

# Python Packages

# Roadmap

---

- 1 **Array Data Structure**
- 2 Data Visualization
- 3 Pandas
- 4 Descriptive Analytics in Python

# Array Data Structure

---



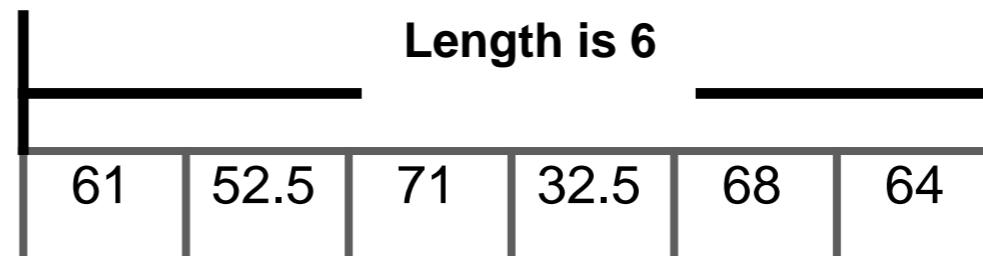
“Fighting with a large army under your command is nowise different from fighting with a small one: it is merely a question of instituting signs and signals.”

—*Sun Tzu on the Art of War*

# Array Data Structure

---

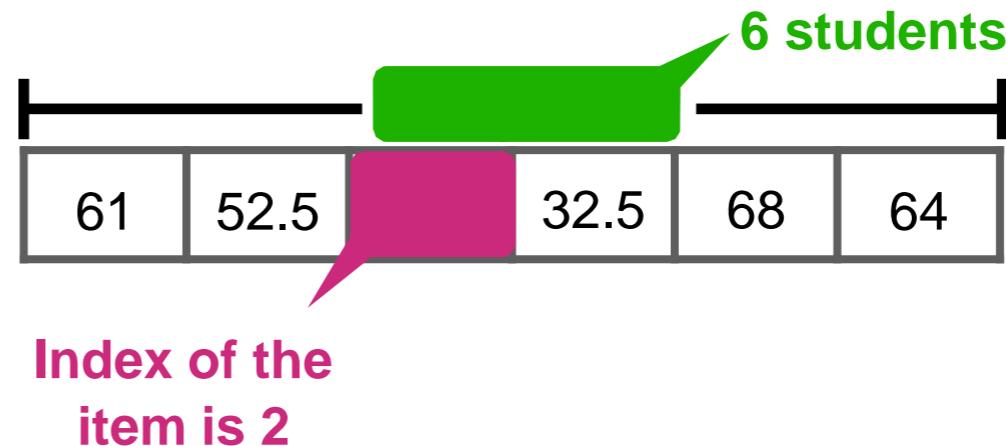
- Concepts
  - ▶ One-dimensional arrays
    - ✓ Final exam grades of six students



# Array Data Structure

---

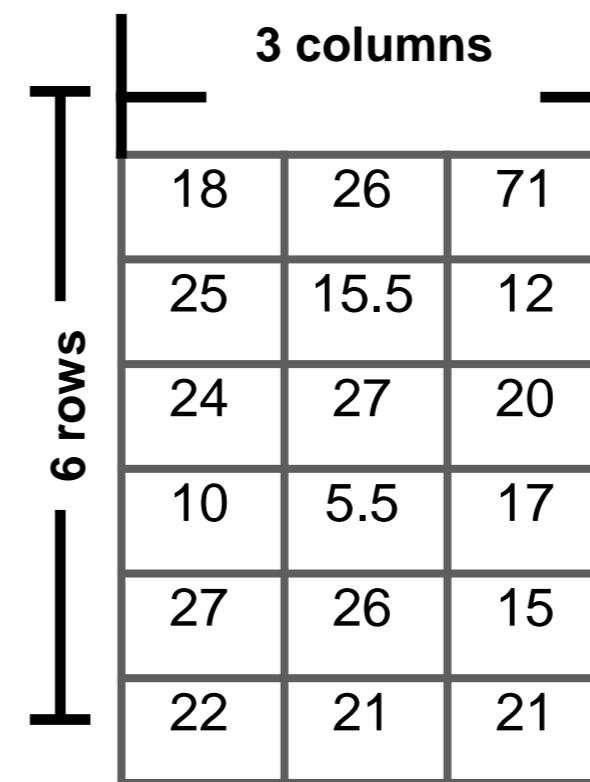
- Concepts
  - ▶ One-dimensional arrays
    - ✓ Final exam grades of six students



# Array Data Structure

---

- Concepts
  - ▶ Two-dimensional arrays
    - ✓ Grades of three components of six students



18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

# Array Data Structure

- Concepts
  - ▶ Two-dimensional arrays
    - ✓ Grades of three components of six students

18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

3 components

3 columns

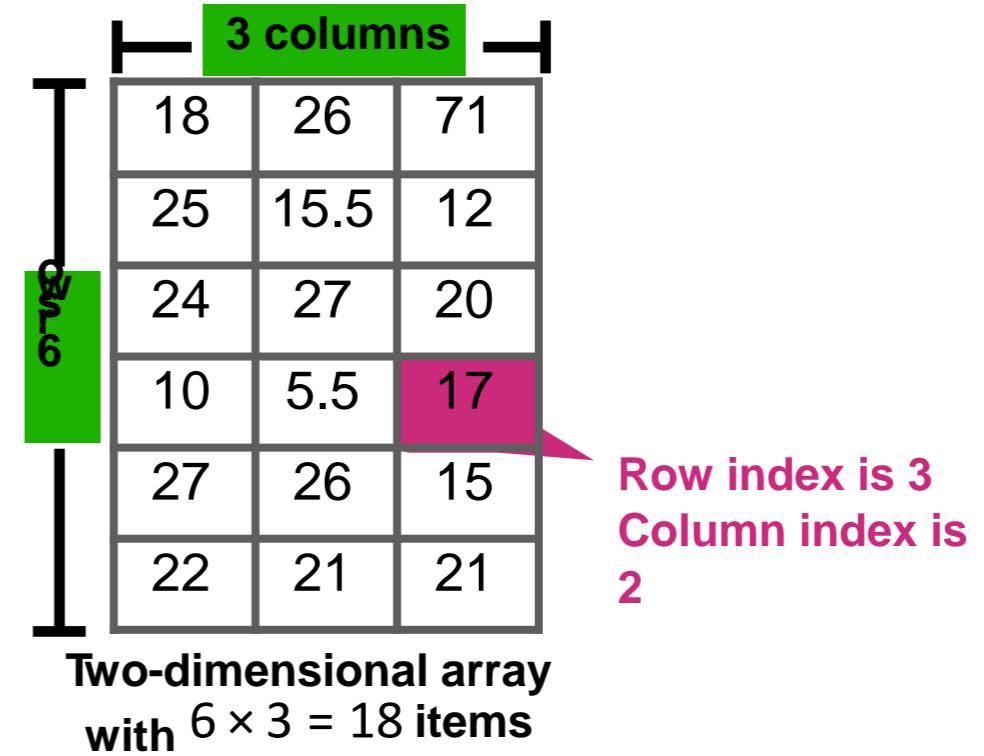
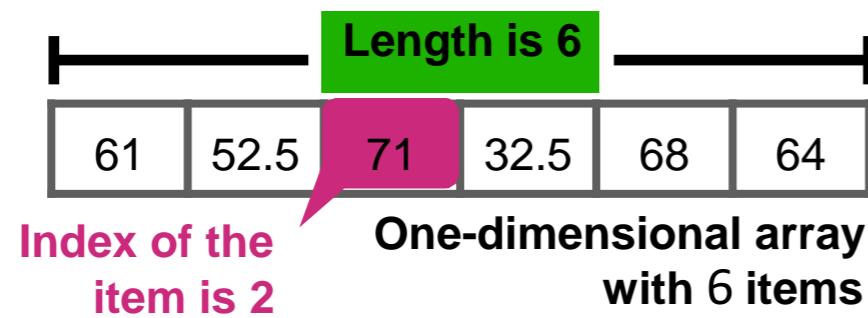
6 students

Row index is 3  
Column index is 2

# Array Data Structure

- Concepts

- $N$ -dimensional arrays



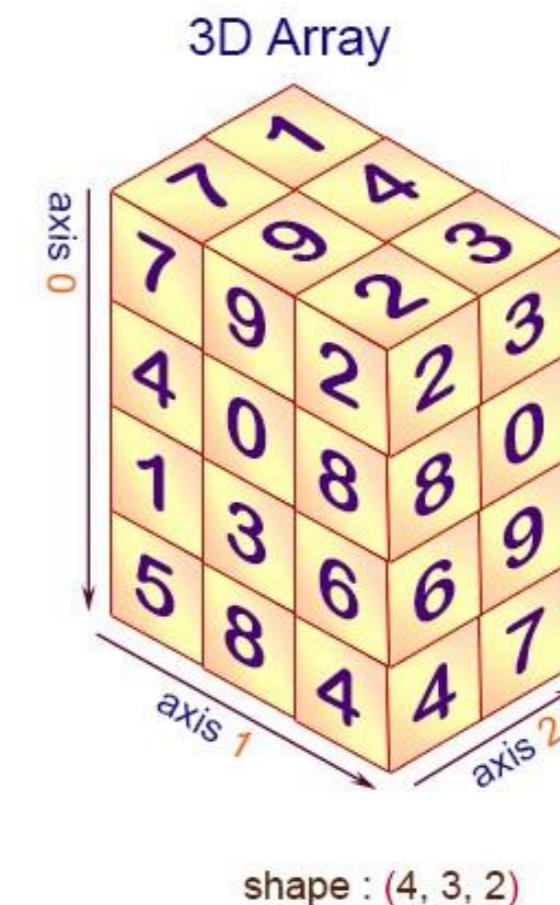
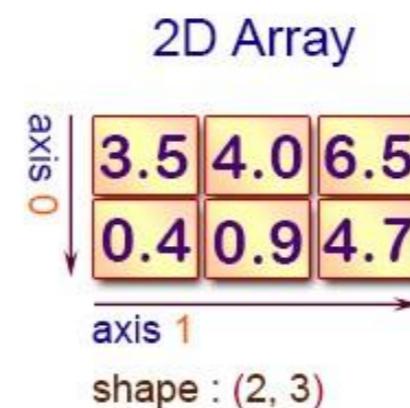
Dimension number	Shape	Index	Item number
1	One number	One index	Shape
2	Two numbers	Two indexes	$\text{Shape}[0] * \text{Shape}[1]$
$N$	$N$ numbers	$N$ indexes	$\text{Shape}[0] * \dots * \text{Shape}[N]$

# Array Data Structure

- Concepts

- ▶  $N$ -dimensional arrays

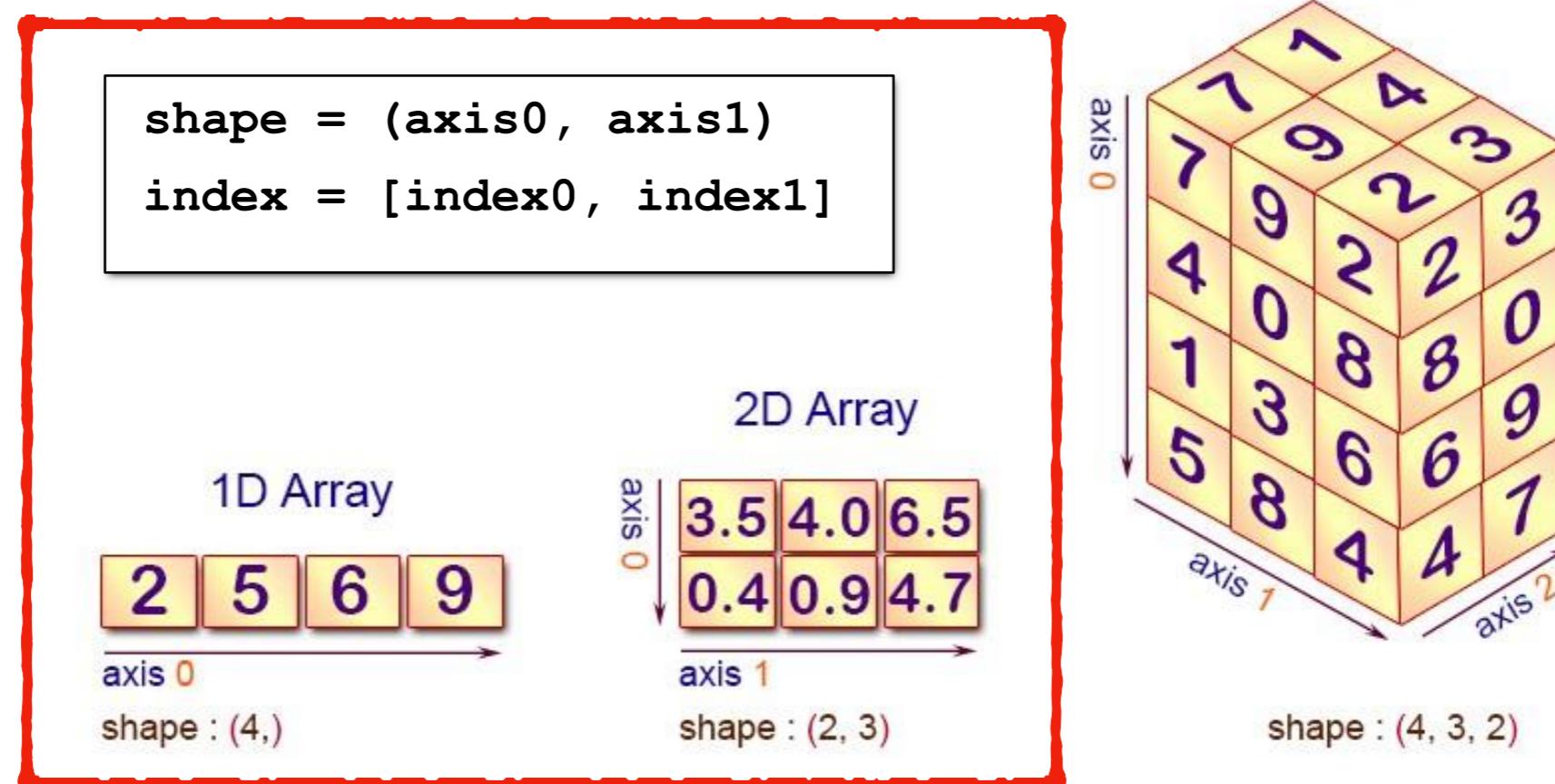
```
shape = (axis0, axis1, ...))  
index = [index0, index1, ...]
```



# Array Data Structure

- Concepts

- ▶  $N$ -dimensional arrays



# Array Data Structure

---

- NumPy arrays
  - ▶ Import the NumPy package

```
import numpy as np  
np.<function name>
```

# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists

✓ Function

```
1 array_1d = np.array([61, 52.5, 71, 32.5, 68, 64])    # This is a 1D array
2 array_2d = np.array([[18, 26, 17],
3                      [25, 15.5, 12],
4                      [24, 27, 20],
5                      [10, 5.5, 17],
6                      [27, 26, 15],
7                      [22, 21, 21]])           # This is a 2D array
```

```
1 print(type(array_1d))
2 print(type(array_2d))
```

```
<class 'numpy.ndarray'>
<class 'numpy.ndarray'>
```

# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists

✓ Function **np.array()**



```
1 array_1d = np.array(61, 52.5, 71, 32.5, 68, 64)      # This is a 1D array
2 array_2d = np.array([[18, 26, 17],
3                      [25, 15.5, 12],
4                      [24, 27, 20],
5                      [10, 5.5, 17],
6                      [27, 26, 15],
7                      [22, 21, 21]])      # This is a 2D array
```

# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists

✓ Function

```
1 array_1d = np.array([61, 52.5, 71, 32.5, 68, 64]) # This is a 1D array
2 array_2d = np.array(
```



18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

*is a 2D array*

# Array Data Structure

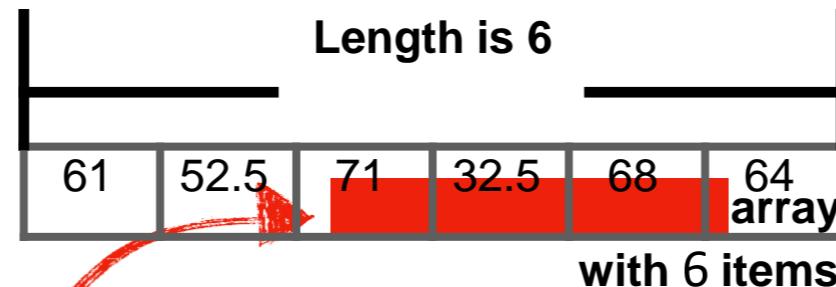
- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays

```
1 # Attributes of the one-dimensional array
2
3 print('Dimension number: ', )
4 print('Array shape: ', )
5 print('Number of data items: ', )
6 print('Data type of items: ', )
```

```
Dimension number: [red]
Array shape: [green]
Number of data items: [pink]
Data type of items: [blue]
```

# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists
  - ✓ Attributes of arrays

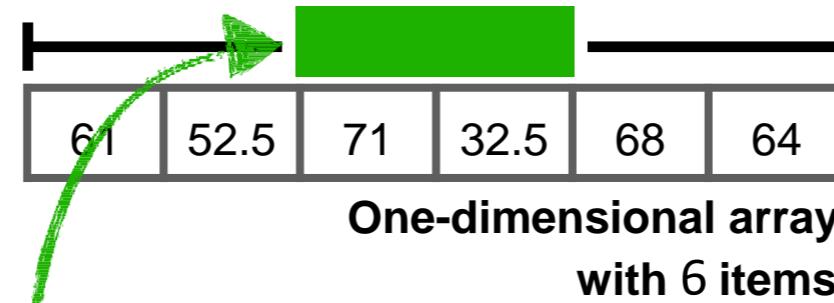


Dimension number:  
Array shape:  
Number of data items:  
Data type of items:

6  
(6,)  
6  
float64

# Array Data Structure

- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays

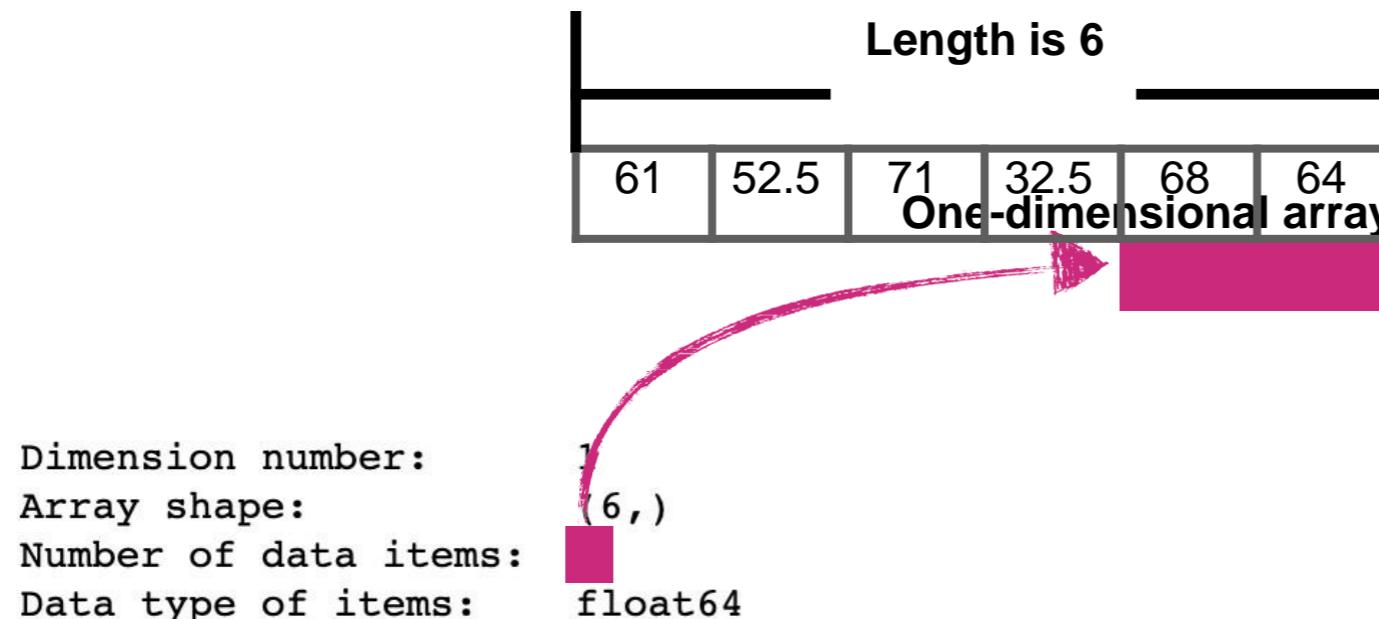


Dimension number: 1  
Array shape: (6,)  
Number of data items: 6  
Data type of items: float64

A tuple with one item

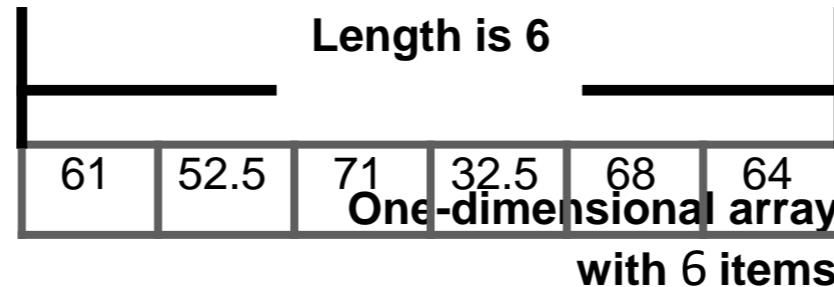
# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists
  - ✓ Attributes of arrays



# Array Data Structure

- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays



Dimension number: 1  
Array shape: (6,)  
Number of data items: 6  
Data type of items:

**Data type: floating point numbers**

# Array Data Structure

- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays

```
1 # Attributes of the two-dimensional array
2
3 print('Dimension number: ', )
4 print('Array shape: ', )
5 print('Number of data items: ', )
6 print('Data type of items: ', )
```

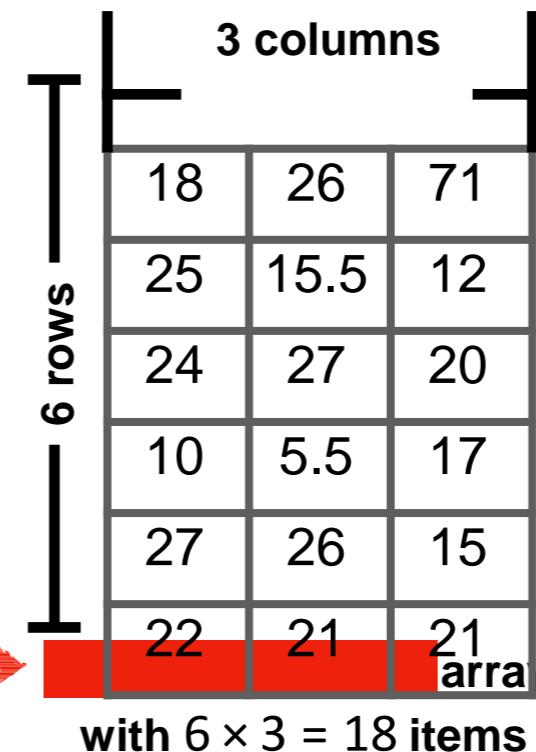
Dimension number:   
Array shape:   
Number of data items:   
Data type of items: 

# Array Data Structure

- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays

Dimension number:  
Array shape:  
Number of data items:  
Data type of items:

