

# DATA SCIENCE

It is the process of analyzing a large set of data points to get answers to questions related to that dataset.



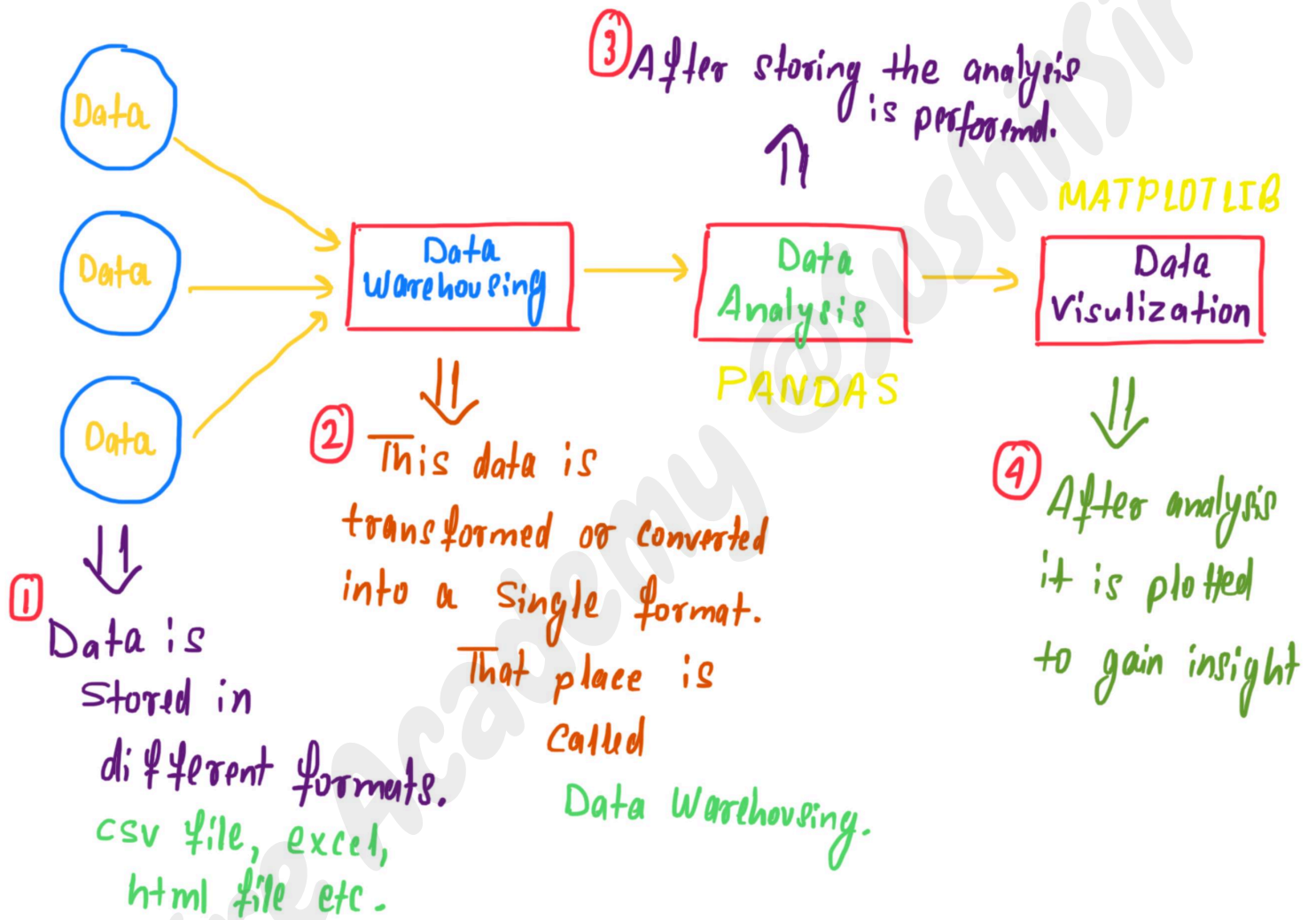
How to handle this big data?

