

# Pandas → for data Analysis.

## Why Pandas?

- i) simple to use
- ii) integrated with many other data science & ML Python tools
- iii) Helps you get your data ready for machine Learning

## Things to cover:

- most useful functions
- pandas Datatypes
- Importing & Exporting data
- Describing data
- Viewing & Selecting data
- Manipulating data

## where to get help?

- i) chatgpt
- ii) stack overflow
- iii) pandas documentation

# Serries, Data frames and CSV's.

## Pandas Data frame Anatomy

	make	colour	odometer	Doors	Price	Column name
0	toyota	white	150043	4	\$4000	
1	hyundai	red	87899	4	\$5000	
2	toyota	blue	32549	3	\$7000	
3	bmw	black	11179	5	\$22000	
4	nissan	white	213095	4	\$3500	

→ index starts at 0  
→ row(axis=1)  
→ Row(axis=0)

→ import pandas as pd

# 2 main data types.

i) **series** = pd.Series(["bmw", "toyota", "honda"])  
    ↳ 1-dimensional

**series**

```
>> 0 bmw
    1 toyota
    2 honda
```

ii) Dataframe = pd. DataFrame ( { "key": value , "key": value } )

↳ 2-dimensional

eg

Clothes = pd. Series ( [ "red", "blue", "green" ] )

Car-data = pd. DataFrame ( { "Car-name": Series , "Colour": Colour } )

	Car-name	Colour
0	bmw	red
1	toyota	blue
2	honda	green

[ row  $\Rightarrow$  axis = 0  
column  $\Rightarrow$  axis = 1 ]

we can also create dataframes from csv, excel

→ import data

eg df = pd.read\_csv ('car-sales.csv')

→ export data frame

→ df. to\_excel ('exported-car-sales.xlsx')

→ df. to\_csv ('exported-car-sales.csv', [ index = False ])

( by default it is True )

add an "unnamed" index col

importing data from URLs:

heart-disease = pd.read\_csv ("https://raw.URL")

↳ some CSV URL

## Describing data with Pandas:

### # attributes

### # functions

i) .dtypes

i) .toCSV()

= car\_sales.dtypes

ii) .describe()

>> Make object

Colour object

odometer int64

Doors int64

Price object

ii) .columns

car\_columns = car\_sales.columns

car\_columns

>> Index(['make', 'colour', 'odometer', 'Doors', 'Price'],  
dtype='object')

iii) index

car\_index = car\_sales.index

car\_index

>> RangeIndex(start=0, stop=10, step=1)

iv) .shape → get dimensions

= car\_sales.shape

→ (10, 5)

## functions

### i) .describe()

car-sales.describe()

	odometer	doors
count	10.0	10.0
mean	26.2	4.0
std	61.2	0.0
min	11	3.0
25%	35	4.0
50%	57	4.0
75%	96	4.0
Max	219	5.0

### ii) .info() → (index combined with dtypes)

car-sales.info()

→ <class 'pandas.core.frame.DataFrame'>

RangeIndex : 10 entries 0 to 9

Data columns (total 5 columns):

#	Column	Non-null count	Dtype
0	make	10 non-null	object
1	colour	10 non-null	object
2	odometer	10 non-null	int64
3	doors	10 non-null	int64
4	price	10 non-null	object

dtypes : int64(2), objects(3)

memory usage: 532.0 bytes

iii) .mean()

car\_sales.mean()

>> Odometer : 78601.0

Doors : 4.0

dtype: float64

iv) .sum()

car\_sales[ "Doors" ].sum()

>> 40

v) .median()

vi) len() (no. of rows)

>> len(car\_sales)

> 10

Selecting and viewing Data with Pandas.

i) .head([n])

↳ returns 5 top rows by default

ii) .tail([n])

↳ returns 5 last rows by default

iii) .loc , .iloc → integer based position (index)

↳ location label based

eg animals = pd.Series(["cat", "dog", "horse", "cow"],  
index=[0, 1, 1, 2])

animals.loc[1]

```
>> 1 dog  
1 horse
```

# iloc  
animals.iloc [1] → at index (ie 0, 1, 2)  
of dog horse

```
>> dog
```

# we can slice with loc & iloc.

eg animals.iloc [:2]  
→ ["cat", "dog"]

# .loc slice

```
animals.loc [:1]
```

```
0 cat  
1 dog  
2 horse
```

# selecting columns

dataframe ["ColumnName"] → preferred way

eg car\_sales ["model"] AND car\_sales. make

↳ causes issues if col-name has spaces.