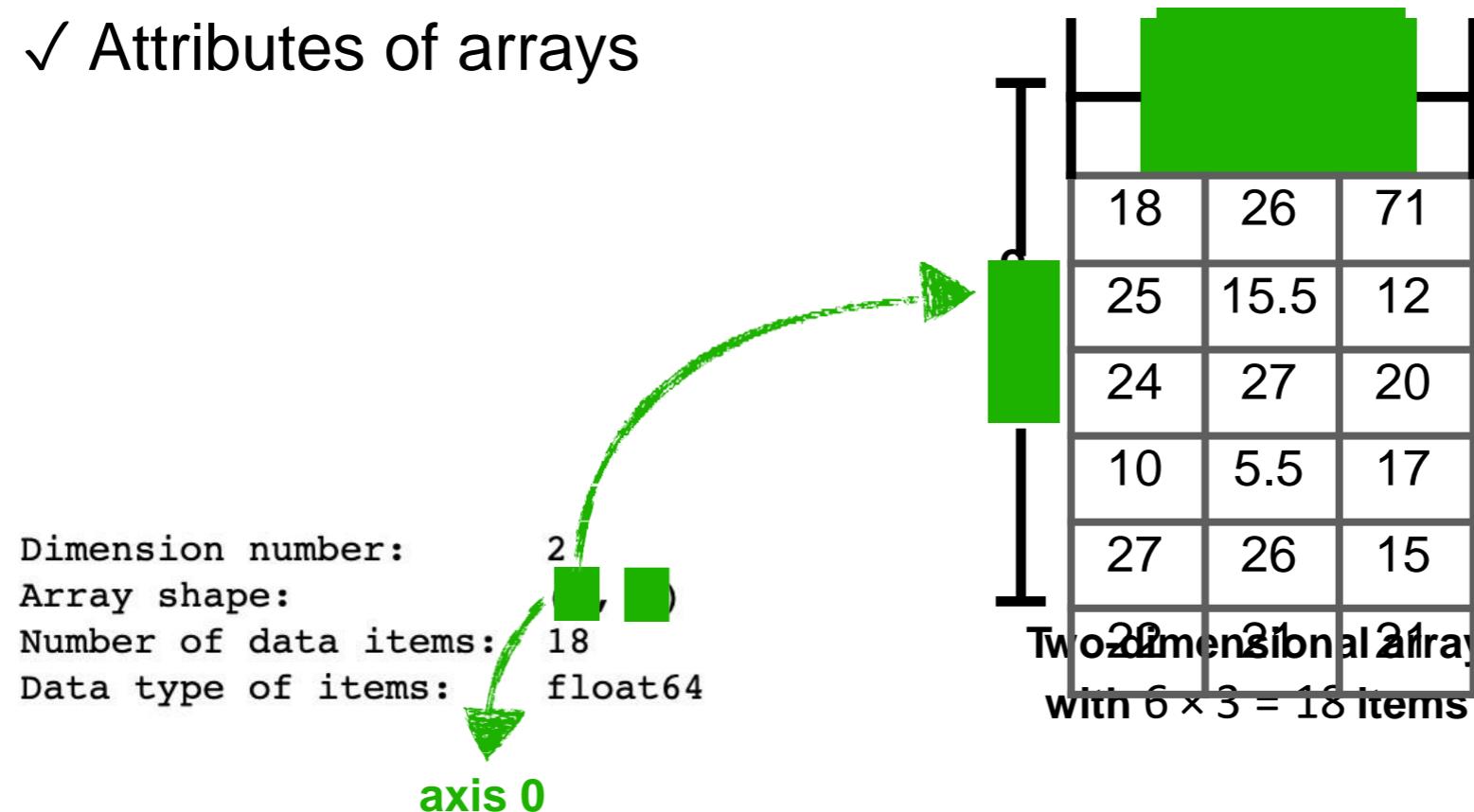
1  
(6, 3)  
18  
float64



# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists

## ✓ Attributes of arrays



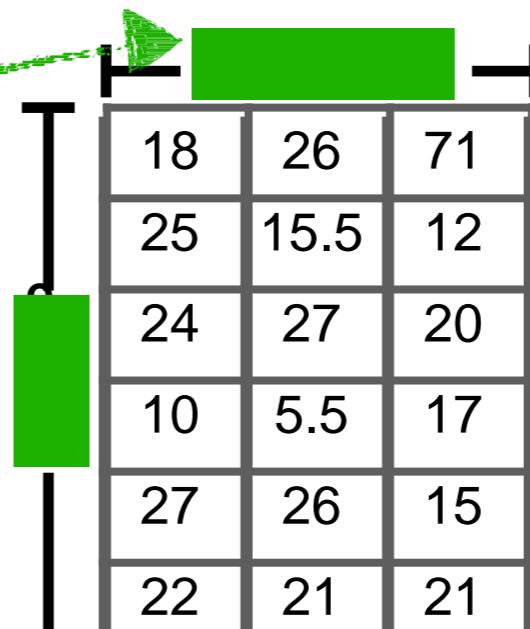
# Array Data Structure

- NumPy arrays
  - Create arrays from Python lists

✓ Attributes of arrays

Dimension number:  
Array shape:  
Number of data items:  
Data type of items:

2  
18  
float64



18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

Two-dimensional array  
with  $6 \times 3 = 18$  items

# Array Data Structure

- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays

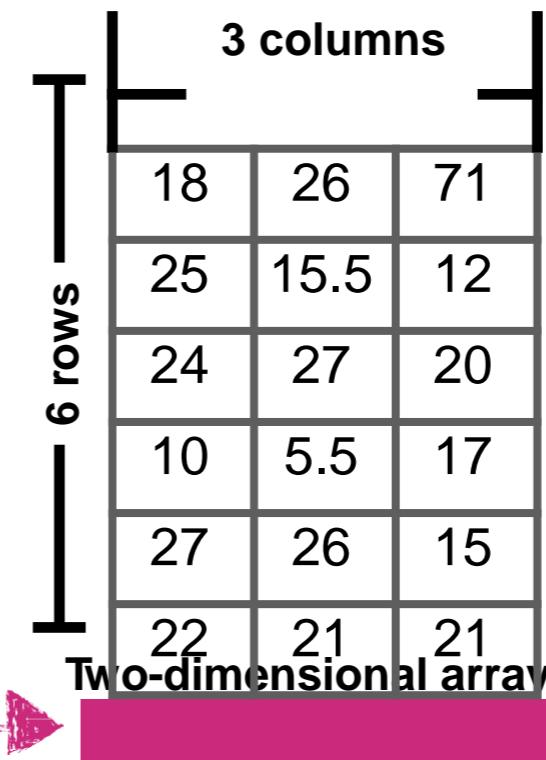
Dimension number:

Array shape:

Number of data items:

Data type of items:

2  
(6, 3)  
float64



# Array Data Structure

- NumPy arrays
  - ▶ Create arrays from Python lists
  - ✓ Attributes of arrays

Dimension number: 2  
Array shape: (6, 3)  
Number of data items: 18  
Data type of items:

A 2D array diagram with 6 rows and 3 columns. The columns are labeled '3 columns' and the rows are labeled '6 rows'. The array contains the following data:

18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

Two-dimensional array with  $6 \times 3 = 18$  items

Data type: floating point numbers

# Array Data Structure

- NumPy arrays

- Create special

- arrays**

```
1 ones_2d = np.ones( )  
2 ones_2d
```

```
array([[1., 1., 1., 1., 1.],  
       [1., 1., 1., 1., 1.],  
       [1., 1., 1., 1., 1.]])
```

A two-dimensional array:  
3 rows and 5 columns

```
1 zeros_1d = np.zeros()  
2 zeros_1d
```

```
array([0., 0., 0., 0., 0.])
```

A one-dimensional array: length is  
5

```
1 range_array = np.arange()  
2 range_array
```

```
array([2. , 2.5, 3. , 3.5, 4. , 4.5])
```

(start, stop, step)

# Array Data Structure

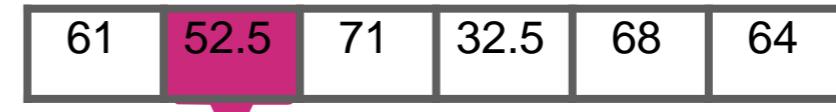
- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
9 print()
10 print(array_2d[0, 1])
11 print(array_1d[2:])
12 print(array_2d[3:, :])
13 print(array_2d[[0, 2, 1]])
```

```

26.0
[71. 32.5 68. 64. ]
[[10. 5.5 17. ]
 [27. 26. 15. ]
 [22. 21. 21. ]]
[[18. 26. 17. ]
 [24. 27. 20. ]
 [25. 15.5 12. ]]
```



Index: 1

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

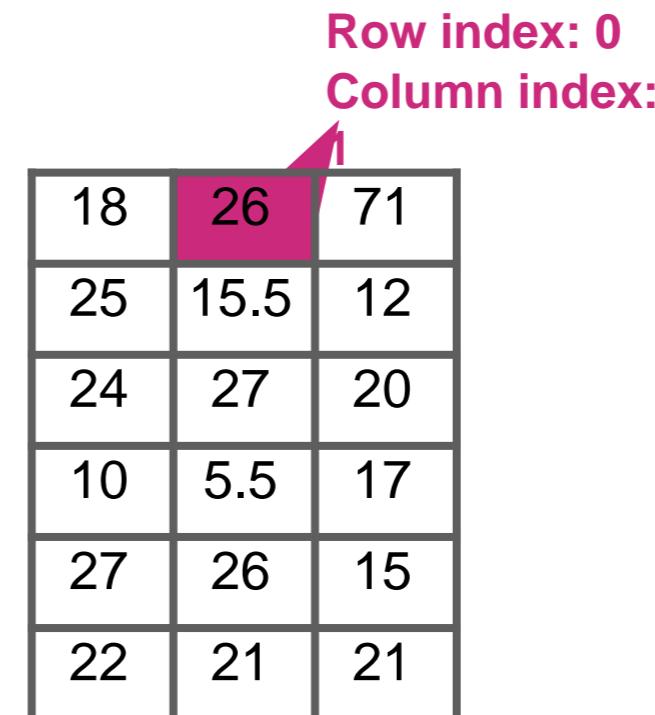
✓ Similar to  type data

```
9 print(array_1d[1])
10 print(
11 print(array_1d[2:])
12 print(array_2d[3:, :])
13 print(array_2d[[0, 2, 1]])
```

```
52.5

[71. 32.5 68. 64. ]
[[10. 5.5 17. ]
 [27. 26. 15. ]
 [22. 21. 21. ]]
[[18. 26. 17. ]
 [24. 27. 20. ]
 [25. 15.5 12. ]]
```

Row index: 0  
Column index:



18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
9 print(array_1d[1])
10 print(array_2d[0, 1])
11 print(
12 print(array_2d[3:, :])
13 print(array_2d[[0, 2, 1]])
```

```
52.5
26.0
[[10.  5.5 17. ]
 [27.  26.  15. ]
 [22.  21.  21. ]]
[[18.  26.  17. ]
 [24.  27.  20. ]
 [25.  15.5 12. ]]
```



Index:  
2:6:1

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
9 print(array_1d[1])
10 print(array_2d[0, 1])
11 print(array_1d[2:])
12 print(
13 print(array_2d[[0, 2, 1]]))
```

```
52.5
26.0
[71. 32.5 68. 64. ]

[[18. 26. 17. ]
 [24. 27. 20. ]
 [25. 15.5 12. ]]
```

18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

Row index: 3:6:1  
Column index:  
0:3:1

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
9 print(array_1d[1])
10 print(array_2d[0, 1])
11 print(array_1d[2:])
12 print(array_2d[3:, :])
13 print()
```

```
52.5
26.0
[71. 32.5 68. 64. ]
[[10. 5.5 17. ]
 [27. 26. 15. ]
 [22. 21. 21. ]]
```

18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

Row index: [0, 2, 1]  
Column index:  
0:3:1

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
9 print(array_1d[1])
10 print(array_2d[0, 1])
11 print(array_1d[2:])
12 print(array_2d[3:, :])
13 print( 0:3:1)
```

```
52.5
26.0
[71. 32.5 68. 64. ]
[[10. 5.5 17. ]
 [27. 26. 15. ]
 [22. 21. 21. ]]
```

Take all columns if  
column index is  
not specified

18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

Row index: [0, 2, 1]  
Column index:

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
1 array_2d =   
2 array_2d  
  
array([[18. , 26. , 17. ],  
       [25. , 15.5, 12. ],  
       [25. , 28. , 21. ],  
       [10. , 5.5, 17. ],  
       [27. , 26. , 15. ],  
       [22. , 21. , 21. ]])
```

```
1 array_1d[2] = 74  
2 array_1d
```

```
array([61. , 52.5, 74. , 32.5, 68. , 64. ])
```

18	26	71
25	15.5	12
24	27	20
10	5.5	17
27	26	15
22	21	21

Items to be  
changed

# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
1 array_2d = 
2 array_2d

array([[18. , 26. , 17. ],
       [25. , 15.5, 12. ],
         


```
1 array_1d[2] = 74
2 array_1d

array([61. , 52.5, 74. , 32.5, 68. , 64. ])
```


```

18	26	71
25	15.5	12
25	28	21
10	5.5	17
27	26	15
22	21	21

 New item values

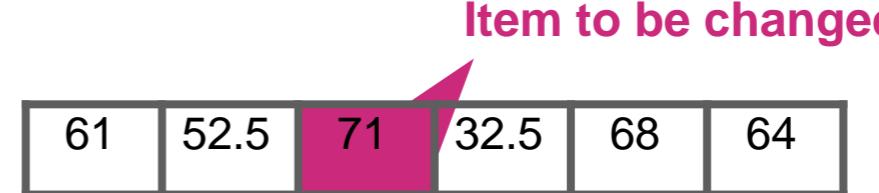
# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
1 array_2d[2, :] = np.array([25, 28, 21])
2 array_2d
```

```
array([[18. , 26. , 17. ],
       [25. , 15.5, 12. ],
       [25. , 28. , 21. ],
       [10. , 5.5, 17. ],
       [27. , 26. , 15. ],
       [22. , 21. , 21. ]])
```



```
1 array_1d[2] = 74
2 array_1d
```

```
array([61. , 52.5, 74. , 32.5, 68. , 64. ])
```

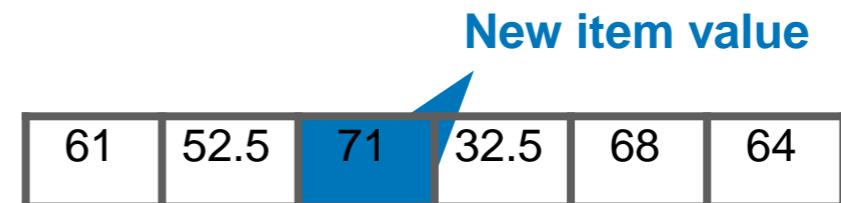
# Array Data Structure

- NumPy arrays
  - Indexing and slicing of arrays

✓ Similar to  type data

```
1 array_2d[2, :] = np.array([25, 28, 21])
2 array_2d
```

```
array([[18. , 26. , 17. ],
       [25. , 15.5, 12. ],
       [25. , 28. , 21. ],
       [10. , 5.5, 17. ],
       [27. , 26. , 15. ],
       [22. , 21. , 21. ]])
```



```
1 array_1d[2] = 71
2 array_1d
```

```
array([61. , 52.5, 71. , 32.5, 68. , 64. ])
```

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations
    - ✓ An array and a scalar

1	2
2	3.5
5	6.5

$$\begin{array}{c} + \quad 3 \quad = \end{array} \quad \begin{array}{c} 4 \quad 5 \\ 5 \quad 6.5 \\ 8 \quad 9.5 \end{array}$$

1	2
2	3.5
5	6.5

$$\begin{array}{c} * \quad 2 \quad = \end{array} \quad \begin{array}{c} 2 \quad 4 \\ 4 \quad 7 \\ 10 \quad 13 \end{array}$$

```
1 array_2d = np.array([[1, 2],  
2 [2, 3.5],  
3 [5, 6.5]])  
4  
5  
6  
7 print(array_2d + array_2d)  
8 print(array_2d * array_2d)
```

  
[[ 2. 4.]  
 [ 4. 7.]  
 [10. 13.]]  
[[ 1. 4.]  
 [ 4. 12.25]  
 [25. 42.25]]

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations
    - ✓ An array and a scalar

1	2
2	3.5
5	6.5

1	2
2	3.5
5	6.5

+

2	4
4	7
10	13

1	2
2	3.5
5	6.5

1	2
2	3.5
5	6.5

\*

1	4
4	12.25
25	42.25

```
1 array_2d = np.array([[1, 2],  
2                      [2, 3.5],  
3                      [5, 6.5]])  
4  
5 print(array_2d + 3)  
6 print(array_2d * 2)  
7  
8
```

```
[[4.  5. ]  
 [5.  6.5]  
 [8.  9.5]]  
[[ 2.  4.]  
 [ 4.  7.]  
 [10. 13.]]
```

# Array Data Structure

- NumPy arrays
  - ▶ Element-wise arithmetic operations

## Example 1

usd is a list of containing five money transactions in US dollars.  
Transfer each transaction into Singapore dollars.

```
import numpy as np

usd = [2, 3.60, 2.05, 13.50, 18.90]
exchange_rate = 1.37

usd_array = np.array(usd)                      # Convert list into a numpy array
sgd_array = usd_array * exchange_rate           # Element-wise multiplication

print(sgd_array)

[ 2.74    4.932   2.8085  18.495   25.893 ]
```

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

2	3.60	2.05	13.50	18.90
---	------	------	-------	-------

```
import numpy as np

exchange_rate = 1.37

sgd_array = usd_array * exchange_rate      # Convert list into a numpy array
                                            # Element-wise multiplication

print(sgd_array)

[ 2.74      4.932     2.8085  18.495   25.893 ]
```

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

2	3.60	2.05	13.50	18.90	* 1.37
---	------	------	-------	-------	--------

```
import numpy as np

usd = [2, 3.60, 2.05, 13.50, 18.90]
exchange_rate = 1.37

usd_array = np.array(usd)          # Convert list into a numpy array
                                    # Element-wise multiplication

print(sgd_array)

[ 2.74    4.932   2.8085  18.495  25.893 ]
```

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

2.74	4.932	2.8085	18.495	25.893
------	-------	--------	--------	--------

```
import numpy as np

usd = [2, 3.60, 2.05, 13.50, 18.90]
exchange_rate = 1.37

usd_array = np.array(usd)                      # Convert list into a numpy array
sgd_array = usd_array * exchange_rate           # Element-wise multiplication

print(sgd_array)
```

# Array Data Structure

---

- NumPy arrays
  - ▶ Element-wise arithmetic operations
    - ✓ More concise
    - ✓ Easier to read
    - ✓ Faster to execute than loops

# Array Data Structure

---

- NumPy arrays
  - ▶ Element-wise arithmetic operations

## Example 2

The discrete distribution information of two types of newspapers is stored in a two-dimensional array `distr`. The first row represents the probabilities of each weather condition, the second and third rows give the corresponding newspaper demands. Calculate 1) the expected newspaper demands; 2) the standard deviations of newspaper demands; and 3) the expected total profit with the given parameters

# Array Data Structure

---

- NumPy arrays
  - Element-wise arithmetic operations

	Sunny	Cloudy		Thunderstorm	Haze
Probabilities	0.315	0.226			
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

	Sunny	Cloudy		Thunderstorm	Haze
Probabilities	0.315	0.226	Raining	0.289	0.087
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array

```
1 exp1 = np.sum(distr[0] * distr[1])
2 var1 = np.sum(distr[0] * (distr[1] - exp1)**2)
3 std1 = var1 ** 0.5
4
5 print('Expected demand: {:.3f}'.format(exp1))
6 print('Standard deviation: {:.3f}'.format(std1))
```

Expected demand: 449.249  
Standard deviation: 118.908

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

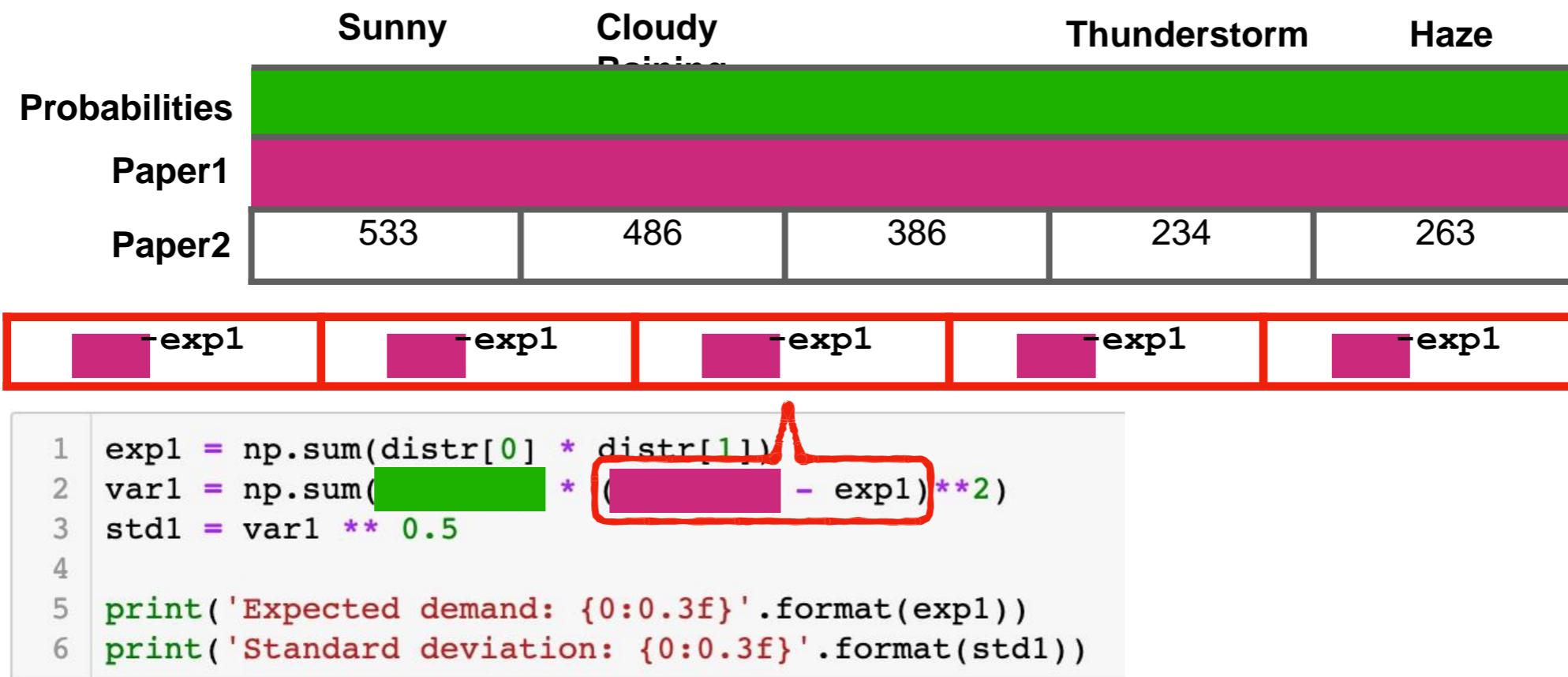
	Sunny	Cloudy	Thunderstorm	Haze
Probabilities	0.35	0.30	0.20	0.15
Paper1	533	486	386	234
Paper2	533	486	386	234
Two-dimensional array				

```
1 exp1 = np.sum(distr[0] * distr[1])
2 var1 = np.sum(distr[0] * (distr[1] - exp1)**2)
3 std1 = var1 ** 0.5
4
5 print('Expected demand: {:.3f}'.format(exp1))
6 print('Standard deviation: {:.3f}'.format(std1))
```

Expected demand: 449.249  
Standard deviation: 118.908

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations



Expected demand: 449.249  
Standard deviation: 118.908

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

	Sunny	Cloudy	Thunderstorm	Haze
Probabilities	0.35	0.30	0.20	0.15
Paper1	533	486	386	234
Paper2	533	486	386	234
	$(p_i - \text{exp1})^{**2}$	$(p_i - \text{exp1})^{**2}$	$(p_i - \text{exp1})^{**2}$	$(p_i - \text{exp1})^{**2}$
	<img alt="A red box highlights the calculation of squared differences for the first four categories (Sunny, Cloudy, Thunder			

# Array Data Structure

- NumPy arrays

## ► Element-wise arithmetic operations

A horizontal bar chart comparing the number of news items for four weather conditions (Sunny, Cloudy, Rainy, Thunderstorm) across two different news sources (Paper1 and Paper2). The y-axis represents the 'Probabilities' of the news items. The x-axis categories are 'Sunny', 'Cloudy', 'Rainy', 'Thunderstorm', and 'Haze'. Paper1 is represented by a green bar, and Paper2 by a pink bar. The values for Paper1 are 560, 530, 389, 202, and 278 respectively. The values for Paper2 are 450, 400, 350, 200, and 150 respectively.