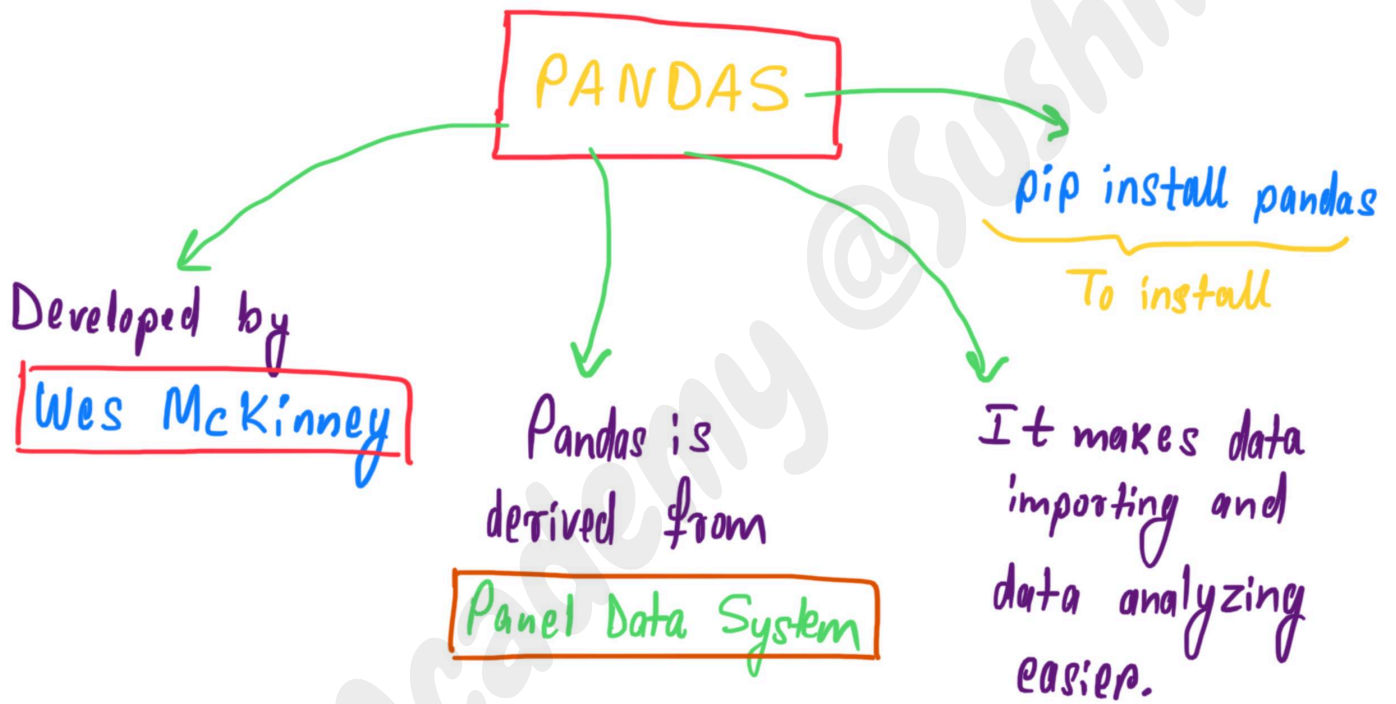


It can be handled through

Data Life Cycle (DLC)

# Data Life Cycle





# Pandas module doesn't come bundled with Standard Python.

# PIP stands for "Preferred Installer Program"

# PyPI stands for "Python Package Index".

# To use pandas → "import pandas as pd"

## Data Representation:

It can easily represent data in a form naturally suited for data analysis through

Data Frame

Series

PANDAS

It deliver fast performance and can be speed up more by making use of

Cython

(C-extension of Python)