# selecting columns with filters

```
dataframe [dataframe["column"] == "Toyota"]
```

eg car\_sales [ car\_sales ["make"] == "Toyota"]

>>	make	colour	odometer	doors	price
0	Toyota	white	1500	4	\$ 4000
2	Toyota	blue	32	3	\$ 7000
5	Toyota	green	99	4	\$ 4500
8	Toyota	white	60	4	\$ 1250

eg car\_sales[car\_sales["odometer"] > 100]

>>	make	colour	odometer	doors	price
0	Toyota	white	1500	4	\$ 4000

RegEx :- regular expressions.

car\_sales["Price"] = car\_sales["price"].str.replace('^\\$[1-9][0-9]\*\$', regex=True)

documentation

pandas.Series.str.replace



.crosstab( ) —> Compare two columns

P. crosstab(car\_sales["make"], car\_sales["doors"])

Doors 3 4 5

make

Bmw 0 0 1

Monda 0 3 0

Nissen 0 2 0

Toyota 1 3 0

.groupby()

g = car\_sales.groupby(["make"]).mean(numeric\_only=True)

make	odometer	doors
Benz	111	5.00
Mazda	627	4.00
Nissan	1123	4.00
Toyota	854	3.75

filtering groupby

car\_sales.groupby(["make"]).mean(numeric\_only=True).loc["Toyota"]

>>

odometer      854  
doors          3.75

dtype : Toyota    dtype: float64

plot?

car\_sales["odometer"].plot()



if matplot lib don't work

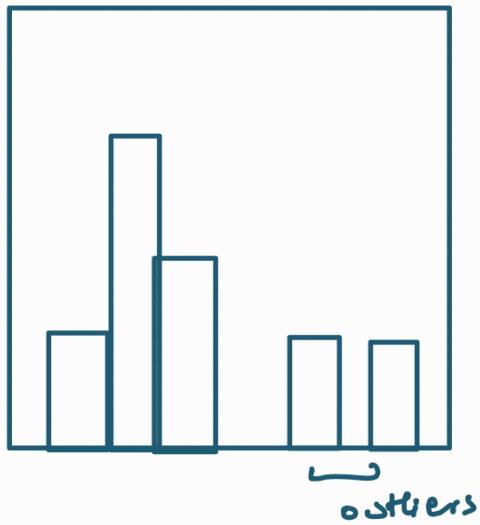
(%matplotlib inline

import matplotlib.pyplot as plt)

histogram → a fancy word for spread

eg

car\_sales["odometer"].hist()



---

Manipulating Data:

- i) .str.lower()
- ii) .str.capitalize()
- iii) .str.upper()
- iv) .fillna()
- v) .dropna()
- vi) .drop("Col\_name", axis=1)

eg

car\_sales["Miles"] = car\_sales["Miles"].str.capitalize()

eg: fill 0 inplace of NaN

car\_sales["odometer"] = car\_sales["odometer"].fillna(0.0).  
~~astype(float)~~

#

inplace parameter:

used to update data in place instead of copying and reassigning

e.g. car\_sales["odometer"].fillna(0.0, inplace=True) ~~.astype(float)~~

car\_sales.dropna(inplace=True)

↳ removes all rows that has NaN values

Create data from existing data.

# Column from series

seats\_column = pd.Series([5, 5, 5, 5, 5])

new column called seats

car\_sales["seats"] = seats\_column

car\_sales

```
>>   make    colour  odometer  owners  price  seats
      0    toyota     red       100      4    1000      5
      1    honda      blue      99      4    2000      5
      2    bmw        grey      101      4   10000      5
      3    honda      white     200      4    5000      5
      4    nissan     silver     100      4    3000      5
```

↳ car\_sales["seats"].fillna(5, inplace)

inserting values into df with lists

eg  $\text{len(df)} \rightarrow 5$

$\text{fuel_economy} = [10, 15, 12, 13, 14]$

$\text{car_sales}["\text{fuel_economy}"] = \text{fuel_economy}$

$\text{car_sales}$

make	colour	odometer	doors	Price	seats	fuel economy
0	toyota	100	4	1000	5	10
1	honda	99	4	2000	5	15
2	bmw	101	4	10000	5	12
3	honda	200	4	5000	5	13
4	nissan	100	4	3000	5	14

---

# creating a column from single value

eg  $\text{car_sales}["\text{wheels}"] = 4$

$\text{car_sales}$

>>

make	colour	odometer	doors	Price	seats	fuel economy	wheels
0	toyota	100	4	1000	5	10	4
1	honda	99	4	2000	5	15	4
2	bmw	101	4	10000	5	12	4
3	honda	200	4	5000	5	13	4
4	nissan	100	4	3000	5	14	4

---

drop column:

$\text{car_sales} = \text{car_sales}.drop("col_name", axis=1)$

→ random sample

.sample (frac=0-1)

eg :-

car-sales-shuffled = car-sales.sample (frac=1)

# only select 20% of data

car-sales-shuffled.sample (frac=0.2)

→ reset index

car-sales-shuffled.reset\_index (drop=False, inplace=True)

---

.apply() → applies a function to a column

car-sales['odometer (km)'] = car-sales['odometer (km)'].apply  
(lambda x : x / 1.6)

---

.rename → to rename columns

df = df.rename (columns = {'odometer (km)': 'odometer (miles)'})

---