Weather Condition	Paper1	Paper2
Sunny	560	450
Cloudy	530	400
Rainy	389	350
Thunderstorm	202	200
Haze	278	150

```
1 exp2 = np.sum( * )
2 var2 = np.sum( * ( - exp2) ** 2 )
3 std2 = var2 ** 0.5
4
5 print('Expected demand: {:.3f}'.format(exp2))
6 print('Standard deviation: {:.3f}'.format(std2))
```

Expected demand: 431.472  
Standard deviation: 101.316

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations

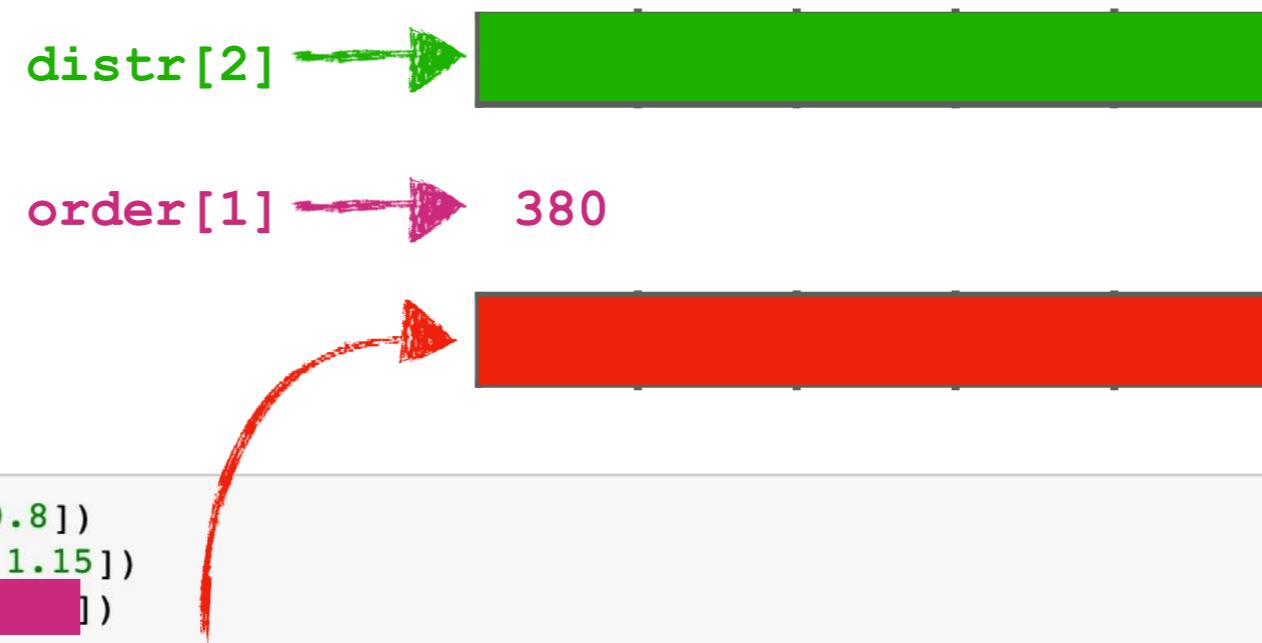


```
1 cost = np.array([0.6, 0.8])
2 price = np.array([1.0, 1.15])
3 order = np.array([430, 380])
4
5 exp_sold1 = np.sum(distr[0] * np.minimum(order[0], distr[1]))
6 exp_sold2 = np.sum(distr[0] * np.minimum(order[1], distr[2]))
7 exp_profit = (price[0]*exp_sold1 + price[1]*exp_sold2) - np.sum(order*cost)
8 print('Expected profit: ${0:0.2f}'.format(exp_profit))
```

Expected profit: \$234.92

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations



```
1 cost = np.array([0.6, 0.8])
2 price = np.array([1.0, 1.15])
3 order = np.array([430, 380])
4
5 exp_sold1 = np.sum(distr[0] * np.minimum(order[0], distr[1]))
6 exp_sold2 = np.sum(distr[0] * np.minimum(order[1], distr[1]))
7 exp_profit = (price[0]*exp_sold1 + price[1]*exp_sold2) - np.sum(order*cost)
8 print('Expected profit: ${0:0.2f}'.format(exp_profit))
```

Expected profit: \$234.92

# Array Data Structure

- NumPy arrays
  - Element-wise arithmetic operations



Expected profit: \$234.92

# Array Data Structure

---

- NumPy arrays
  - ▶ Element-wise arithmetic operations

```
1 cost = np.array([0.6, 0.8])
2 price = np.array([1.0, 1.15])
3 order = np.array([430, 380])
4
5 exp_sold1 = np.sum(distr[0] * np.minimum(order[0], distr[1]))
6 exp_sold2 = np.sum(distr[0] * np.minimum(order[1], distr[2]))
7
8 print('Expected profit: ${0:0.2f}'.format(exp_profit))
```

Expected profit:

# Array Data Structure

---

- NumPy arrays

- Broadcasting

- ✓ An array and a scalar
    - ✓ Two arrays of the same shape

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 2 & 3.5 \\ \hline 5 & 6.5 \\ \hline \end{array} + 3 = \begin{array}{|c|c|} \hline 4 & 5 \\ \hline 5 & 6.5 \\ \hline 8 & 9.5 \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 2 & 3.5 \\ \hline 5 & 6.5 \\ \hline \end{array} + \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 2 & 3.5 \\ \hline 5 & 6.5 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 2 & 4 \\ \hline 4 & 7 \\ \hline 10 & 13 \\ \hline \end{array}$$

# Array Data Structure

- NumPy arrays
    - Broadcasting
- ✓ Two arrays of different same shapes

```
1 np.ones((3, 3)) + np.arange(3)
```

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline \end{array}$$

`np.arange(3)`

`np.ones((3, 3))`

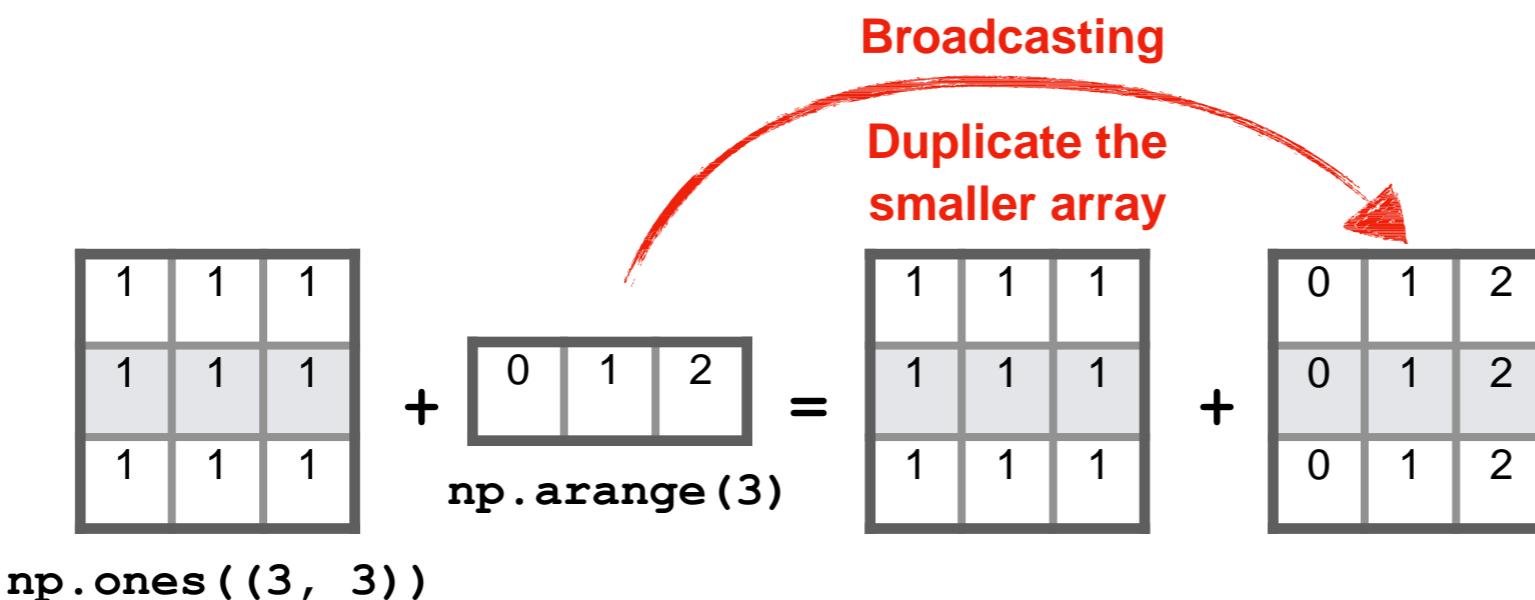
# Array Data Structure

- NumPy arrays

- Broadcasting

- ✓ Two arrays of different same shapes

```
1 np.ones((3, 3)) + np.arange(3)
```



# Array Data Structure

- NumPy arrays

- Broadcasting

- ✓ Two arrays of different same shapes

```
1 np.ones((3, 3)) + np.arange(3)
```

```
array([[1., 2., 3.],
       [1., 2., 3.],
       [1., 2., 3.]])
```

$$\begin{array}{c} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} & + & \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline \end{array} & = & \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} & + & \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 0 & 1 & 2 \\ \hline 0 & 1 & 2 \\ \hline \end{array} & = & \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 1 & 2 & 3 \\ \hline 1 & 2 & 3 \\ \hline \end{array} \\ \text{np.ones((3, 3))} & & & & & & & & \end{array}$$

# Array Data Structure

---

- NumPy arrays
    - ▶ Broadcasting
- ✓ Two arrays of different same shapes

```
1 np.arange(3).reshape((3, 1)) + np.arange(3)
```

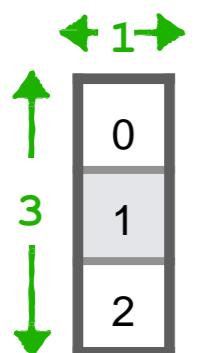
0	1	2
---	---	---

`np.arange(3)`

# Array Data Structure

- NumPy arrays
  - Broadcasting
    - ✓ Two arrays of different same shapes

```
1 np.arange(3).reshape((3, 1)) + np.arange(3)
```



Shape is changed  
from (3, ) to (3, 1)

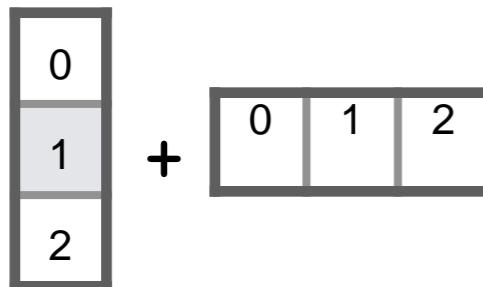
```
np.arange(3).reshape( )
```

# Array Data Structure

---

- NumPy arrays
    - Broadcasting
- ✓ Two arrays of different same shapes

```
1 np.arange(3).reshape((3, 1)) + np.arange(3)
```



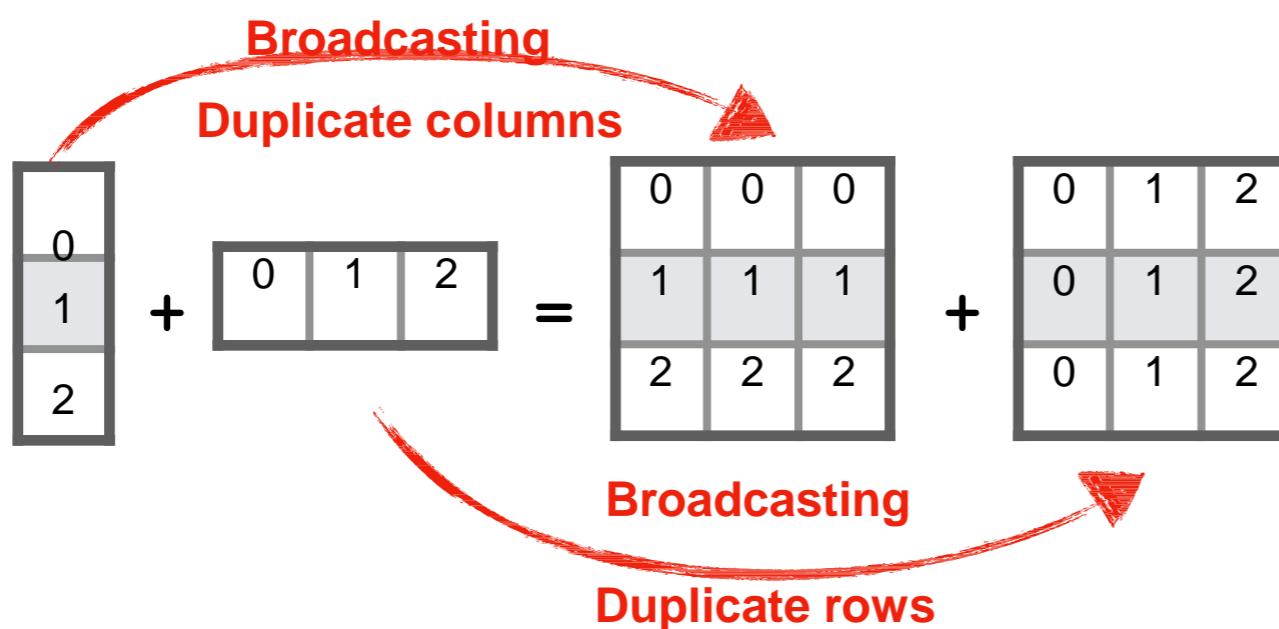
# Array Data Structure

- NumPy arrays

- Broadcasting

- ✓ Two arrays of different same shapes

```
1 np.arange(3).reshape((3, 1)) + np.arange(3)
```



# Array Data Structure

- NumPy arrays

- Broadcasting

- ✓ Two arrays of different same shapes

```
1 np.arange(3).reshape((3, 1)) + np.arange(3)
```

```
array([[0, 1, 2],  
       [1, 2, 3],  
       [2, 3, 4]])
```

$$\begin{array}{c|c|c|c|c|c|c|c|c|c} & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ \hline & 0 & & & & & & & & \\ \hline & & 1 & & & & & & & \\ \hline & & & 2 & & & & & & \\ \hline \end{array} + \begin{array}{c|c|c|c|c|c|c|c|c|c} & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ \hline & 0 & & 1 & & 2 & & & & \\ \hline & & & & & & & & & \\ & & & & & & & & & \\ \hline \end{array} = \begin{array}{c|c|c|c|c|c|c|c|c|c} & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ \hline & 0 & & 0 & & 0 & & & & \\ \hline & & 1 & & 1 & & 1 & & & \\ \hline & & & 2 & & & & & & \\ \hline & 2 & & 2 & & 2 & & & & \\ \hline \end{array} + \begin{array}{c|c|c|c|c|c|c|c|c|c} & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ \hline & 0 & & 1 & & 2 & & & & \\ \hline & & 0 & & 1 & & 2 & & & \\ \hline & & & 0 & & & & & & \\ \hline & 0 & & 1 & & 2 & & & & \\ \hline & & 1 & & 2 & & 3 & & & \\ \hline & & & 2 & & 3 & & & & \\ \hline & 2 & & 3 & & 4 & & & & \\ \hline \end{array} = \begin{array}{c|c|c|c|c|c|c|c|c|c} & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ \hline & 0 & & 1 & & 2 & & & & \\ \hline & & 1 & & 2 & & 3 & & & \\ \hline & & & 2 & & 3 & & & & \\ \hline & 2 & & 3 & & 4 & & & & \\ \hline \end{array}$$

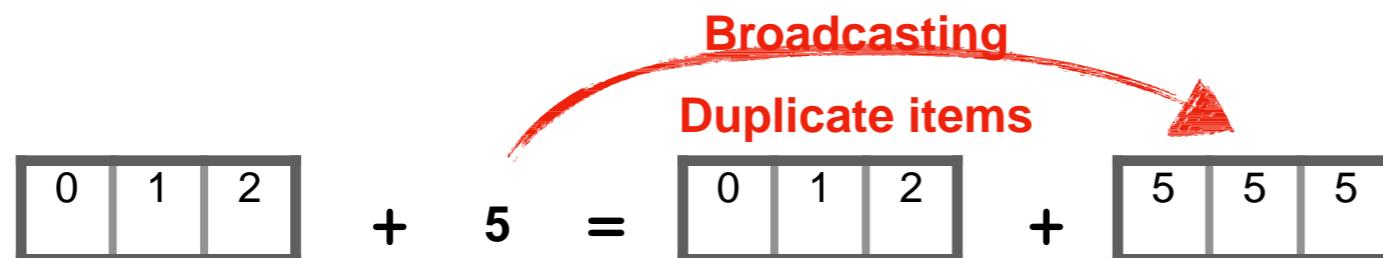
# Array Data Structure

- NumPy arrays

- Broadcasting

- ✓ An array and a scalar, interpreted by broadcasting

```
1 np.arange(3) + 5
```



# Array Data Structure

---

- NumPy arrays

- Broadcasting

- ✓ An array and a scalar, interpreted by broadcasting

```
1 np.arange(3) + 5
```

```
array([5, 6, 7])
```

$$\begin{bmatrix} 0 & 1 & 2 \end{bmatrix} + 5 = \begin{bmatrix} 0 & 1 & 2 \end{bmatrix} + \begin{bmatrix} 5 & 5 & 5 \end{bmatrix} = \begin{bmatrix} 5 & 6 & 7 \end{bmatrix}$$

# Array Data Structure

- NumPy arrays
  - Broadcasting

$$\text{np.arange}(3)+5$$

A 1D array of three blue cubes labeled 0, 1, 2 is added to a scalar 5, represented as a 1D array of three blue cubes labeled 5. The result is a 1D array of three blue cubes labeled 5, 6, 7.

$$\text{np.ones}(3, 3)+\text{np.arange}(3)$$

A 3x3 matrix of ones (3 rows, 3 columns) is added to a 1D array of range(3) (3 rows, 1 column). The result is a 3x3 matrix where each row is [1, 2, 3].

$$\text{np.ones}(3, 1)+\text{np.arange}(3)$$

A 3x1 column of ones (3 rows, 1 column) is added to a 1D array of range(3) (1 row, 3 columns). The result is a 3x1 column where each element is [0, 1, 2].

# Array Data Structure

---

- NumPy arrays
  - Functions and array methods
    - ✓ Math functions

```
1 print(np.log(3))          # Natural logarithm of 3
2 print(np.exp(np.arange(4))) # Exponential of 0, 1, 2, 3
3 print(np.square(np.arange(3))) # Squares of 0, 1, 2
4 print(np.power(np.arange(3), 3)) # Cubes of 0, 1, 2
```

```
1.0986122886681098
[ 1.           2.71828183  7.3890561  20.08553692]
[0 1 4]
[0 1 8]
```

# Array Data Structure

---

- NumPy arrays
  - ▶ Functions and array methods
    - ✓ Aggregation methods

```
1 array_2d = np.array([[1, 2],  
2                      [2, 3.5],  
3                      [5, 6.5]])      # Create a two-dimensional array  
4  
5 print(array_2d.sum())           # Sum of all item  
6 print(array_2d.max())          # The maximum item in the array  
7 print(array_2d.min())          # The minimum item in the array
```

```
20.0  
6.5  
1.0
```

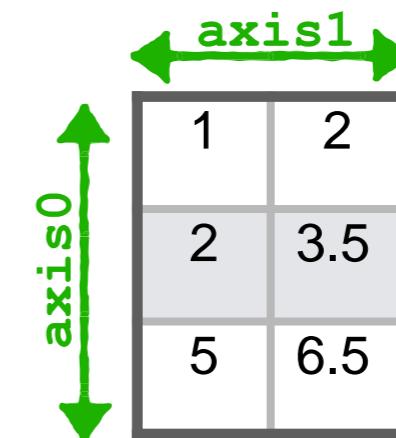
# Array Data Structure

- NumPy arrays
  - Functions and array methods

## ✓ Aggregation methods

```
1 print(array_2d.sum(axis=0))  
2 print(array_2d.sum(axis=1))  
3 print(array_2d.max(axis=0))  
4 print(array_2d.min(axis=1))
```

```
[ 8. 12.]  
[ 3. 5.5 11.5]  
[5. 6.5]  
[1. 2. 5.]
```

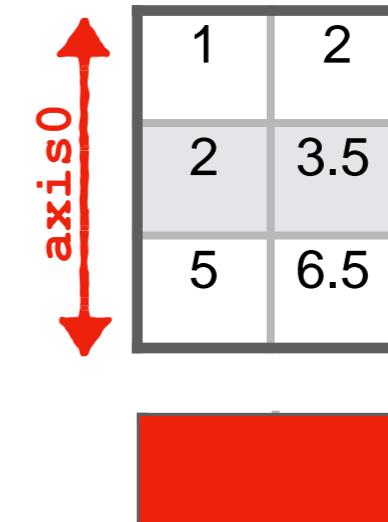


# Array Data Structure

- NumPy arrays
  - ▶ Functions and array methods
    - ✓ Aggregation methods

```
1 print(array_2d.sum())
2 print(array_2d.sum(axis=1))
3 print(array_2d.max(axis=0))
4 print(array_2d.min(axis=1))
```

```
[ 3.  5.5 11.5]
[5.  6.5]
[1.  2.  5.]
```



# Array Data Structure

- NumPy arrays
  - Functions and array methods

## ✓ Aggregation methods

```
1 print(array_2d.sum(axis=0))  
2 print(array_2d.sum())  
3 print(array_2d.max(axis=0))  
4 print(array_2d.min(axis=1))
```

```
[ 8. 12.]
```

```
[5. 6.5]
```

```
[1. 2. 5.]
```



A 3x2 grid of numbers. The first column contains 1, 2, and 5. The second column contains 2, 3.5, and 6.5. A red double-headed arrow labeled "axis1" spans the width of the grid, indicating the axis of aggregation for sum and max operations.

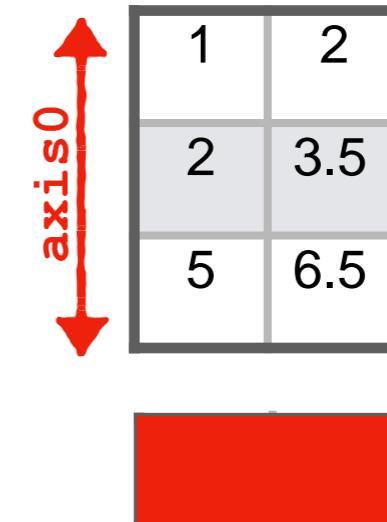
1	2
2	3.5
5	6.5

# Array Data Structure

- NumPy arrays
  - ▶ Functions and array methods
    - ✓ Aggregation methods

```
1 print(array_2d.sum(axis=0))  
2 print(array_2d.sum(axis=1))  
3 print(array_2d.max())  
4 print(array_2d.min(axis=1))
```

```
[ 8. 12.]  
[ 3. 5.5 11.5]  
[ 1. 2. 5.]
```



# Array Data Structure

- NumPy arrays
  - ▶ Functions and array methods

## ✓ Aggregation methods

```
1 print(array_2d.sum(axis=0))  
2 print(array_2d.sum(axis=1))  
3 print(array_2d.max(axis=0))  
4 print(array_2d.min())
```

```
[ 8. 12.]  
[ 3. 5.5 11.5]  
[5. 6.5]
```



# Array Data Structure

---

- NumPy arrays
  - ▶ Functions and array methods

The discrete distribution information of two types of newspapers is stored in a two-dimensional array `distr`. The first row represents the probabilities of each weather condition, the second and third rows give the corresponding newspaper demands. Calculate 1) the expected newspaper demands; 2) the standard deviations of newspaper demands; and 3) the expected total profit with the given parameters

# Array Data Structure

---

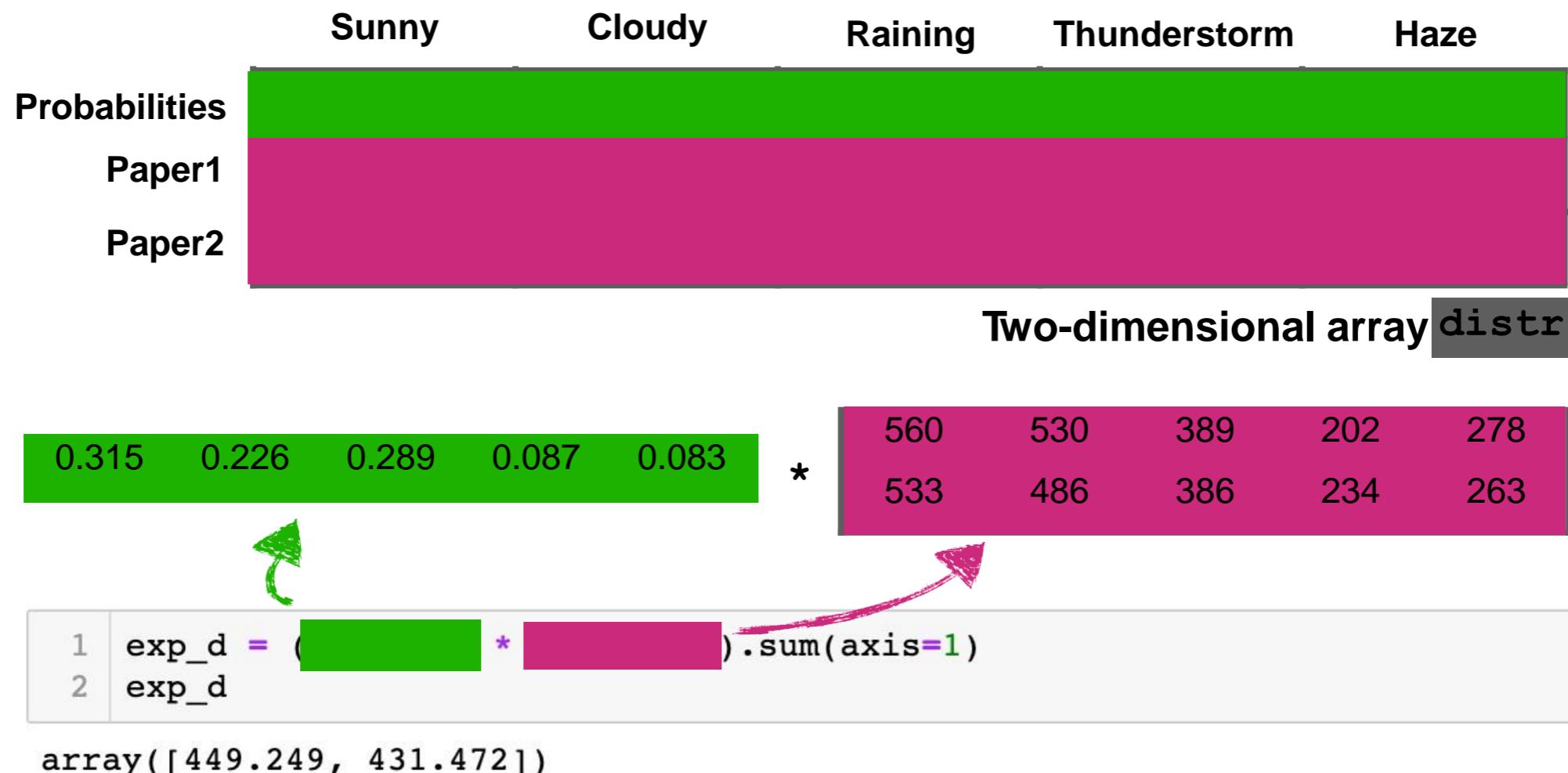
- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array

# Array Data Structure

- NumPy arrays
  - Functions and array methods



# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities					
Paper1					
Paper2					
Two-dimensional array <code>distr</code>					
0.315	0.226	0.289	0.087	0.083	560
0.315	0.226	0.289	0.087	0.083	533
					530
					486
					389
					202
					234
					278
					263

**Broadcasting**

```
1 exp_d = ( [ 0.315, 0.226, 0.289, 0.087, 0.083 ] * [ 560, 530, 389, 202, 278 ] ).sum(axis=1)
2 exp_d
```

array([ 449.249, 431.472])