3 well-established Python libraries

Numpy → Numerical Python  
Pandas → data analysis  
matplotlib → visualization

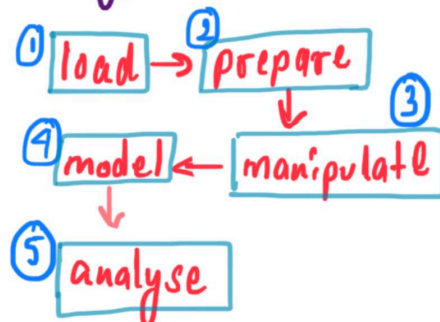
high-performance  
open source library  
for data analysis

developed by

Wes McKinney

in 2008.

using it we can accomplish five typical steps





# DATA STRUCTURE IN PANDAS



## Definition of DS

Data structure is defined as the storage and management of data for its efficient and easy access in the future where the data is collected, modified and the various types of operations are performed on data.



Pandas provide two data structures

import pandas as pd

### SERIES

- 1-D (Similar to array, list)
- assign label/index to each item
- By default (index) → 0 to N
- $N = \text{series length} - 1$
- `pd.Series(data, index, dtype, copy)`
- Label of series are called index.
- Store homogeneous data
- Data mutable (can be modified)
- Size immutable (cannot be modified)

### DATA FRAME

- tabular data structure comprised of row and column.
- two different indexes ↗ row  
↘ column
- 2-D (excel sheet)
- stores heterogeneous data
- `pd.DataFrame(data, index, columns, dtype, copy)`
- data & size (both mutable)
- `pd.DataFrame()`

import pandas as pd

# SERIES

```
S1 = np.arange(10, 15)  
S = pd.Series(S1**2,  
              index=S)
```

pd.Series (data, index = idx)

① List

```
l = [1, 2, 3, 4]  
i = ['a', 'b', 'c', 'd']
```

```
s = pd.Series(l, index=i)
```

```
S['a']  
S[['a', 'b', 'c']]
```

Access  
value

```
S[0]  
S[:3]  
S[-3:]
```

Slicing

```
S.iloc[1:4]  
S.loc['b':'d']
```

position-based

name-based

② ndarray

③ scalar value

```
• pd.Series(10, index = range(0, 3))  
• pd.Series(15, index = range(1, 6, 2))
```

④ dictionary

```
dict = {'Jan': 31, 'Feb': 28, 'Mar': 31}
```

```
s = pd.Series(dict)
```

```
s.name = "Days"
```

```
s.index.name = "Month"
```

