plt.pie()
 - Scatter 13) plt.scatter(c=colors, cmap=' ')
 - histogram 14) plt.hist() visualis, green etc.
 - 15) plt.subplot(nrows, ncols, index)
 - Saves chart 16) plt.savefig("chart.png")
 - clears figure 17) plt.clf()
 - density plot 18) plt.plot(kind='Kde')

plt.subplot(2, 1, 1)
 ↓ ↓ first index
 2 rows 1st column

```

x = ['ColName']
mylabels = [' ', ' ', ' ', ' ']
mycolors = ['r', 'hotpink', 'b', '#4CAF50']
myexplode = [0.2, 0, 0, 0]

```

**

Code `plt.figure(figsize=(4,4))`

```

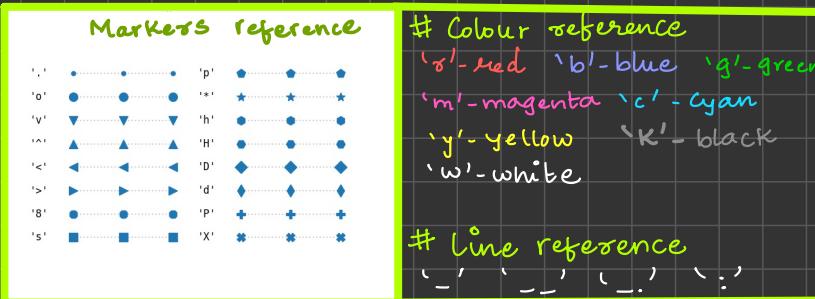
sns.pie(plt.pie(x=' ', labels=mylabels, colors=mycolors, autopct='%.1f%%',
                  startangle=90, explode=myexplode))

```

```

plt.title(' ')
plt.legend(title=' ')
plt.xlabel(' ')
plt.ylabel(' ')
plt.xticks(rotation=15)
plt.show()

```



Pandas

- directory to be same
 - 1) `df = pd.read_csv("Path name" / "filename")`
 - 2) `df.iloc[[0:4]]`
 - 3) `df.iloc[[0, 4]]`
- 1st & rows
 - 4) `df.loc[RowNo, 'ColName']`
 - 5) `df[['ColName']]`
- 1st and 5th row
 - 6) `pd.options.display.max_rows = 10,000`
 - 7) `df.head()`
 - 8) `df.info()`
 - 9) `df.describe()`
 - 10) `df.rename(columns = {'oldValue': 'newValue'})`
 - 11) `df.corr()`
- Column of 2nd row
 - 12) `df = pd.read_json("file path" / "file name")`
 - 13) `df.to_string()`
- Column as DataFrame
 - 14) `pd.DataFrame()`
- set no. of rows

- converts JSON file to DataFrame
- saves as CSV file
- converts to DataFrame
- multiplication of matrix

- 1) `arr.copy()`] Owns the array
- 2) `arr.view()`] data
- 3) `arr.reshape()`] mirror
- 4) `arr.flatten()`] image
- 5) `arr = np.array()`
- 6) `np.arange()`
- 7) `np.linspace()`
- 8) `np.dot()`
- 9) `np.transpose()`
- 10) `np.char.add()`
- 11) `np.char.replace()`
- 12) `np.char.upper()`
- 13) `np.char.lower()`

- concat the string
- replace the value
- splits where space
- splits where newline

- 16) `np.mean()`
- 17) `np.median()`
- 18) `np.std()`
- 19) `np.sum()`
- 20) `np.min() / max()`

$$\text{Speed} = [98, 97, 95, 96, 103]$$

$$\text{mean} = \frac{98 + 97 + 95 + 96 + 103}{5}$$

median = 95, 96, 97, 98, 103
(soot the ascending)
odd \rightarrow central value
even \rightarrow mean of 2 central values

different than this
 subplot(1, 2, 1)
 $\downarrow \quad \downarrow$ Index
 rows cols
 $\equiv 11$

subplot (overlapped)

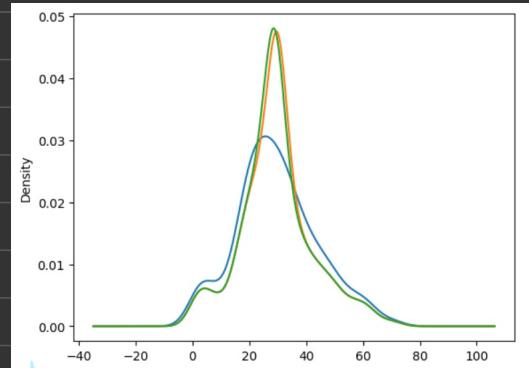
① `fig = plt.figure()`
`ax = fig.add_subplot(111)`

`df['Age'].plot(kind='Kde', ax=ax)`
`df['Age_mean'].plot(kind='Kde', ax=ax)`
`df['Age_median'].plot(kind='Kde', ax=ax)`

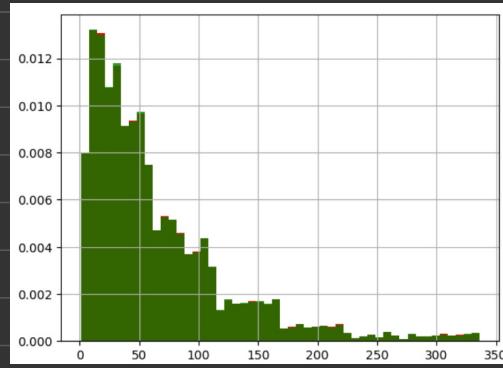
} can select
 colour as
 well

② `fig = plt.figure()`
`ax = fig.add_subplot(111)`

`df['colName'].hist(bins=50, ax=ax, density=True, color='red')`
`new_df['colName'].hist(bins=50, ax=ax, density=True, color='red')`



variation



— overlaps
 (No diff.)

Subplots (grid of plots)

`cols = ['colName1', 'colName2']`

`fig, axes = plt.subplots(2, 2, figsize=(12, 8))`
 rows cols

`axes = axes.flatten()`

for i, col in enumerate(cols):

`df[col].hist(bins=50, ax=axes[i], density=True, alpha=0.2, color='red',
 label='original-{col}')`

`new_df[col].hist(bins=50, ax=axes[i], density=True, alpha=0.2, color='red',
 label='imputed-{col}')`

`axes[i].x_label(col)`

`axes[i].y_label()`

`axes[i].legend()`

`plt.tight_layout()`

`plt.show()`