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

**Two-dimensional array `distr`**

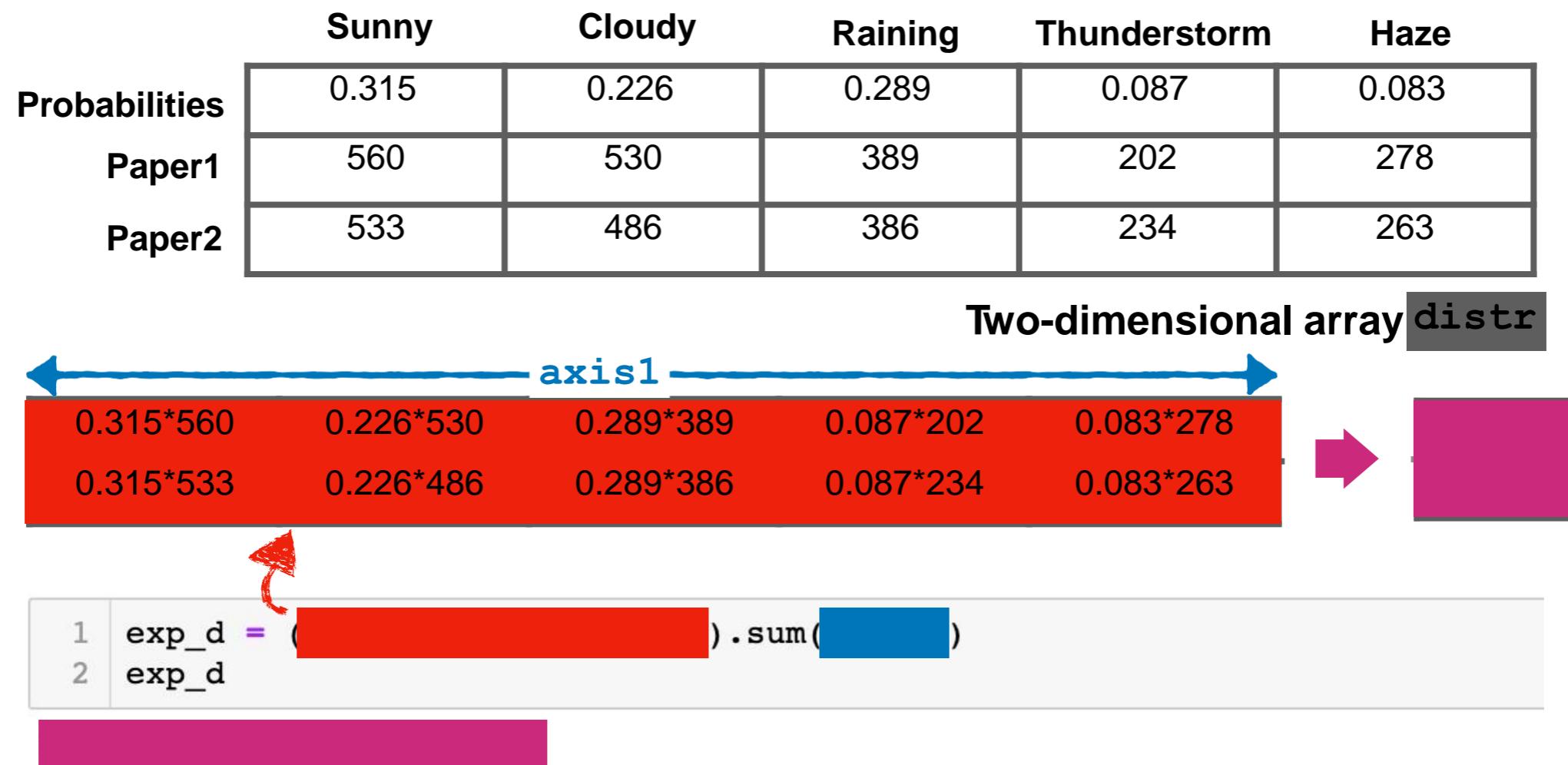
A diagram showing a 2D array with axis1 highlighted. The array is represented as a grid of numbers. A horizontal double-headed arrow labeled "axis1" spans the width of the grid, indicating the first dimension. The grid is enclosed in a red border. Above the grid, the text "Two-dimensional array `distr`" is centered. Below the grid, there is a code snippet with a red arrow pointing to the first line.

```
1 exp_d = ( ) . sum( )  
2 exp_d
```

array([449.249, 431.472])

# Array Data Structure

- NumPy arrays
  - Functions and array methods

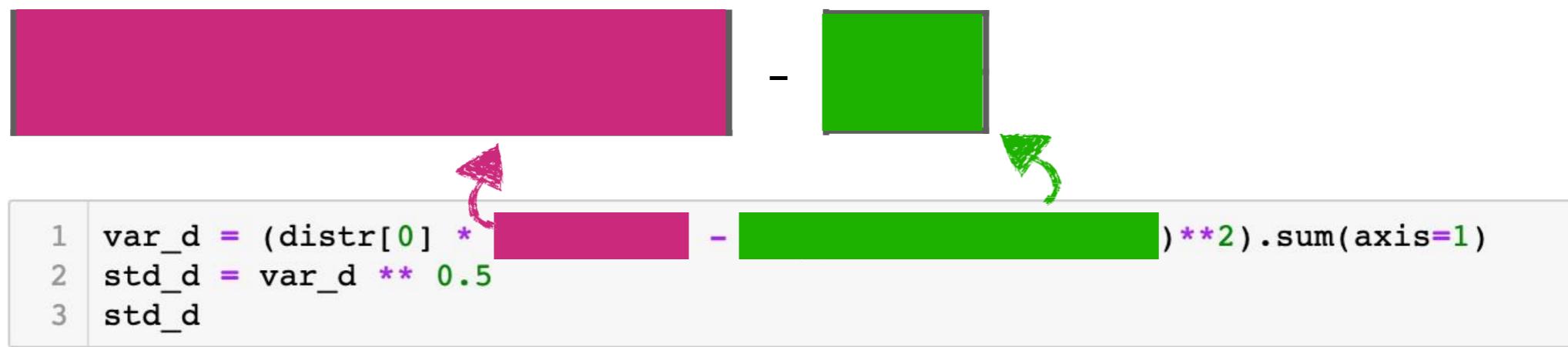


# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array



# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array

560	530	389	202	278	-	449.249	449.249	449.249	449.249	449.249
533	486	386	234	263	-	431.472	431.472	431.472	431.472	431.472

Broadcasting

1	var_d = (distr[0] * [	-	)**2).sum(axis=1)
2	std_d = var_d ** 0.5	-	
3	std_d	-	

array([118.90763222, 101.3157797 ])

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array

```
1 var_d = (distr[0] * (distr[1] - distr[0])**2).sum(axis=1)
2 std_d = var_d ** 0.5
3 std_d
```

```
array([118.90763222, 101.3157797 ])
```

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array

```
431.472)**2 (263-431.472)**2
1 var_d = (distr[0] *
2 std_d = var_d ** 0.5
3 std_d
```

```
array([118.90763222, 101.3157797 ])
```

# Array Data Structure

- NumPy arrays
    - ▶ Functions and array methods

Two-dimensional array `distr`

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities					
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

axis1

```
1 var_d = ( (560-449.249)**2 + (530-449.249)**2 + (389-449.249)**2 + (202-449.249)**2 + (278-449.249)**2 + (533-431.472)**2 + (486-431.472)**2 + (386-431.472)**2 + (234-431.472)**2 + (263-431.472)**2 ) .sum() * 0.5
2 std_d = var_d ** 0.5
3 std_d
```

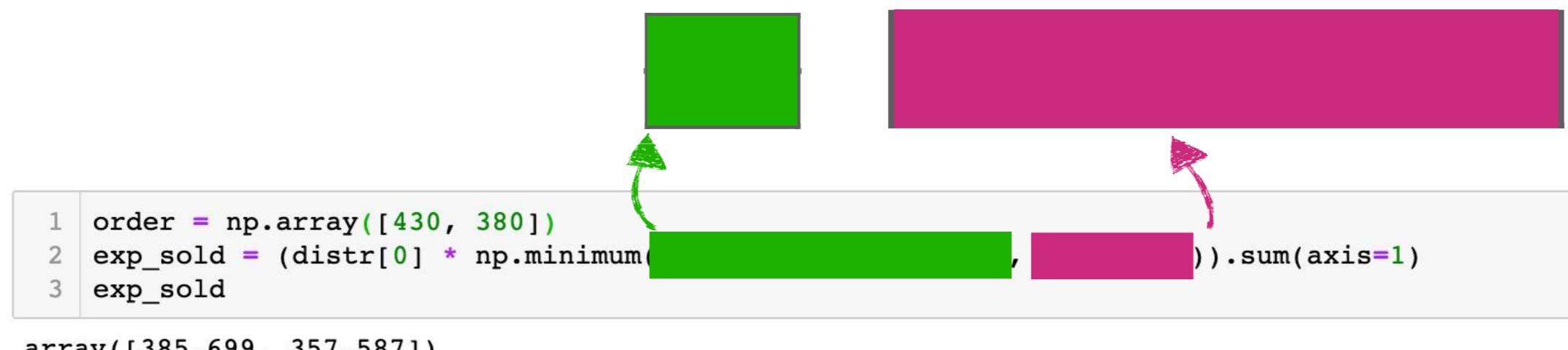
```
array([118.90763222, 101.3157797 ])
```

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1					
Paper2					

**Two-dimensional array**



# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1					
Paper2					

Two-dimensional array



Broadcasting

```
1 order = np.array([430, 380])
2 exp_sold = (distr[0] * np.minimum(distr[1], order)).sum(axis=1)
3 exp_sold
```

array([385.699, 357.587])

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities					
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array `distr`

$$\begin{bmatrix} 560 & 530 & 389 & 202 & 278 \end{bmatrix} * \begin{bmatrix} 430 & 430 & 389 & 202 & 278 \\ 380 & 380 & 380 & 234 & 263 \end{bmatrix}$$

```
1 order = np.array([430, 380])
2 exp_sold = order * distr
3 exp_sold
```

array([385.699, 357.587])

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities					
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array `distr`

$$\begin{bmatrix} 0.315 & 0.226 & 0.289 & 0.087 & 0.083 \\ 0.315 & 0.226 & 0.289 & 0.087 & 0.083 \end{bmatrix} * \begin{bmatrix} 430 & 430 & 389 & 202 & 278 \\ 380 & 380 & 380 & 234 & 263 \end{bmatrix}$$

Broadcasting

```
1 order = np.array([430, 380])
2 exp_sold = order * distr
3 exp_sold
array([385.699, 357.587])
```

# Array Data Structure

- NumPy arrays
  - Functions and array methods

	Sunny	Cloudy	Raining	Thunderstorm	Haze
Probabilities	0.315	0.226	0.289	0.087	0.083
Paper1	560	530	389	202	278
Paper2	533	486	386	234	263

Two-dimensional array `distr`



```
1 order = np.array([430, 380])
2 exp_sold = (
3     exp_sold
)
array( )
```

array( )

order = np.array([430, 380])  
exp\_sold = (  
 exp\_sold  
)  
array( )

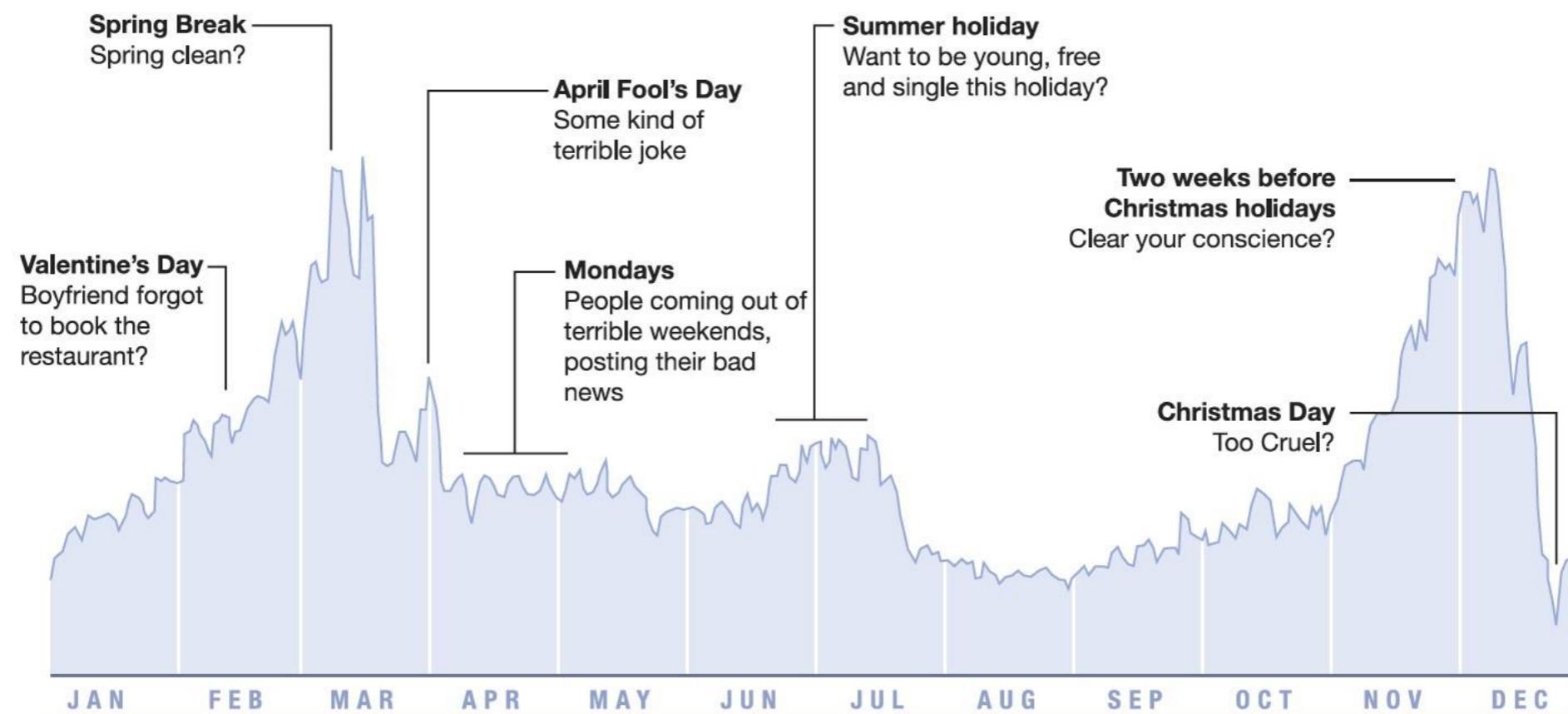
# Roadmap

---

- 1 Array Data Structure
- 2 **Data Visualization**
- 3 Pandas
- 4 Descriptive Analytics in Python

# Data Visualization

## Peak Break-up Times According to Facebook status updates



# Data Visualization

---

- Data Visualization Packages

- ▶ Import the Matplotlib package

```
import matplotlib.pyplot as plt  
plt.<function name>
```

- ▶ Other package

✓ Plotly

✓ Seaborn

# Data Visualization

---

- Functions `plot`, `line`, and `contour`  
)

- Line plot

✓ `plot` information of `line`

## Parameters

---

`x, y :`

The horizontal / vertical coordinates of the data points.  
`*x*` values are optional and default to `range(len(y))`.

Commonly, these parameters are 1D arrays.

They can also be scalars, or two-dimensional (in that case, the columns represent separate data sets).

These arguments cannot be passed as keywords.



# Data Visualization

---

- Functions  $\boxed{\quad}$ ,  $\boxed{\quad}$ , and  $\boxed{\quad}$   
    ► Line plot

## Example 3

Plot the graph of an exponential function  $y = \exp(x)$ , where  $x = 0, 1, 2, 3, 4, 5$ .

# Data Visualization

- Functions `plot`, `scatter`, and `bar`  
)

- Line plot

## Example 3

Plot the graph of an exponential function  $y = \exp(x)$ , where  $x = 0, 1, 2, 3, 4, 5$ .

```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array 0, 1, 2, 3, 4, 5
3 y = np.exp(x)                      # An array of exp(0), exp(1), ..., exp(5)
4
5 plt.plot(x, y)                    # Specify the line plot by the x and y data
6 plt.xlabel('x')                   # Label for x data
7 plt.ylabel('y')                   # Label for y data
8 plt.show()                        # Ready for display (optional for Jupyter)
```

# Data Visualization

---

- Functions `plot`, `scatter`, and `bar`

- ▶ Line plot

`x`

```
1 x = np.linspace(0, 5, 6)          # An array 0, 1, 2, 3, 4, 5
2 y = np.exp(x)                  # An array of exp(0), exp(1), ..., exp(5)
3
4
5 plt.plot(x, y)                  # Specify the line plot by the x and y data
6 plt.xlabel('x')                  # Label for x data
7 plt.ylabel('y')                  # Label for y data
8 plt.show()                      # Ready for display (optional for Jupyter)
```

# Data Visualization

- Functions `plot`, `scatter`, and `bar`

- ▶ Line plot



```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array 0, 1, 2, 3, 4, 5
3 y = np.exp(x)                      # An array of exp(0), exp(1), ..., exp(5)
4
5 plt.plot(x, y)                    # Specify the line plot by the x and y data
6 plt.xlabel('x')                   # Label for x data
7 plt.ylabel('y')                   # Label for y data
8 plt.show()                        # Ready for display (optional for Jupyter)
```

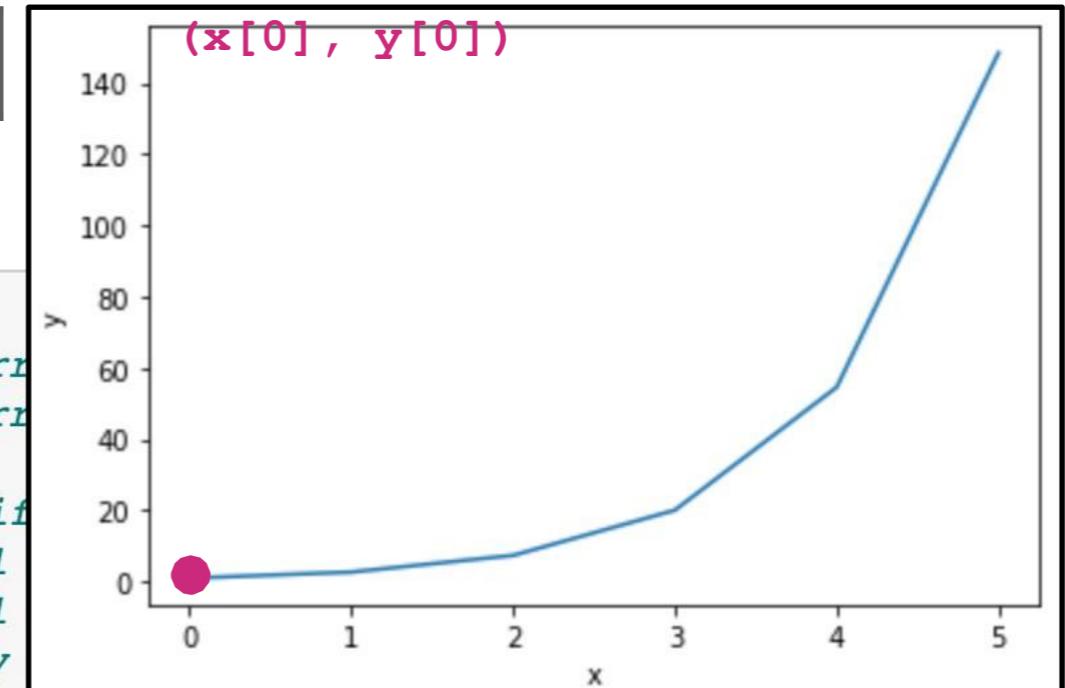
# Data Visualization

- Functions `plot`, `line`, and `bar`)

- Line plot



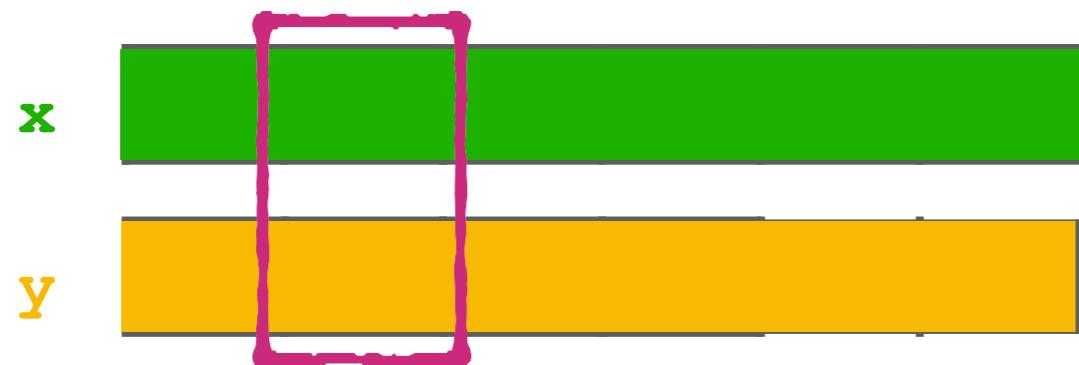
```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array
3 y = np.exp(x)                      # An array
4
5 plt.xlabel('x')                    # Specifying
6 plt.ylabel('y')                    # Label
7 plt.show()                         # Label
8
```



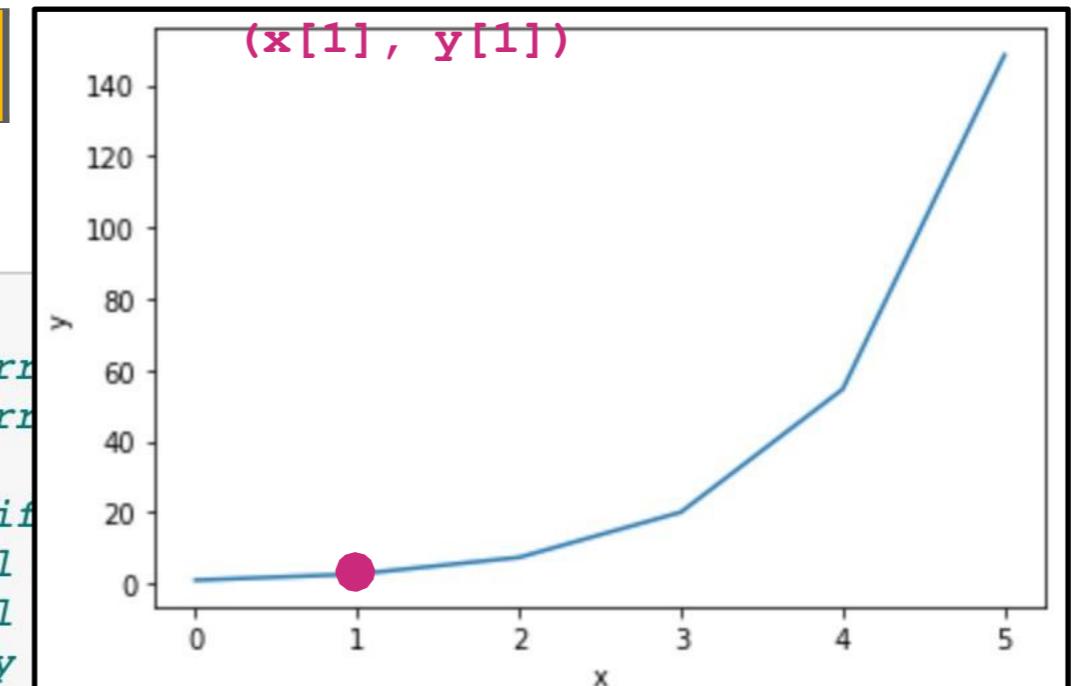
# Data Visualization

- Functions `plot`, `line`, and `bar`  
)

- Line plot



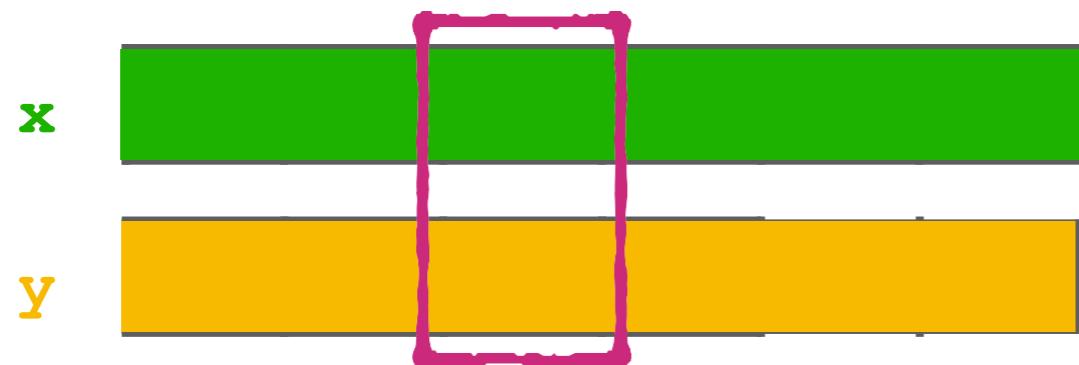
```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array
3 y = np.exp(x)                      # An array
4
5 plt.xlabel('x')                    # Specification
6 plt.ylabel('y')                    # Label
7 plt.show()                         # Label
8
```



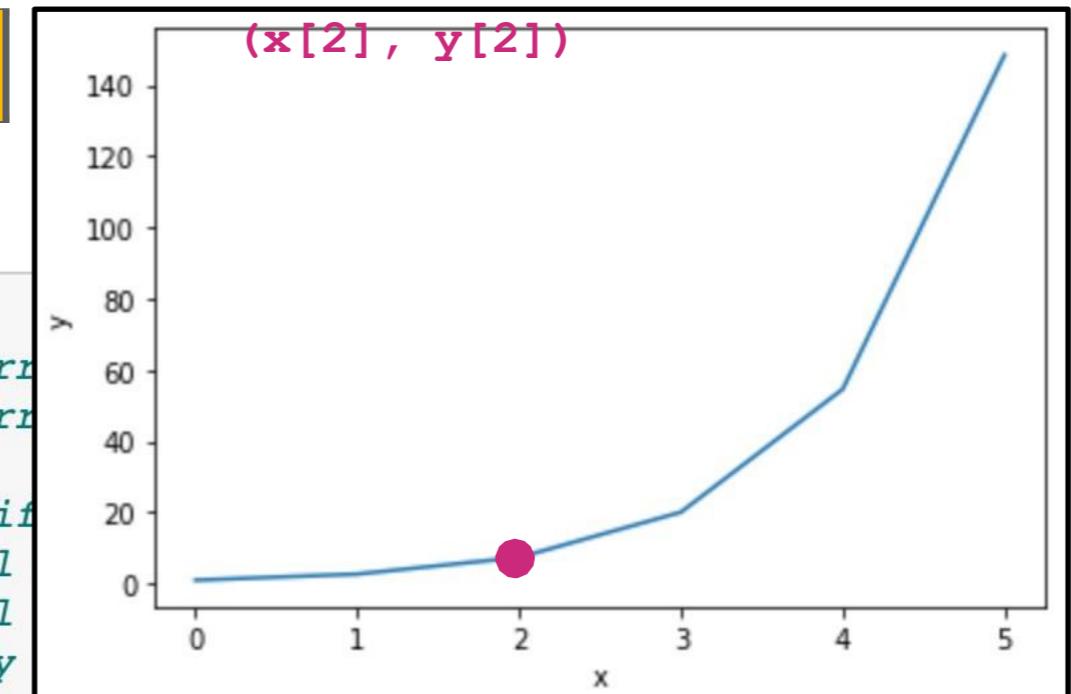
# Data Visualization

- Functions `plot`, `line`, and `bar`

- Line plot



```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array
3 y = np.exp(x)                      # An array
4
5 plt.xlabel('x')                    # Specifying
6 plt.ylabel('y')                    # Label
7 plt.show()                         # Label
8
```