⑤ mathematical function

```
S1 = np.arange(10, 15)
```

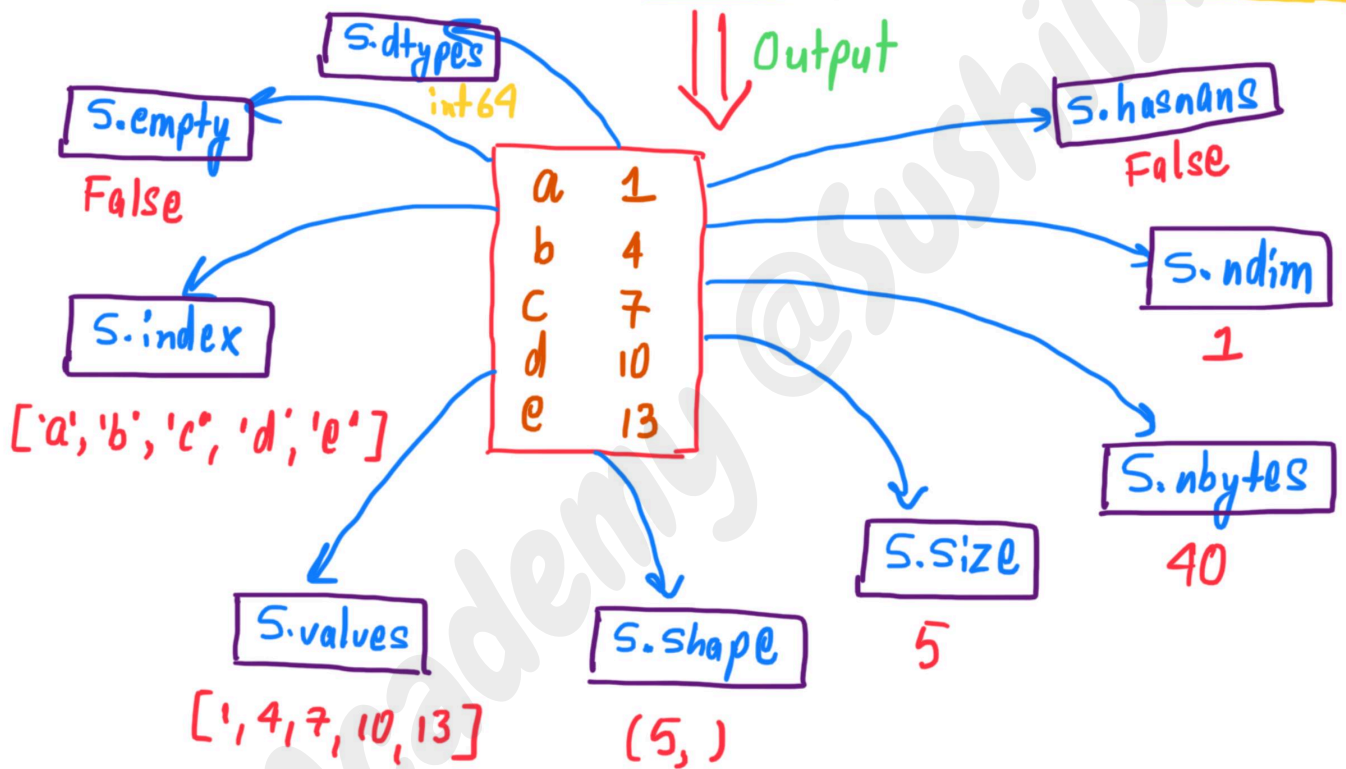
```
S = pd.Series(data = S1*4,  
              index = S1)
```

# SERIES OBJECT ATTRIBUTES

```
import pandas as pd
```

CODE

```
S = pd.Series(range(1,15,3), index=[x for x in 'abcde'])
```



## head() & tail functions

S.head()

It gives first five data by default

S.tail()

It gives last five data by default



# MATHEMATICAL OPERATION ON SERIES

**S**

1	11
2	12
3	13
4	14

pd.Series([11,12,13,14],  
index = range(1,5))

index

data

**S1**

1	21
2	22
3	23
4	24

pd.Series([21,22,23,24],  
index = range(1,5))

**S2**

101	21
102	22
103	23
104	24

pd.Series([21,22,23,24],  
index = range(101,105))