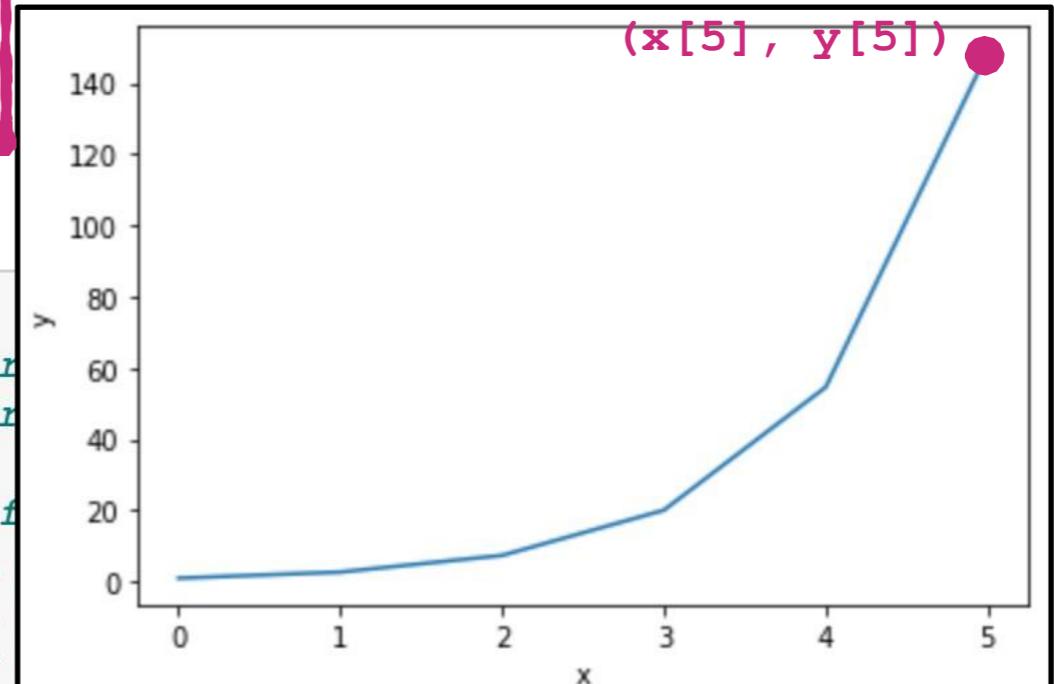
# Data Visualization

- Functions `plot`, `scatter`, and `bar`  
)

- Line plot



```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array
3 y = np.exp(x)                      # An array
4
5 plt.xlabel('x')                    # Specifying
6 plt.ylabel('y')                    # Label
7 plt.show()                         # Label
8
```



# Data Visualization

- Functions `plot`, `show`, and `close`)

- Line plot

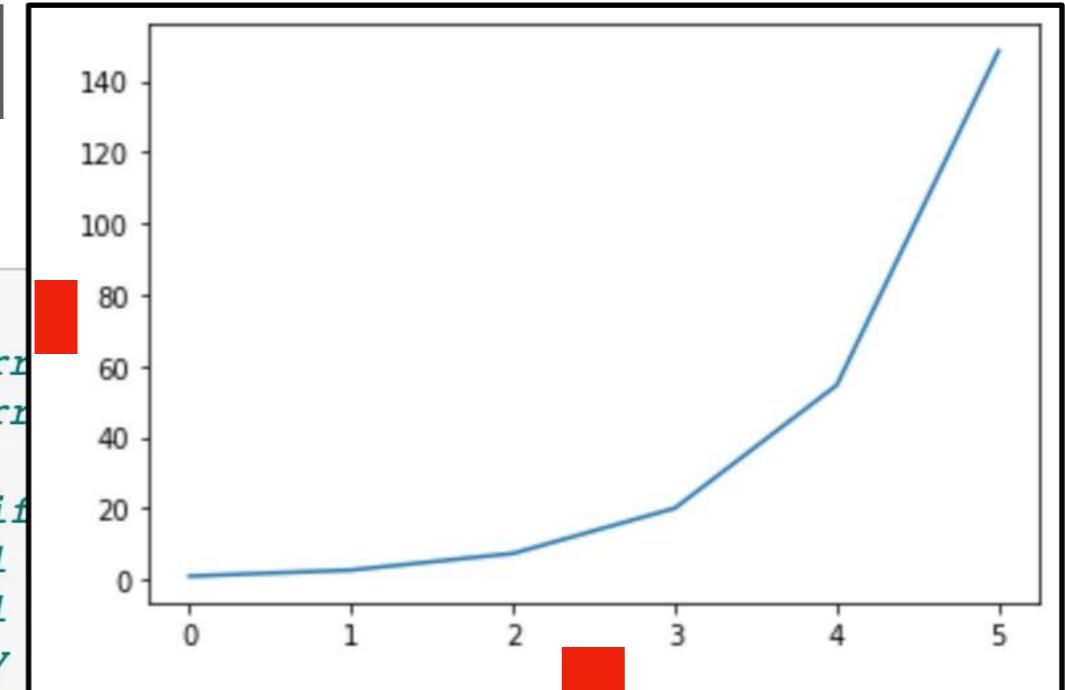
`x`



`y`

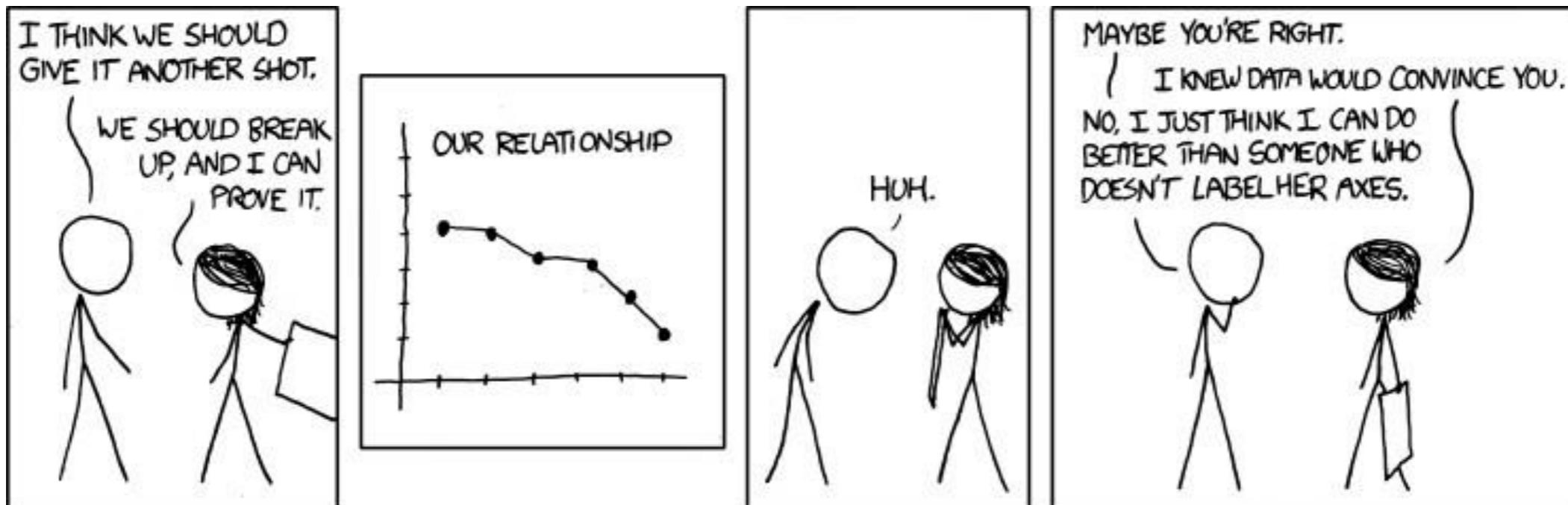


```
1 step = 1
2 x = np.arange(0, 5+step, step)      # An array
3 y = np.exp(x)                      # An array
4
5 plt.plot(x, y)                    # Specifying
6 plt.xlabel('x')                   # Label
7 plt.ylabel('y')                   # Label
8 plt.show()                         # Ready
```



# Data Visualization

- Functions , , and 
- 



# Data Visualization

- Functions `plot`, `line`, and `bar`  
)

- Line plot

`x`

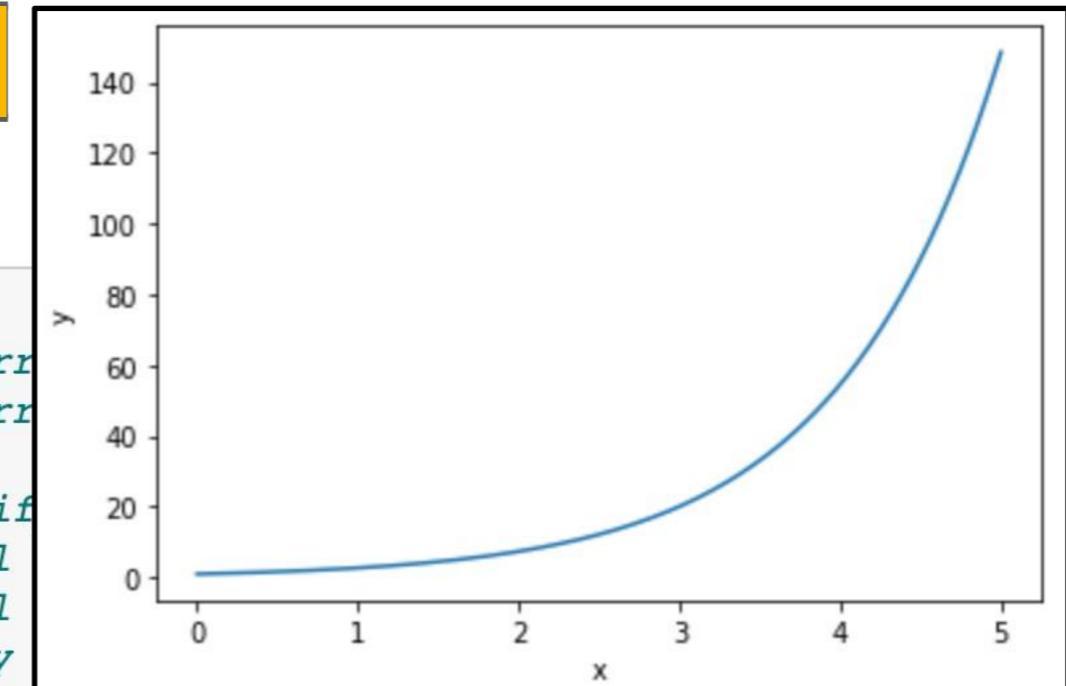


`y`



```
1 # An array
2 # An array
3 # Specifying
4 # Label
5 # Label
6 plt.xlabel('x')
7 plt.ylabel('y')
8 plt.show()
```

```
# An array
# An array
# Specifying
# Label
# Label
# Ready
```



# Data Visualization

- Functions `array`, `array`, and `array`)

- ▶ Scatter plot

`x`

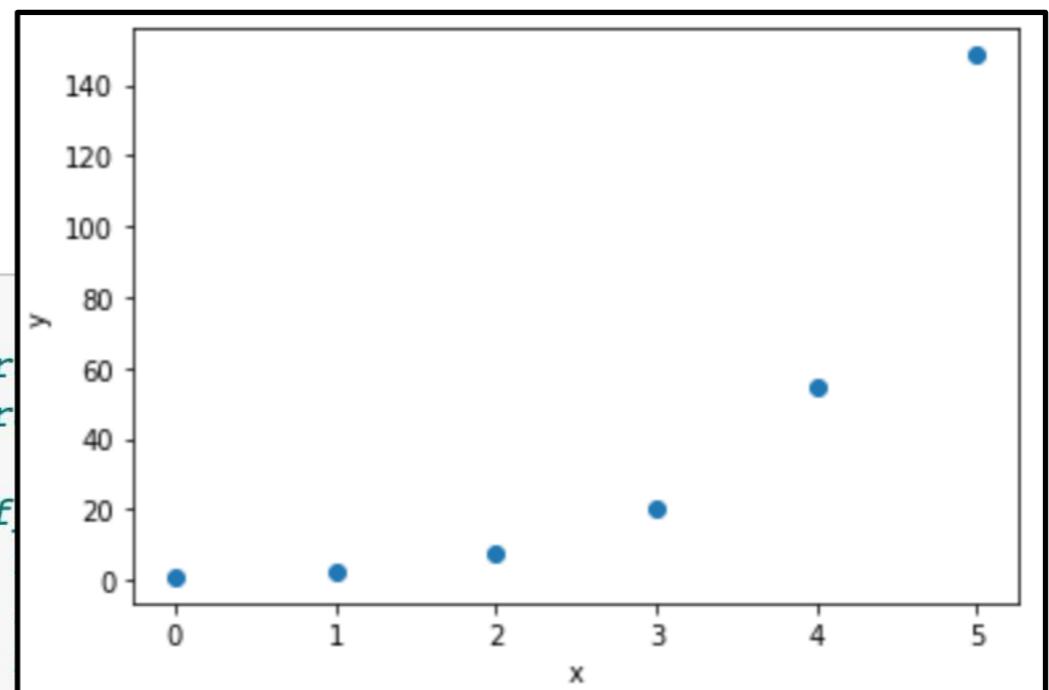


`y`



```
1
2
3
4
5
6 plt.xlabel('x')
7 plt.ylabel('y')
8 plt.show()
```

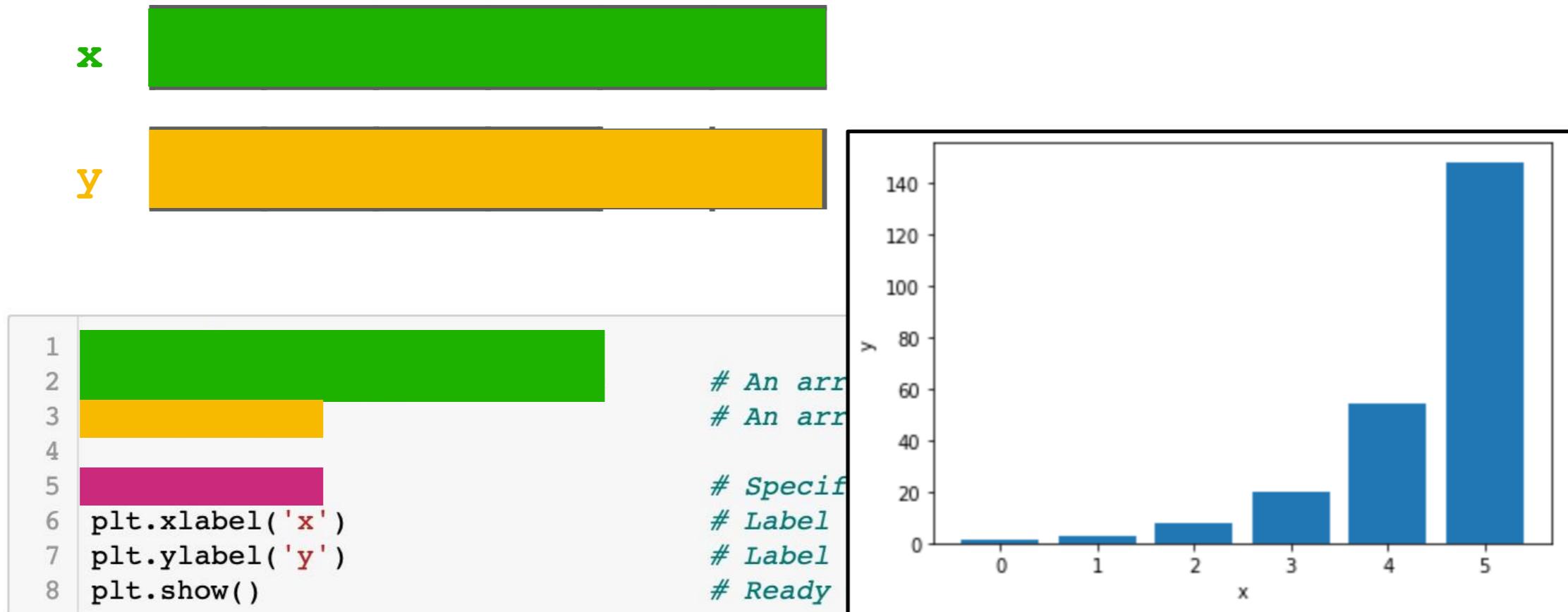
```
# An arr
# An arr
# Specific
# Label
# Label
# Ready
```



# Data Visualization

- Functions `bar()`, `barh()`, and `barc()`

- Bar graph

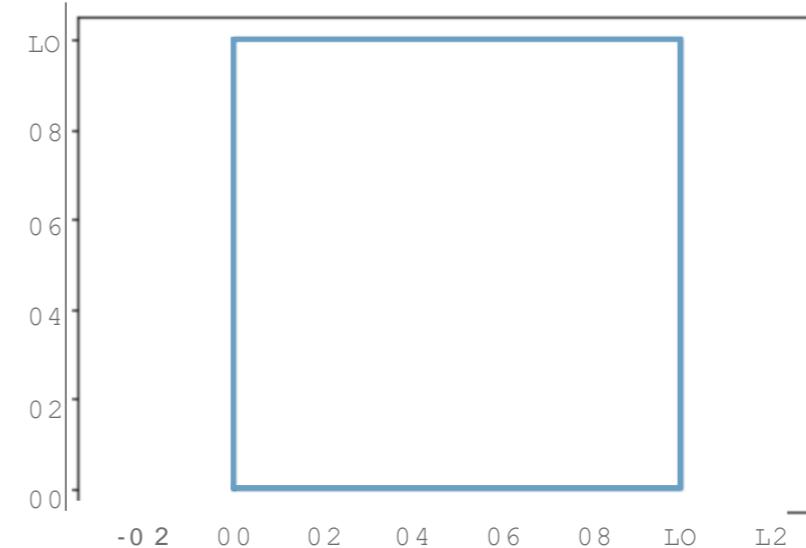


# Data Visualization

## Question

Which of the following code plots a square as the graph shows?

- A. `plt.plot([0, 0, 1, 1, 0], [0, 1, 0, 1, 0])`
- B. `plt.plot([0, 1, 1, 0, 0], [0, 0, 1, 1, 0])`
- C. `plt.plot([0, 1, 1, 0], [0, 0, 1, 1])`
- D. `plt.plot([0, 1, 1, 0], [0, 1, 1, 0])`



# Data Visualization

- Functions , , and

## Example 2

Draw an “elephant” with four parameters.

### essay turning points

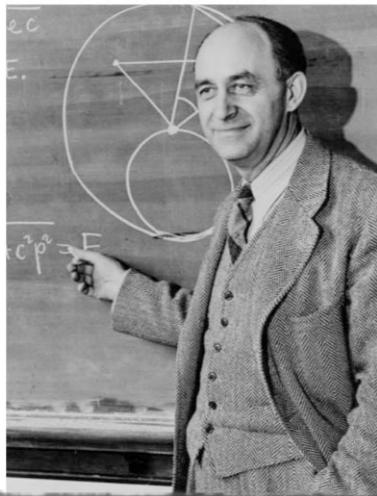
## A meeting with Enrico Fermi

How one intuitive physicist rescued a team from fruitless research.

### Freeman Dyson

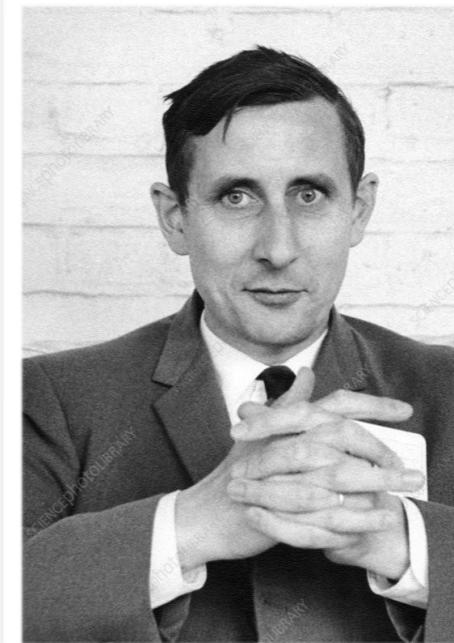
One of the big turning points in my life was a meeting with Enrico Fermi in the spring of 1953. In a few minutes, Fermi politely but ruthlessly demolished a programme of research that my students and I had been pursuing for several years. He probably saved us from several more years of fruitless wandering along a road that was leading nowhere. I am eternally grateful to him for destroying our illusions and telling us the bitter truth.

Fermi was one of the great physicists of our time, outstanding both as a theorist and as an experimenter. He led the team that built the first nuclear reactor in Chicago in 1942. By 1953 he was head of the team that built the Chicago cyclotron, and was using it to explore the strong forces that hold nuclei together. He made the first accurate measurements of the scattering of mesons by protons, an



physical picture, and the forces are so strong that nothing converges. To reach your calculated results, you had to introduce arbitrary cut-off procedures that are not based either on solid physics or on solid mathematics."

In desperation I asked Fermi whether he was not impressed by the agreement between our calculated numbers and his measured numbers. He replied, "How many arbitrary parameters did you use for your calculations?" I thought for a moment about our cut-off procedures and said, "Four." He said, "I remember my friend Johnny von Neumann used to say, with four parameters I can fit an elephant, and with five I can make him wiggle his trunk." With that, the conversation was over. I thanked Fermi for his time and trouble, and sadly took the next bus back to Ithaca to tell the bad news to the students. Because it was important for the students to have their names on a published paper, we did not abandon our calculations immediately. We



Enrico Fermi

Freeman Dyson

# Data Visualization

---

- Functions `ggplot`, `geom_point`, and `stat_smooth`



**“With four parameters I can fit an elephant, and with five I can make him wiggle his trunk”**

*-John von Neumann*

# Data Visualization

---

- Functions  $\text{f1}$ ,  $\text{f2}$ , and  $\text{f3}$   
)

## Drawing an elephant with four complex parameters

American Journal of Physics **78**, 648 (2010); <https://doi.org/10.1119/1.3254017>

Jürgen Mayer

- Max Planck Institute of Molecular Cell Biology and Genetics, Pfotenauerstr. 108, 01307 Dresden, Germany

Khaled Khairy

- European Molecular Biology Laboratory, Meyerhofstraße. 1, 69117 Heidelberg, Germany

Jonathon Howard

- Max Planck Institute of Molecular Cell Biology and Genetics, Pfotenauerstr. 108, 01307 Dresden, Germany

- Functions  $\sin$ ,  $\cos$ , and  $\tan$

## Drawing an elephant with four complex parameters

American Journal of Physics 78, 648 (2010); <https://doi.org/10.1119/1.3254017>

$$\begin{aligned}x &= -30 \sin(t) + 8 \sin(2t) - 10 \sin(3t) - 60 \cos(t) \\ \{y &= -50 \sin(t) - 18 \sin(2t) - 12 \cos(3t) + 14 \cos(5t)\}\end{aligned}$$

where  $t \in [0, 2\pi]$ .