**S + S1**

1	32
2	34
3	36
4	38

**S + S2**

1	NaN
2	NaN
3	NaN
4	NaN
101	NaN
102	NaN
103	NaN
104	NaN

**S \* S1**

1	231
2	264
3	294
4	336

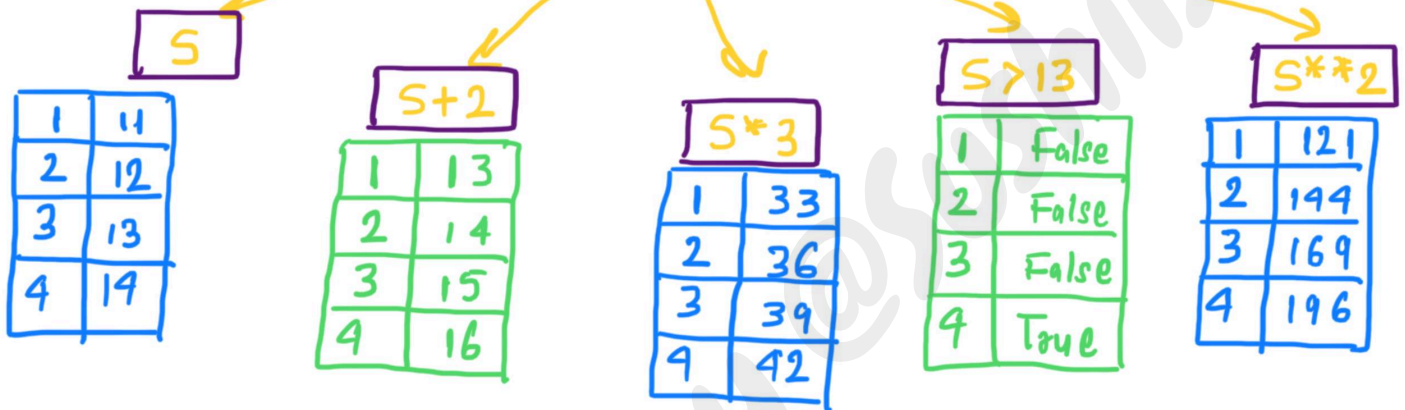
**S / S1**

1	0.523810
2	0.545455
3	0.565217
4	0.583333



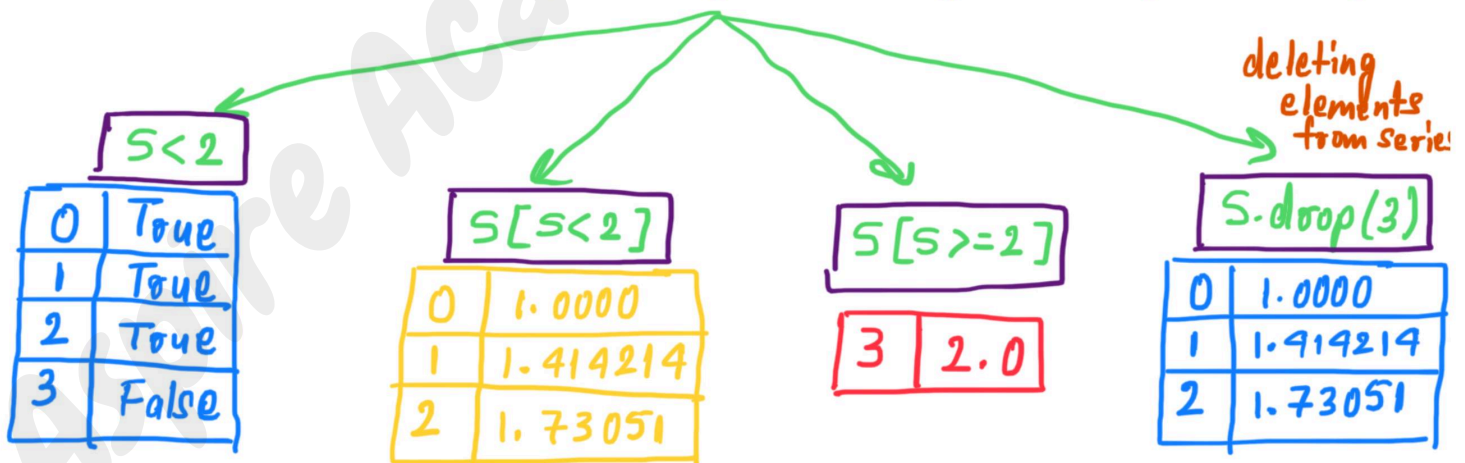
# VECTOR OPERATIONS ON SERIES

`S = pd.Series([11, 12, 13, 14], index=range(1, 5))`



# RETRIEVING VALUES USING CONDITIONS

`S = pd.Series([1.0000, 1.414214, 1.73051, 2.0000])`



x End of Series DS x