# Data Visualization

- Functions  $\sin$ ,  $\cos$ , and  $\tan$

$$x = -30 \sin(t) + 8 \sin(2t) - 10 \sin(3t) - 60 \cos(t)$$

$$\{y = -50 \sin(t) - 18 \sin(2t) - 12 \cos(3t) + 14 \cos(5t)$$

where  $t \in [0, 2\pi]$ .

```
pi = np.pi
x = -30*np.sin(t) + 8*np.sin(2*t) - 10*np.sin(3*t) - 60*np.cos(t)
y = -50*np.sin(t) - 18*np.sin(2*t) - 12*np.cos(3*t) + 14*np.cos(5*t)

plt.plot(x, y)      # Plot the curve as the elephant
plt.scatter(20, 20)  # Plot the eye of the elepant as a scatter point

plt.show()
```

An array running from 0 to  $2\pi$

# Data Visualization

- Functions  $\sin$ ,  $\cos$ , and  $\tan$

$$\{y = -50 \sin(t) - 18 \sin(2t) - 12 \cos(3t) + 14 \cos(5t)$$

where  $t \in [0, 2\pi]$ .

```
pi = np.pi
t = np.arange(0, 2*pi+0.01, 0.01)

y = - 50*np.sin(t) - 18*np.sin(2*t) - 12*np.cos(3*t) + 14*np.cos(5*t)

plt.plot(x, y)          # Plot the curve as the elephant
plt.scatter(20, 20)      # Plot the eye of the elepant as a scatter point

plt.show()
```

# Data Visualization

- Functions  $\sin$ ,  $\cos$ , and  $\tan$

)

$$x = -30 \sin(t) + 8 \sin(2t) - 10 \sin(3t) - 60 \cos(t)$$

{  
y

where  $t \in [0, 2\pi]$ .

```
pi = np.pi
t = np.arange(0, 2*pi+0.01, 0.01)

x = -30*np.sin(t) + 8*np.sin(2*t) - 10*np.sin(3*t) - 60*np.cos(t)

plt.plot(x, y)      # Plot the curve as the elephant
plt.scatter(20, 20) # Plot the eye of the elepant as a scatter point

plt.show()
```

# Data Visualization

- Functions  $\sin$ ,  $\cos$ , and  $\tan$

$$x = -30 \sin(t) + 8 \sin(2t) - 10 \sin(3t) - 60 \cos(t)$$

$$\{y = -50 \sin(t) - 18 \sin(2t) - 12 \cos(3t) + 14 \cos(5t)$$

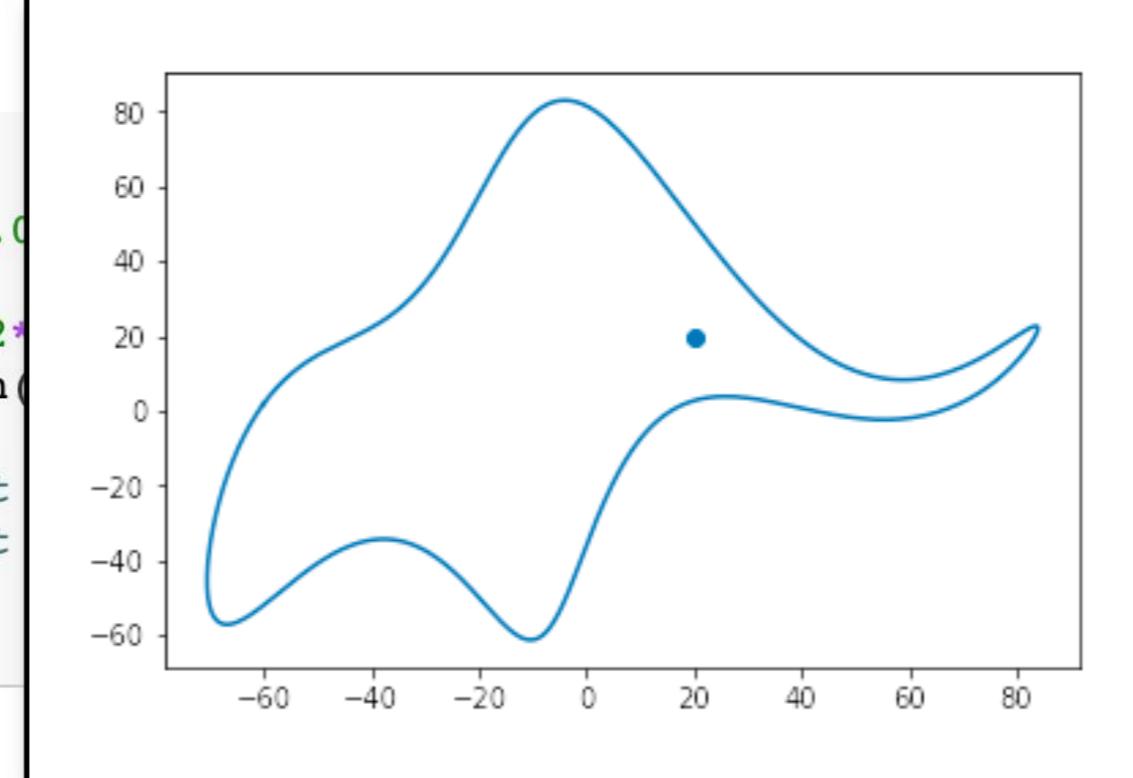
where  $t \in [0, 2\pi]$ .

```
pi = np.pi
t = np.arange(0, 2*pi+0.01, 0.01)

x = -30*np.sin(t) + 8*np.sin(2*t)
y = -50*np.sin(t) - 18*np.sin(3*t)
z = 14*np.cos(5*t)

# Plot
# Plot

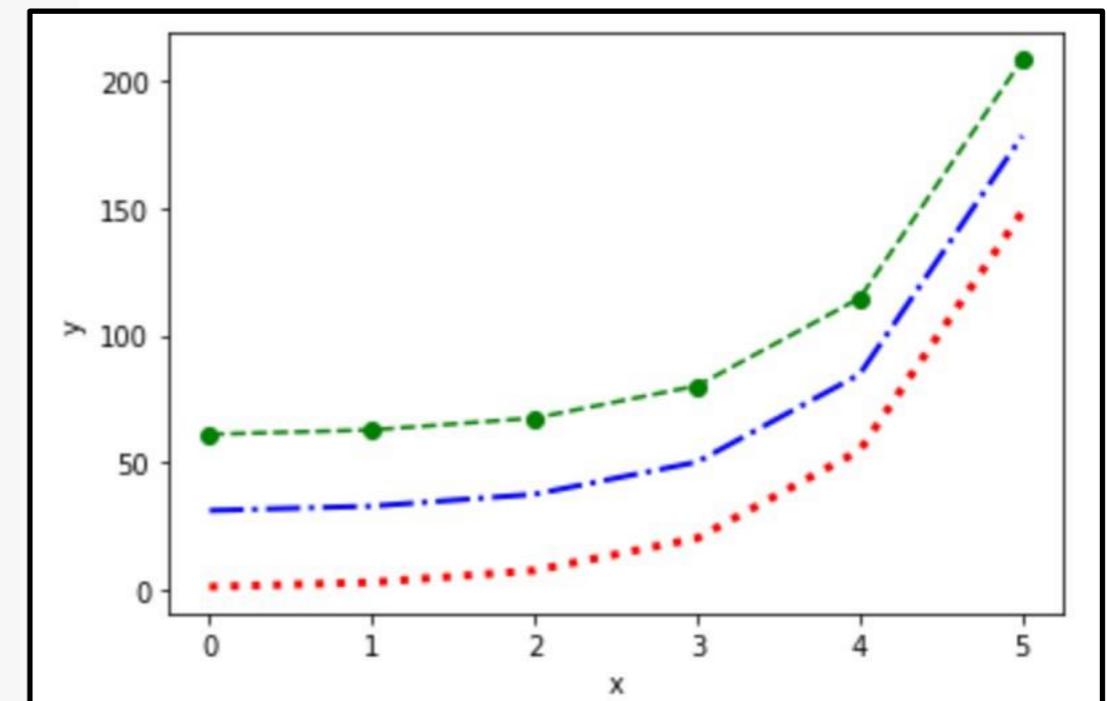
plt.show()
```



# Data Visualization

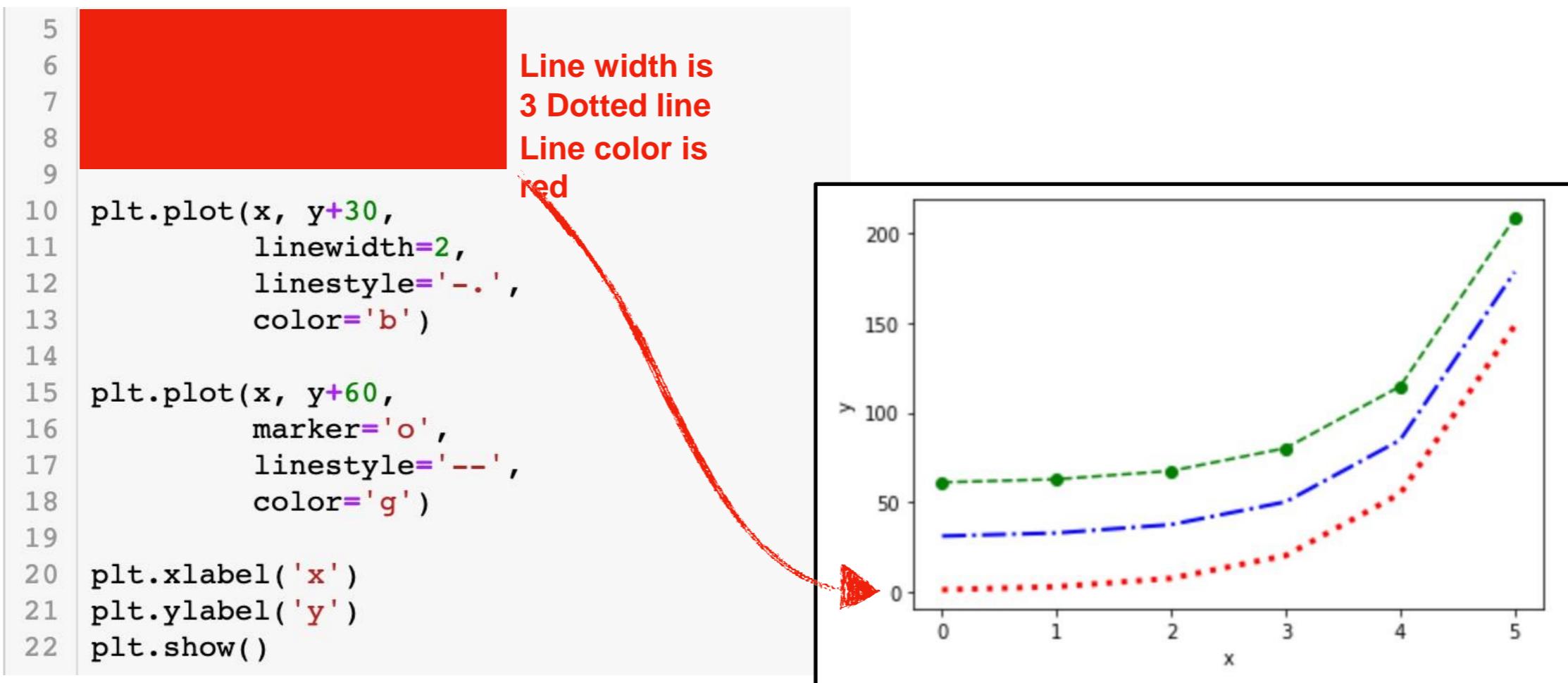
- Configuration of plots and figures
  - Specify the line plot by keyword arguments

```
5 plt.plot(x, y,
6             linewidth=3,
7             linestyle=':',
8             color='r')
9
10 plt.plot(x, y+30,
11             linewidth=2,
12             linestyle='-.',
13             color='b')
14
15 plt.plot(x, y+60,
16             marker='o',
17             linestyle='--',
18             color='g')
19
20 plt.xlabel('x')
21 plt.ylabel('y')
22 plt.show()
```



# Data Visualization

- Configuration of plots and figures
  - Specify the line plot by keyword arguments

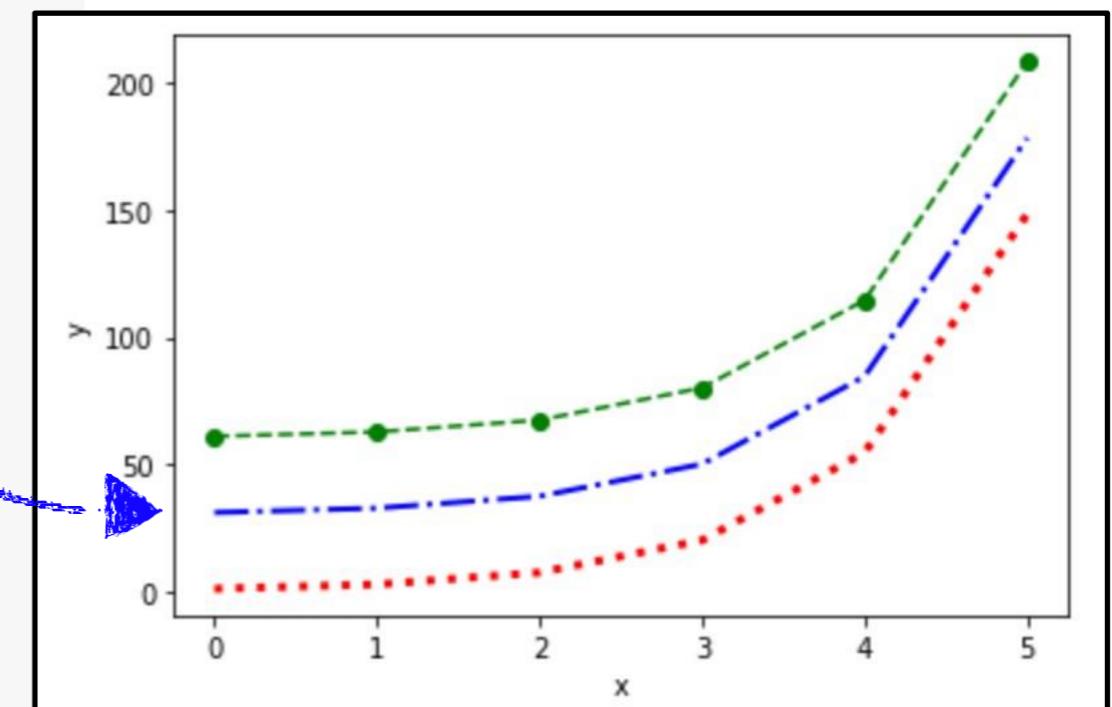


# Data Visualization

- Configuration of plots and figures
  - Specify the line plot by keyword arguments

```
5 plt.plot(x, y,
6             linewidth=3,
7             linestyle=':',
8             color='r')
9
10
11
12
13
14
15 plt.plot(x, y+60,
16             marker='o',
17             linestyle='--',
18             color='g')
19
20 plt.xlabel('x')
21 plt.ylabel('y')
22 plt.show()
```

Line width is 2  
Dash-dotted line  
Line color is  
blue

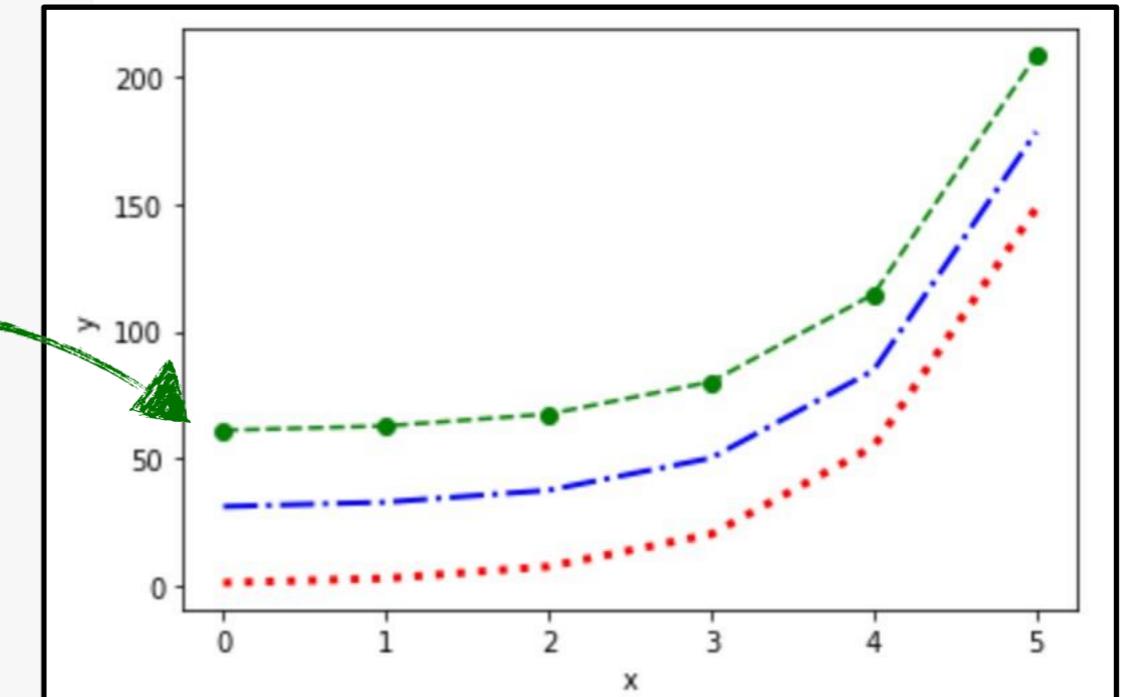


# Data Visualization

- Configuration of plots and figures
  - Specify the line plot by keyword arguments

```
5 plt.plot(x, y,
6             linewidth=3,
7             linestyle=':',
8             color='r')
9
10 plt.plot(x, y+30,
11             linewidth=2,
12             linestyle='-.',
13             color='b')
14
15
16
17
18
19
20 plt.xlabel('x')
21 plt.ylabel('y')
22 plt.show()
```

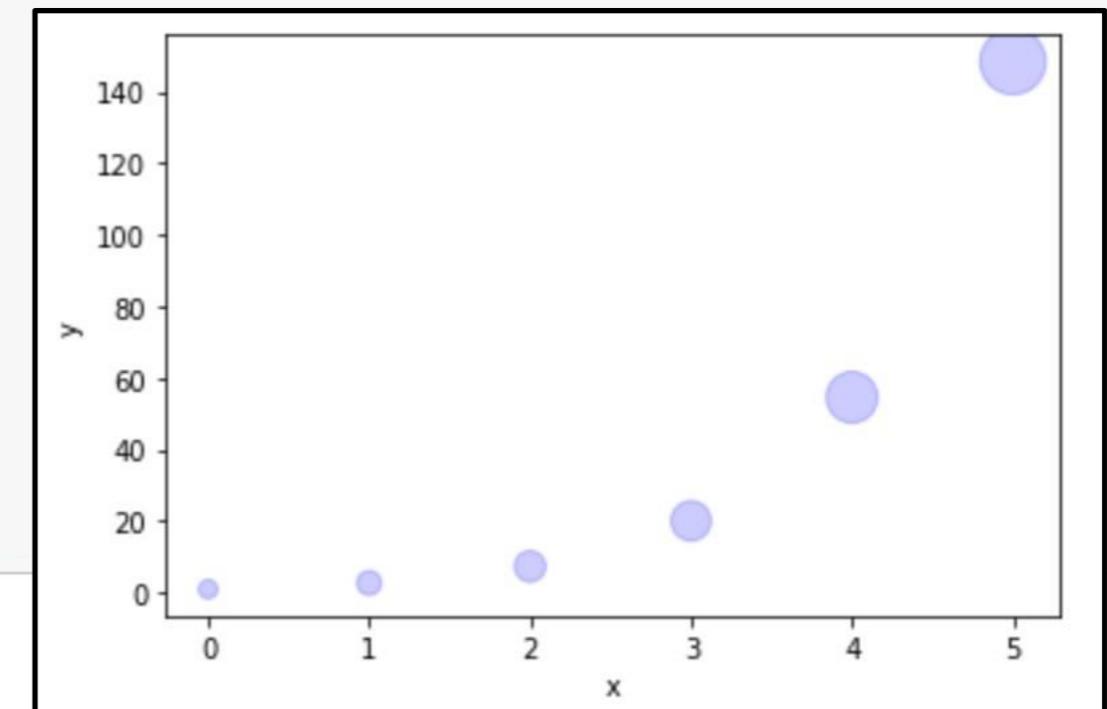
Circle marker  
Dashed line  
Line color is  
green



# Data Visualization

- Configuration of plots and figures
  - Specify the scatter plot by keyword arguments

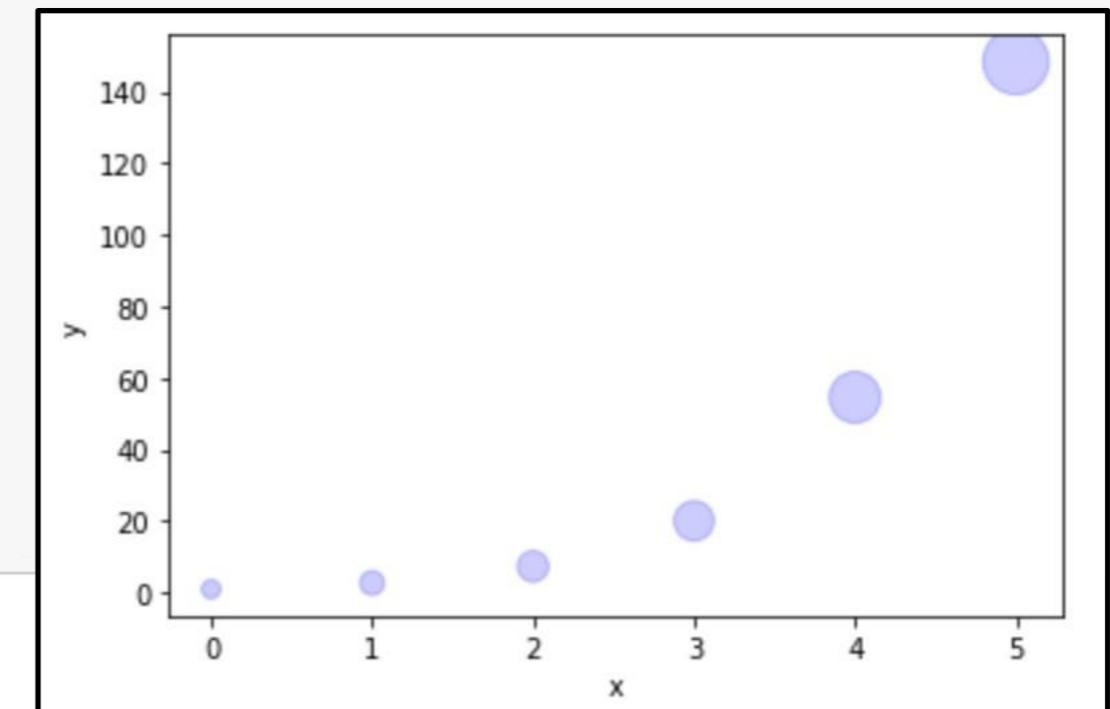
```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4 z = np.sqrt(y)
5
6 plt.scatter(x, y,
7             s=z*50,  Size of scatter dots
8             c='b',
9             alpha=0.2)
10 plt.xlabel('x')
11 plt.ylabel('y')
12 plt.show()
```



# Data Visualization

- Configuration of plots and figures
  - Specify the scatter plot by keyword arguments

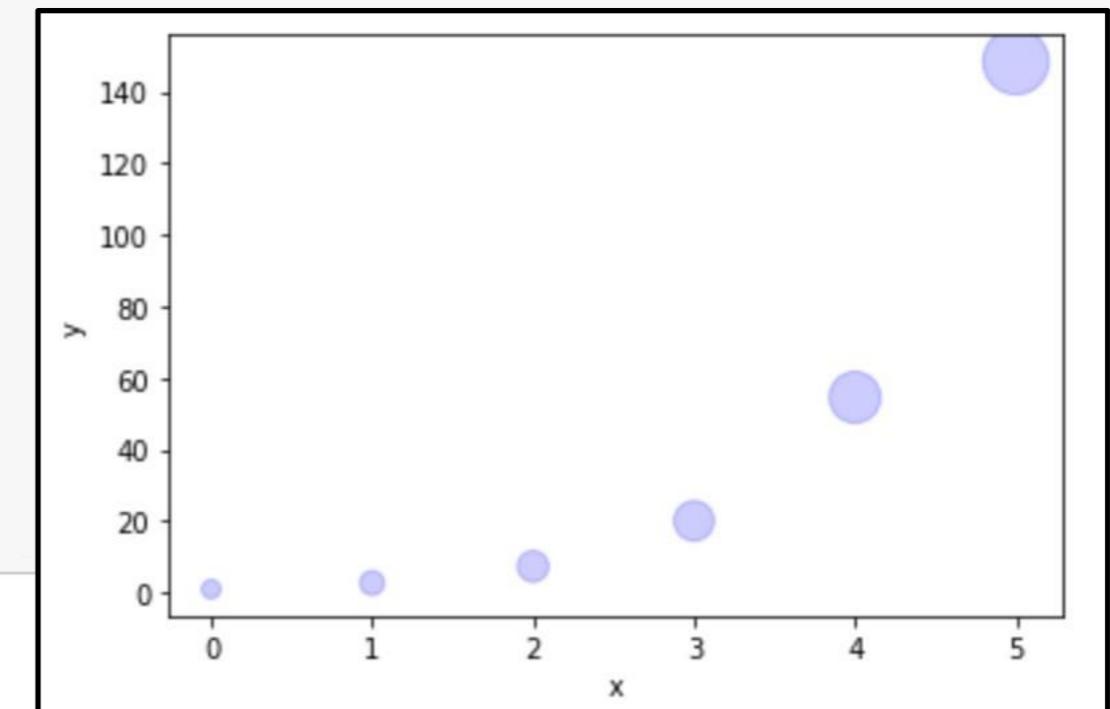
```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4 z = np.sqrt(y)
5
6 plt.scatter(x, y,
7             s=z*50,
8             c='blue', Color of scatter dots
9             alpha=0.2)
10 plt.xlabel('x')
11 plt.ylabel('y')
12 plt.show()
```



# Data Visualization

- Configuration of plots and figures
  - Specify the scatter plot by keyword arguments

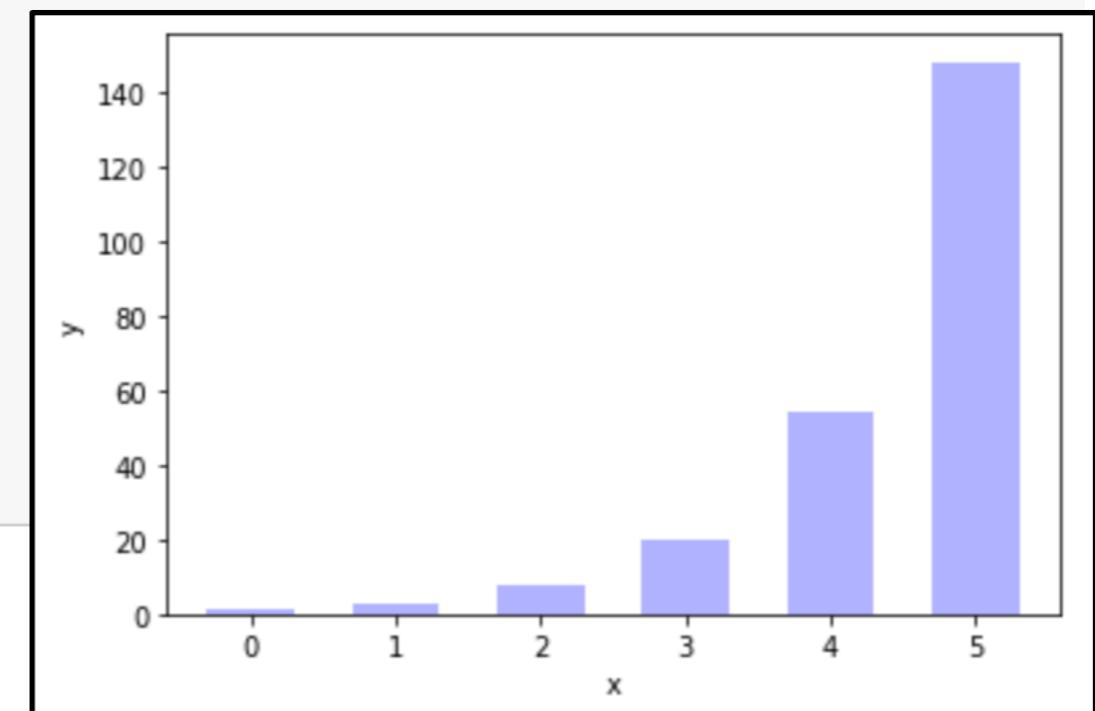
```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4 z = np.sqrt(y)
5
6 plt.scatter(x, y,
7             s=z*50,
8             c='b',
9             alpha=0.5) # Opacity of dots
10 plt.xlabel('x')
11 plt.ylabel('y')
12 plt.show()
```



# Data Visualization

- Configuration of plots and figures
  - Specify the bar graph by keyword arguments

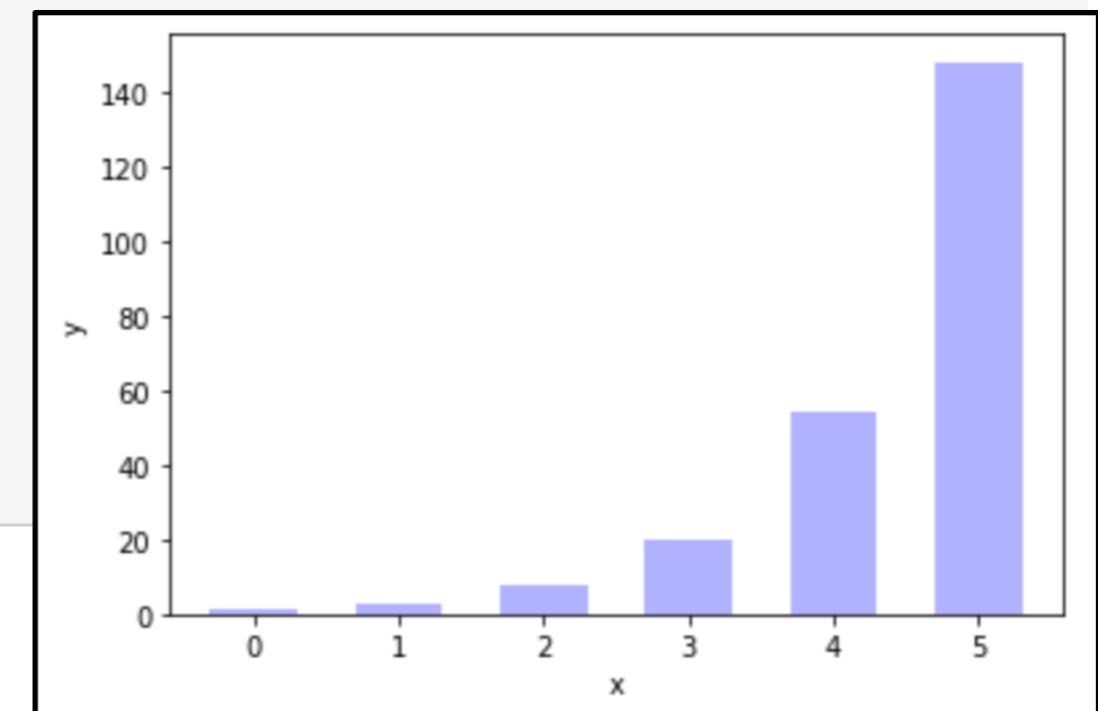
```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4
5 plt.bar(x, y,
6          Bar width is 0.6
7          color='b',
8          alpha=0.3)
9 plt.xlabel('x')
10 plt.ylabel('y')
11 plt.show()
```



# Data Visualization

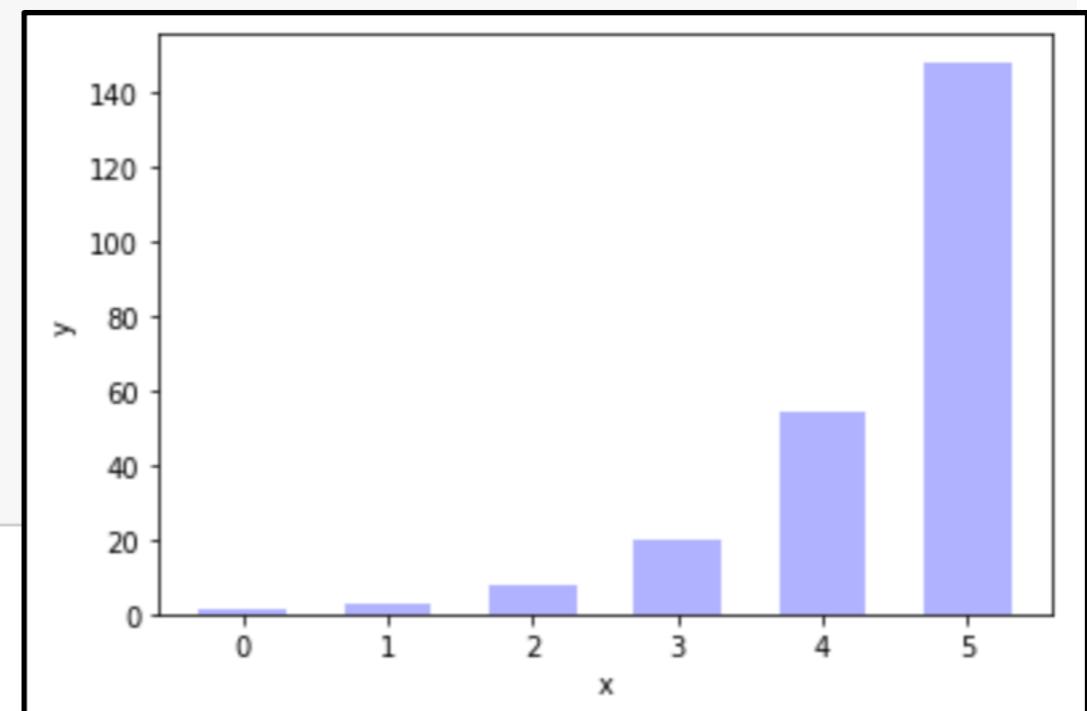
- Configuration of plots and figures
  - Specify the bar graph by keyword arguments

```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4
5 plt.bar(x, y,
6          width=0.6,
7          color=Color of bar graph
8          alpha=0.3)
9 plt.xlabel('x')
10 plt.ylabel('y')
11 plt.show()
```



- Configuration of plots and figures
  - Specify the bar graph by keyword arguments

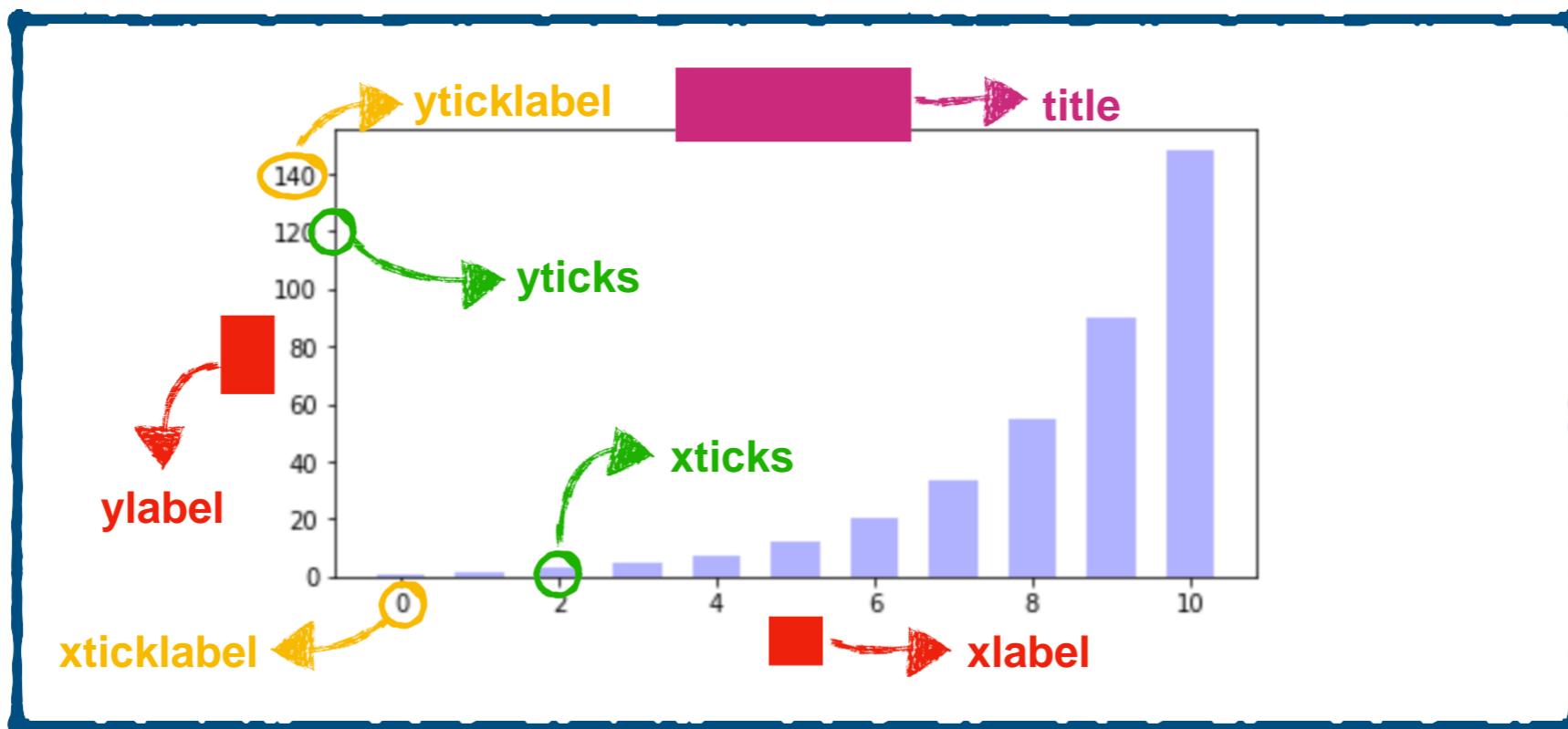
```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4
5 plt.bar(x, y,
6          width=0.6,
7          color='b',
8          alpha=0.5) Opacity of bar graph
9 plt.xlabel('x')
10 plt.ylabel('y')
11 plt.show()
```



# Data Visualization

- Configuration of plots and figures
  - ▶ Components of a figure

Figure

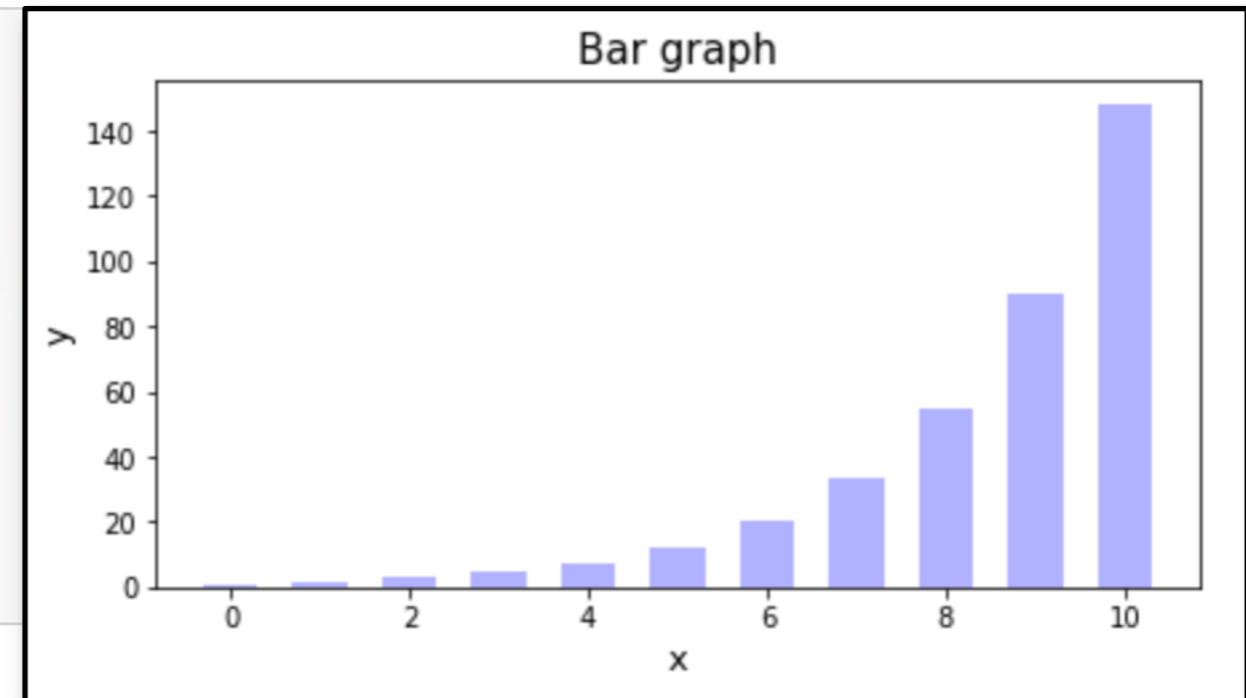


# Data Visualization

- Configuration of plots and figures

- Attributes of a figure

```
1 step = 1
2 x = np.arange(0, 10+step, step)
3 y = np.exp(x/2)
4
5 plt.figure(figsize=(7, 3.5))
6 plt.title('Bar graph', fontsize=15)
7 plt.bar(x, y,
8         width=0.6,
9         color='b',
10        alpha=0.3)
11 plt.xlabel('x', fontsize=13)
12 plt.ylabel('y', fontsize=13)
13 plt.show()
```

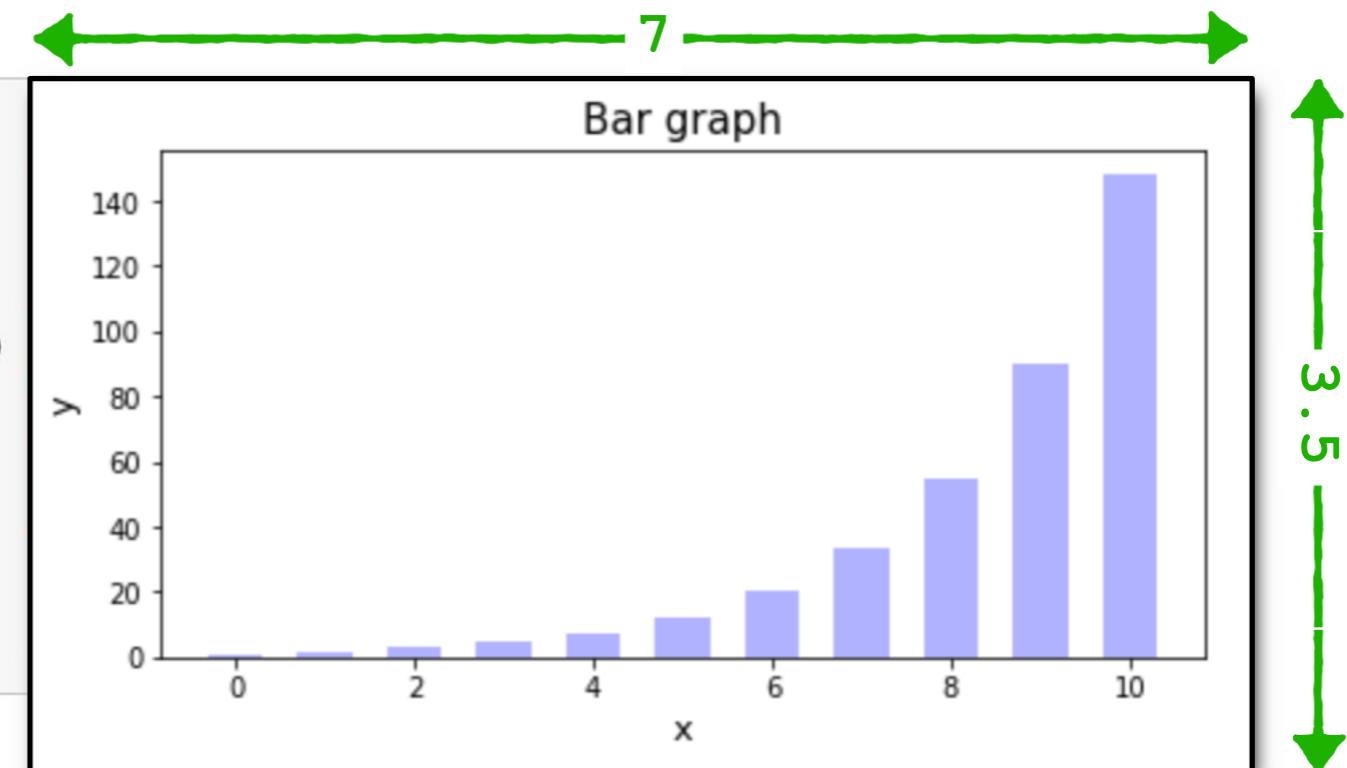


# Data Visualization

- Configuration of plots and figures

- Attributes of a figure

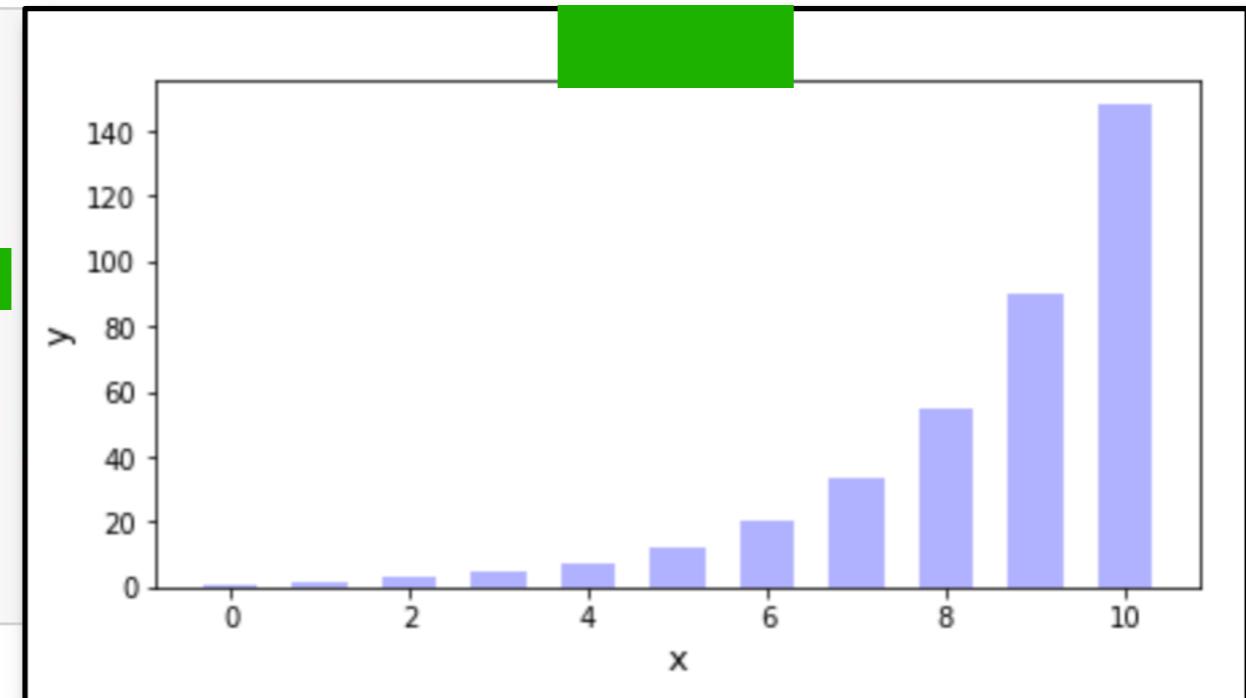
```
1 step = 1
2 x = np.arange(0, 10+step, step)
3 y = np.exp(x/2)
4
5
6 plt.title('Bar graph', fontsize=15)
7 plt.bar(x, y,
8         width=0.6,
9         color='b',
10        alpha=0.3)
11 plt.xlabel('x', fontsize=13)
12 plt.ylabel('y', fontsize=13)
13 plt.show()
```



# Data Visualization

- Configuration of plots and figures
  - Attributes of a figure

```
1 step = 1
2 x = np.arange(0, 10+step, step)
3 y = np.exp(x/2)
4
5 plt.figure(figsize=(7, 3.5))
6
7 plt.bar(x, y,
8         width=0.6,
9         color='b',
10        alpha=0.3)
11 plt.xlabel('x', fontsize=13)
12 plt.ylabel('y', fontsize=13)
13 plt.show()
```

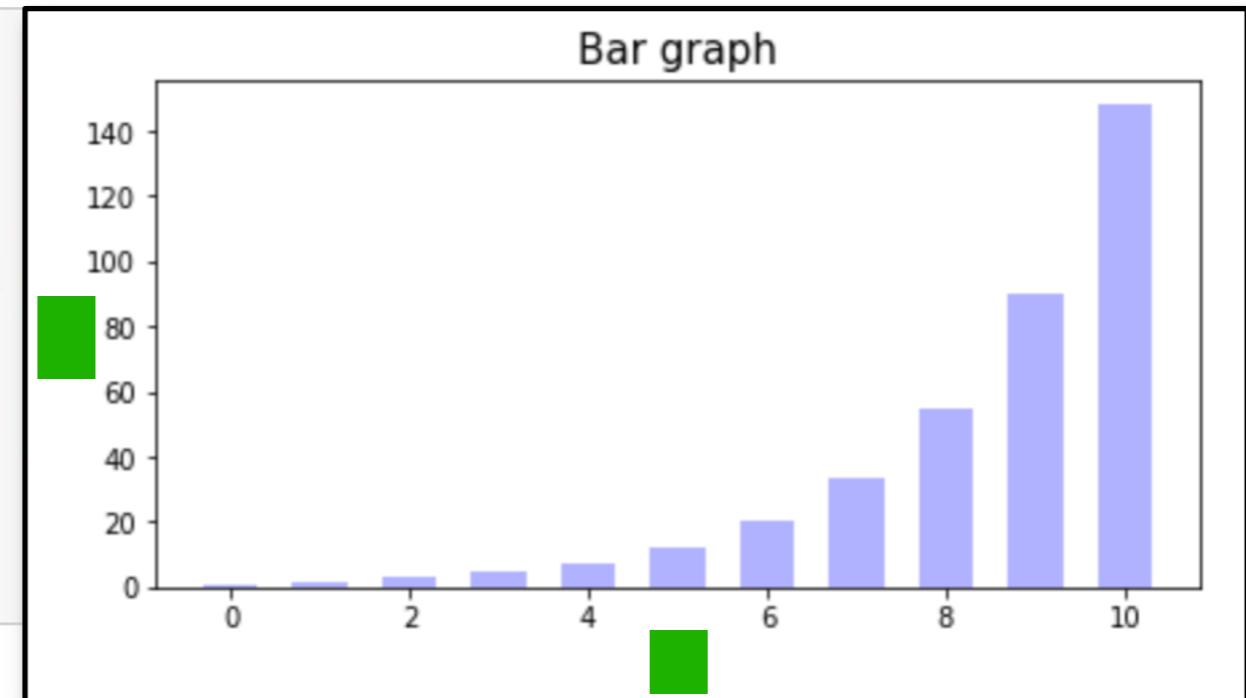


# Data Visualization

- Configuration of plots and figures

- Attributes of a figure

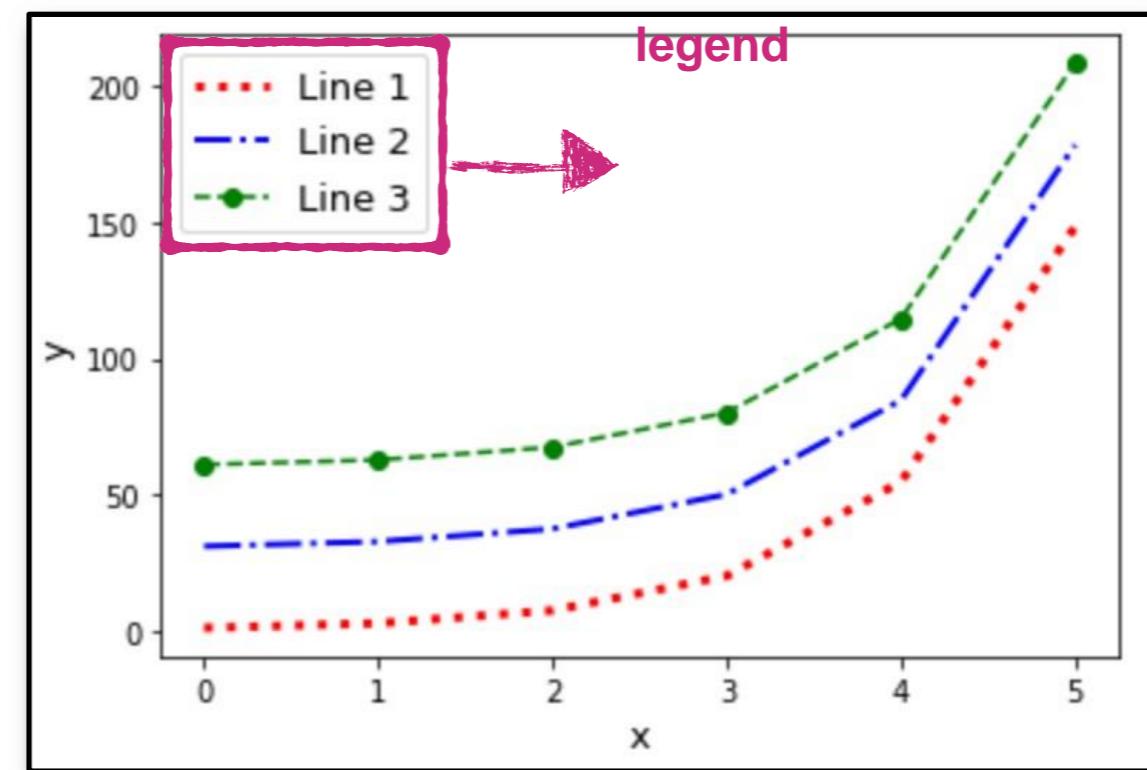
```
1 step = 1
2 x = np.arange(0, 10+step, step)
3 y = np.exp(x/2)
4
5 plt.figure(figsize=(7, 3.5))
6 plt.title('Bar graph', fontsize=15)
7 plt.bar(x, y,
8         width=0.6,
9         color='b',
10        alpha=0.3)
11
12 plt.show()
```



# Data Visualization

- Labeling plots and the function

```
    )
```

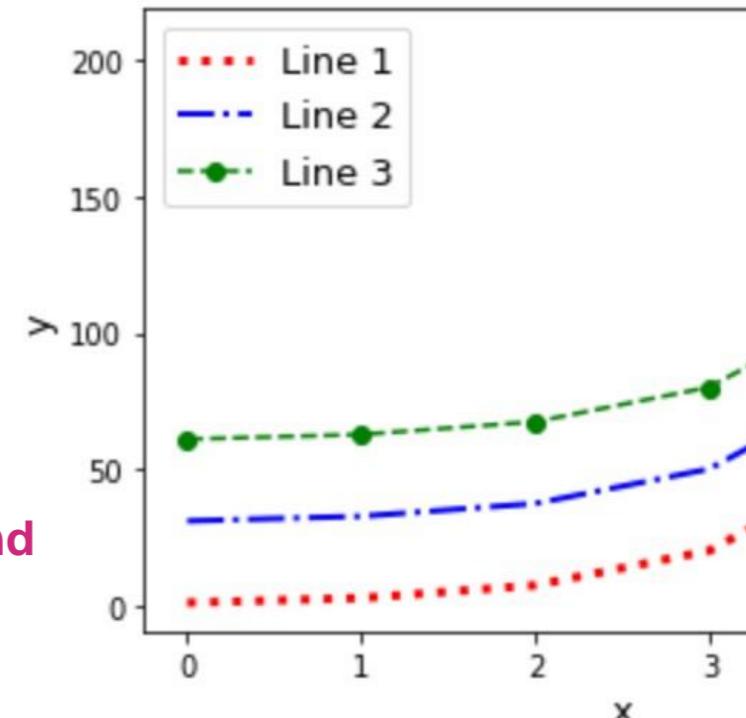


# Data Visualization

- Labeling plots and the function  
    )- Connection between plots and labels

```
1 step = 1
2 x = np.arange(0, 5+step, step)
3 y = np.exp(x)
4
5 plt.plot(x, y, linewidth=3, linestyle=':', color='r',
6           label='Line 1')
7
8 plt.plot(x, y+30, linewidth=2, linestyle='-.', color='b',
9           label='Line 2')
10
11 plt.plot(x, y+60, marker='o', linestyle='--', color='g',
12           label='Line 3')
13
14 plt.legend(fontsize=13)
15 plt.xlabel('x', fontsize=13)
16 plt.ylabel('y', fontsize=13)
17 plt.show()
```

legend



# Data Visualization

---

## Example 5

The discrete distribution information of two types of newspapers is stored in the dictionary `dist_dict`. Use a bar graph to show the demands of these two types of newspapers under different weather conditions.

```
{ 'weather': ['Sunny', 'Cloudy', 'Raining', 'Haze'],
  'Thunderstorm',
  'probs': [0.315, 0.226, 0.289, 0.087, 0.083],
  'paper1': [560, 530, 389, 202, 278],
  'paper2': [533, 486, 386, 234, 263]}
```

# Data Visualization

```
1 Array 0, 1, 2, ..., 4
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.6,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.6,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```

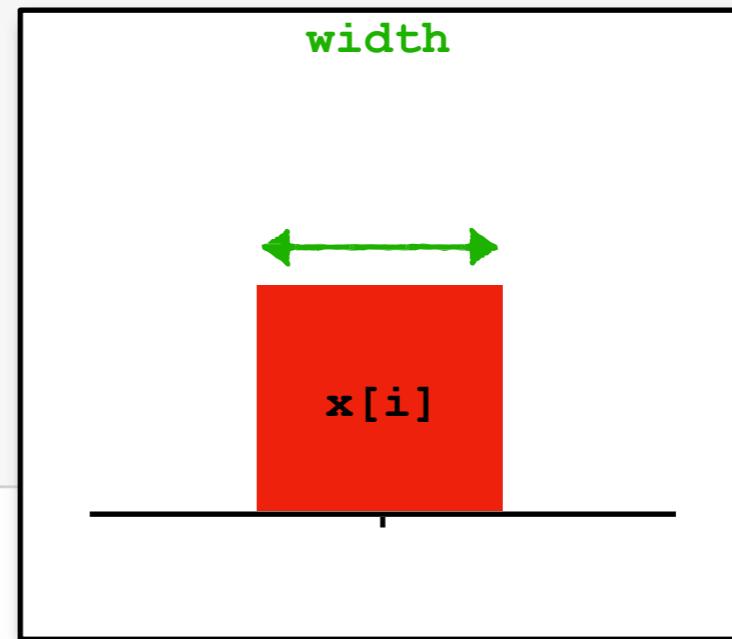
# Data Visualization

---

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4)) Configure the size of figure
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.6,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.6,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```

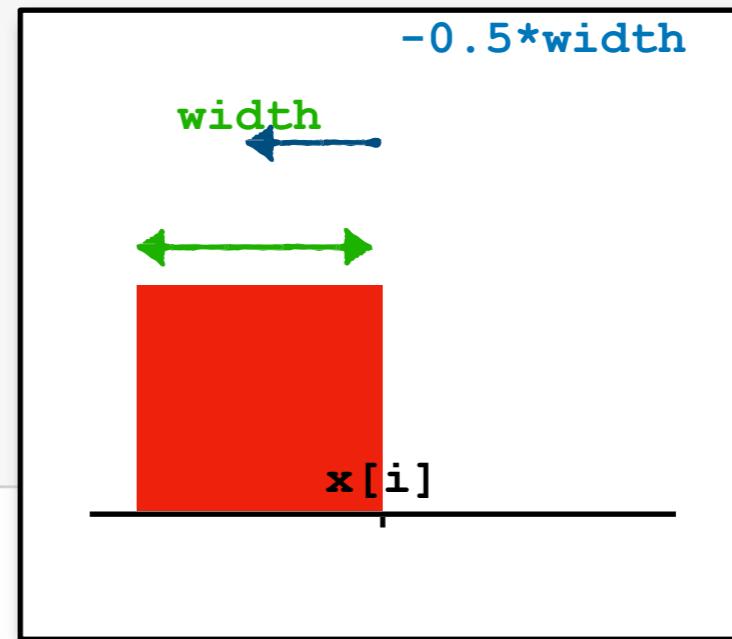
# Data Visualization

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.6,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.6,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```



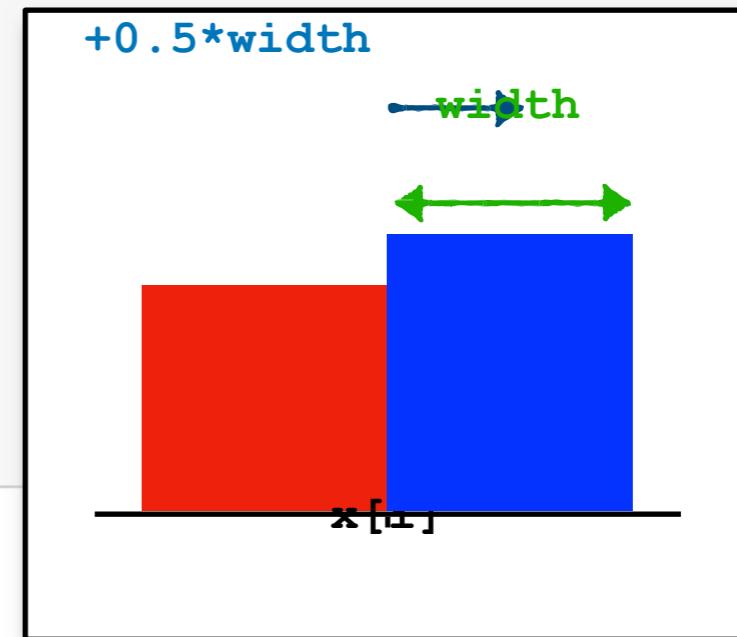
# Data Visualization

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(-0.5*width, paper1,
6         color='r', width=width, alpha=0.6,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.6,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```



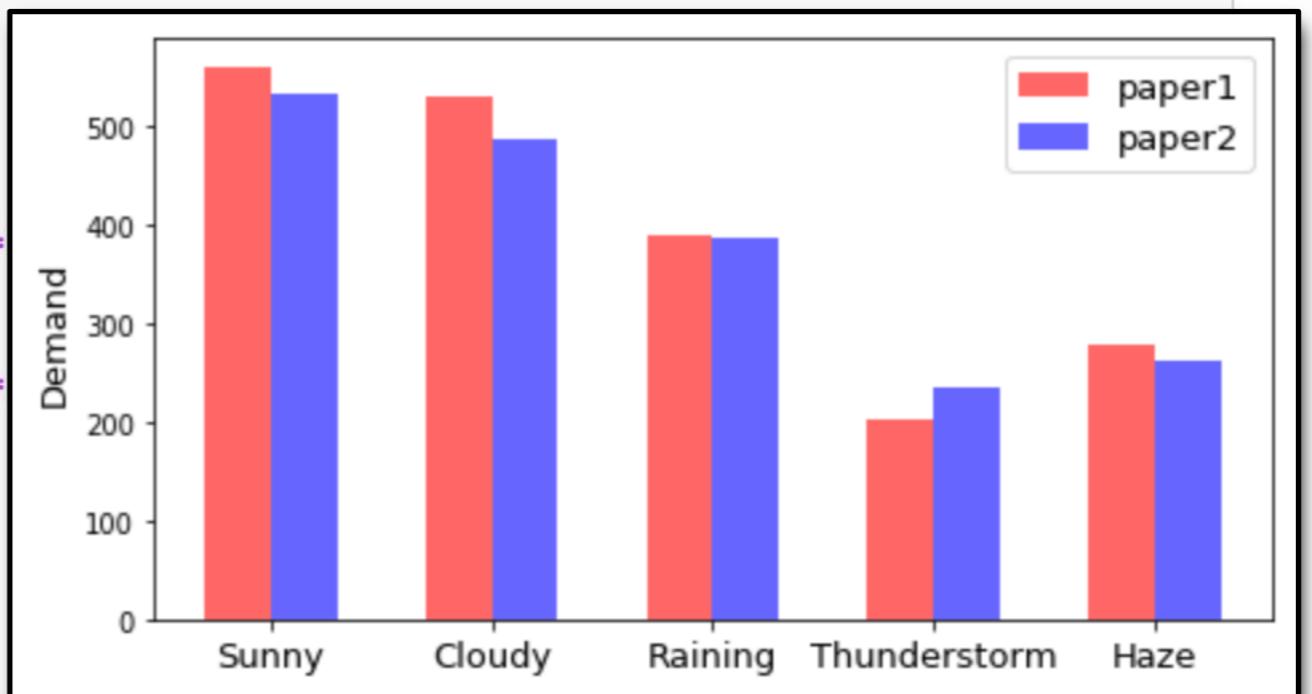
# Data Visualization

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.6,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.6,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```



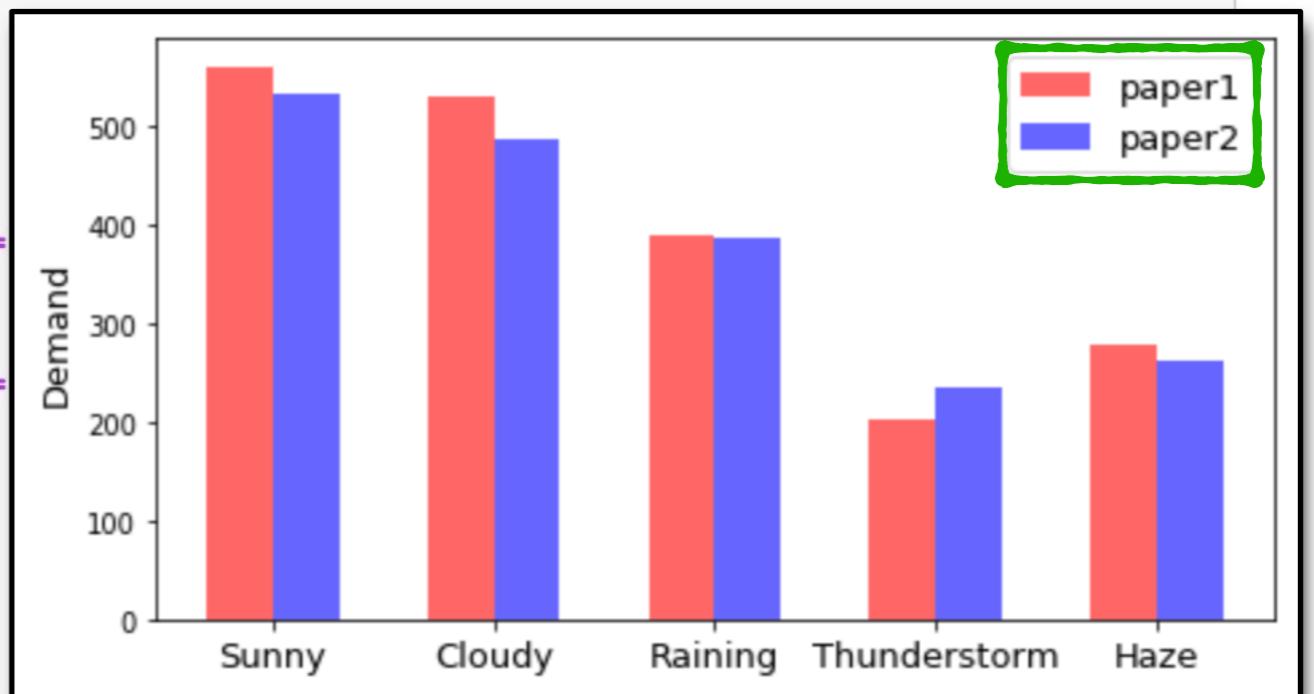
# Data Visualization

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.8,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.8,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```



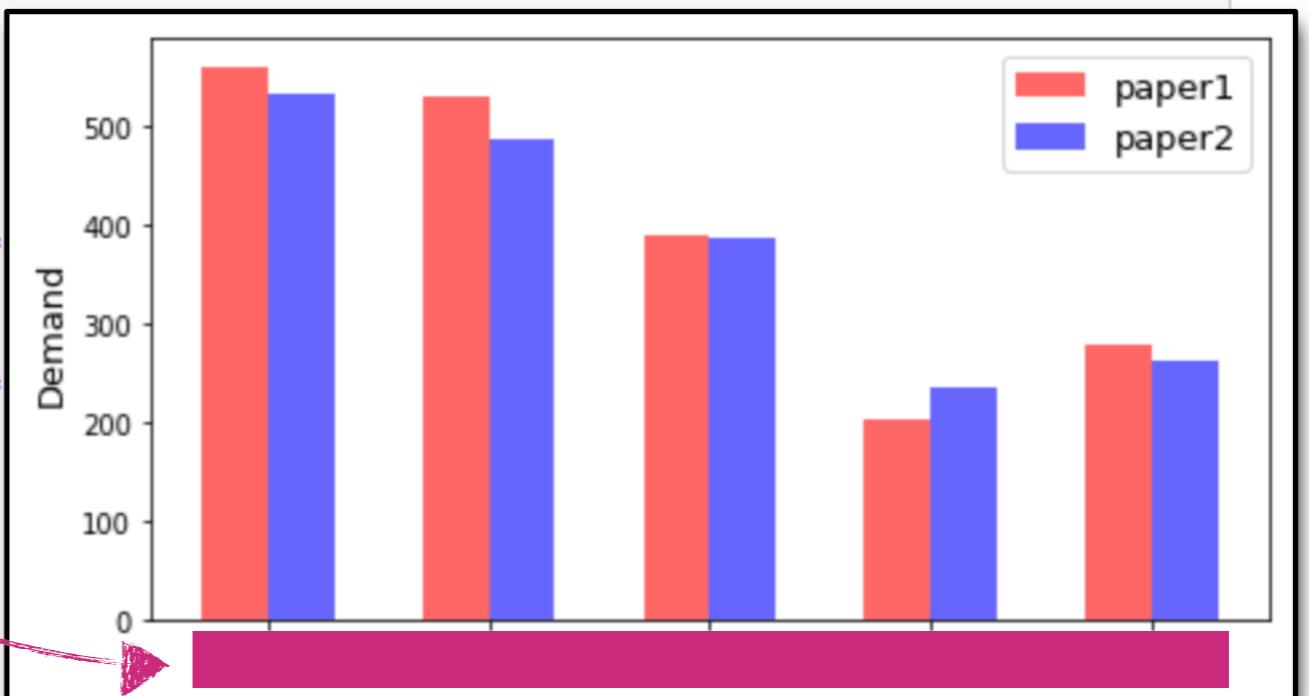
# Data Visualization

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.8,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.8,
10        label='paper2')
11
12
13 plt.xticks(x, weather, fontsize=13)
14 plt.ylabel('Demand', fontsize=13)
15 plt.show()
```



# Data Visualization

```
1 x = np.arange(len(weather))
2 width = 0.3
3
4 plt.figure(figsize=(7.5, 4))
5 plt.bar(x-0.5*width, paper1,
6         color='r', width=width, alpha=0.8,
7         label='paper1')
8 plt.bar(x+0.5*width, paper2,
9         color='b', width=width, alpha=0.8,
10        label='paper2')
11
12 plt.legend(fontsize=13)
13 plt.ylabel('Demand', fontsize=13)
14 plt.show()
```



Replace integer  
x-ticks by strings

# Python Packages

---

- Packages NumPy and Matplotlib



Joan Watson: “How did you know he had an affair?”

—*Elementary, Season 1, Episode 1*

# Python Packages

---

- Packages NumPy and Matplotlib



Joan Watson: “How did you know he had an affair?”

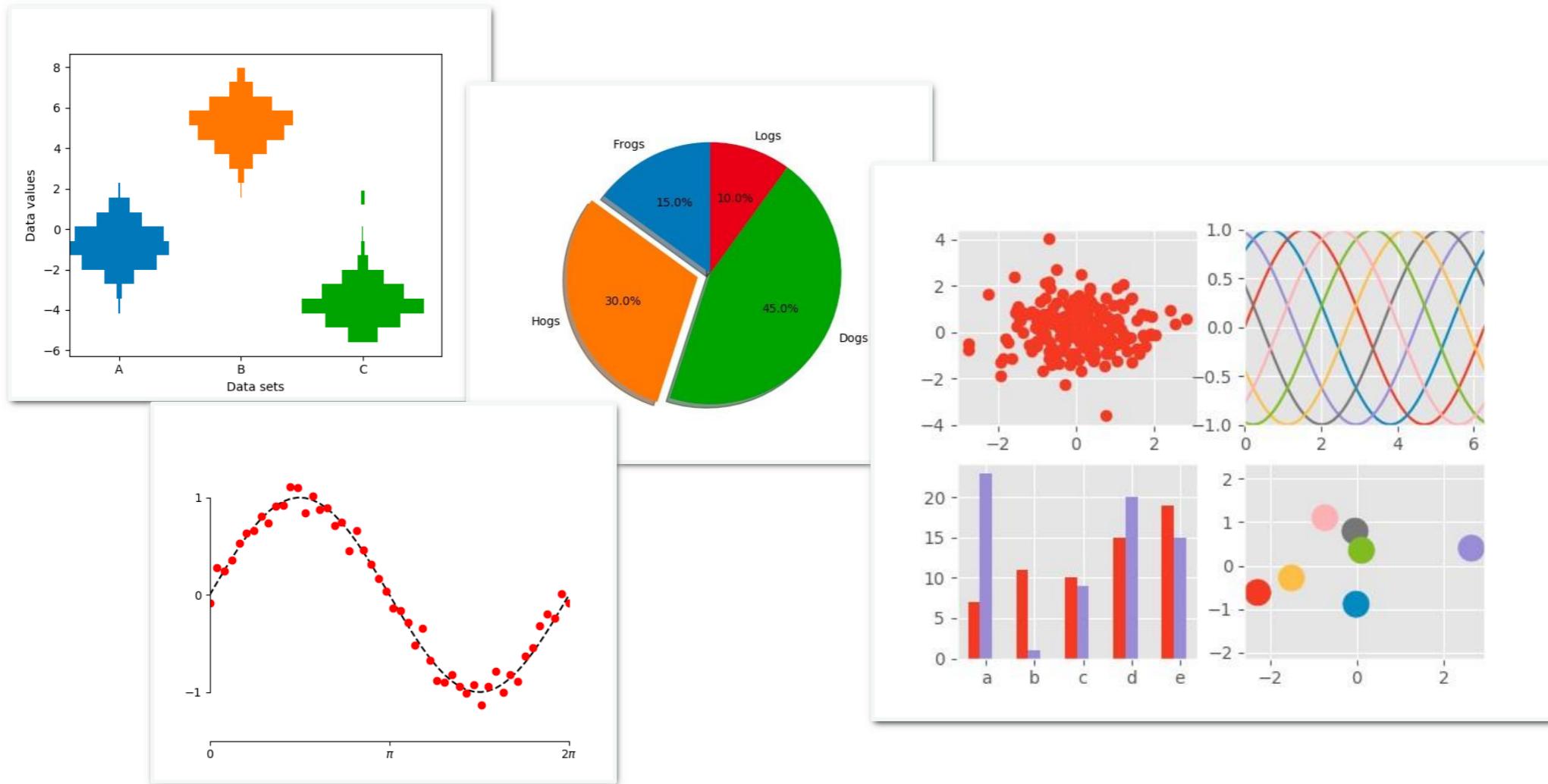
Sherlock Holmes: “Google. Well, not everything is deducible.”

—*Elementary, Season 1, Episode 1*

# Python Packages

- Packages NumPy and Matplotlib

- [Matplotlib thumbnail gallery](#): examples and demos



# Roadmap

---

- 1 Array Data Structure
- 2 Data Visualization
- 3 **Pandas**
- 4 Descriptive Analytics in Python

“It is a capital mistake to theorize before one has data. Insensibly one begins to twist facts to suit theories, instead of theories to suit facts.”

—*Sherlock Holmes, A Scandal in Bohemia*



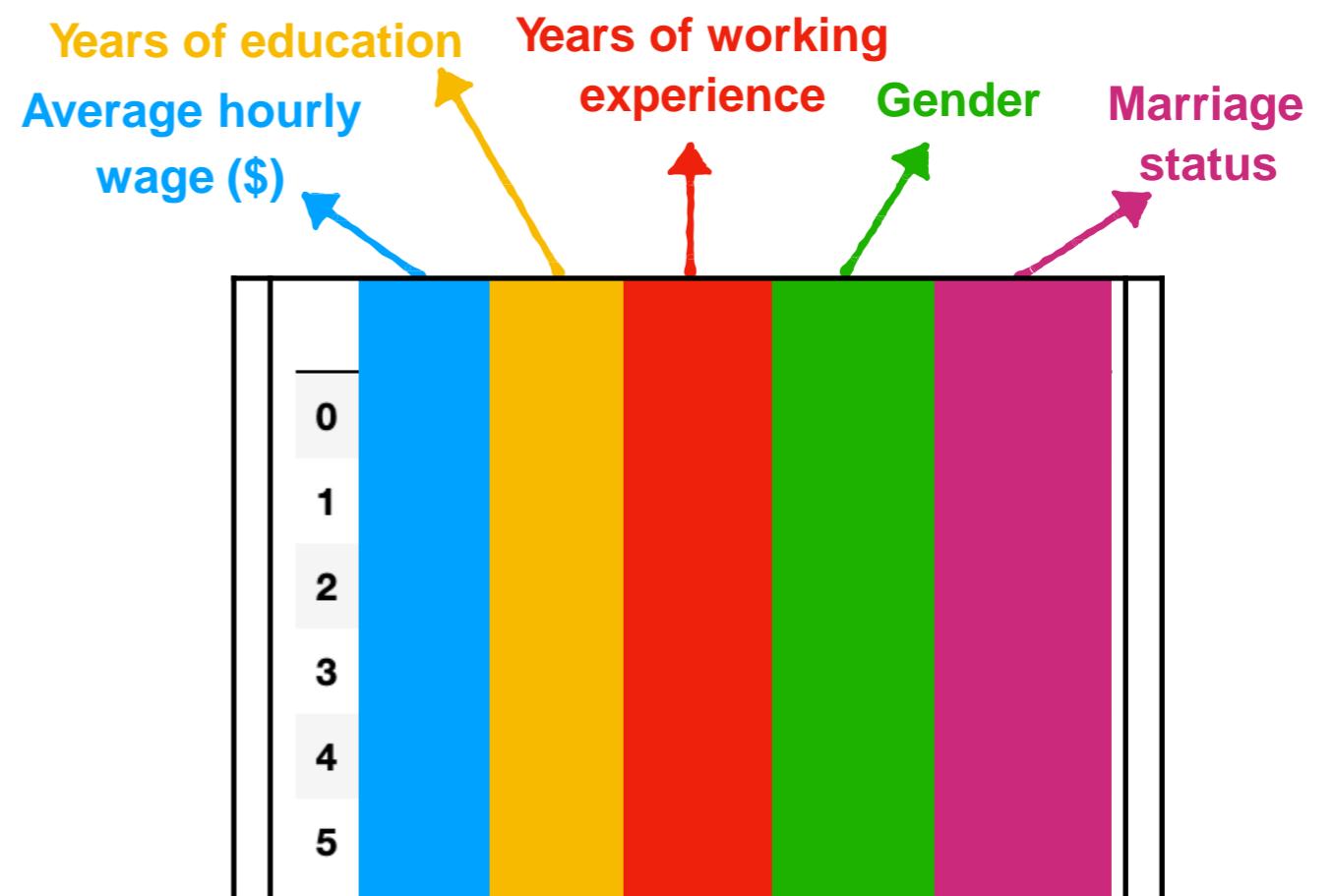
- Data representation

## Example 1

The data table contains a cross-sectional dataset on a number of working individuals for the year 1976. Explore the dataset to gain some basic knowledge of the historical data.

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

- Data representation
  - Variables (fields/attributes) as column labels



# Datasets

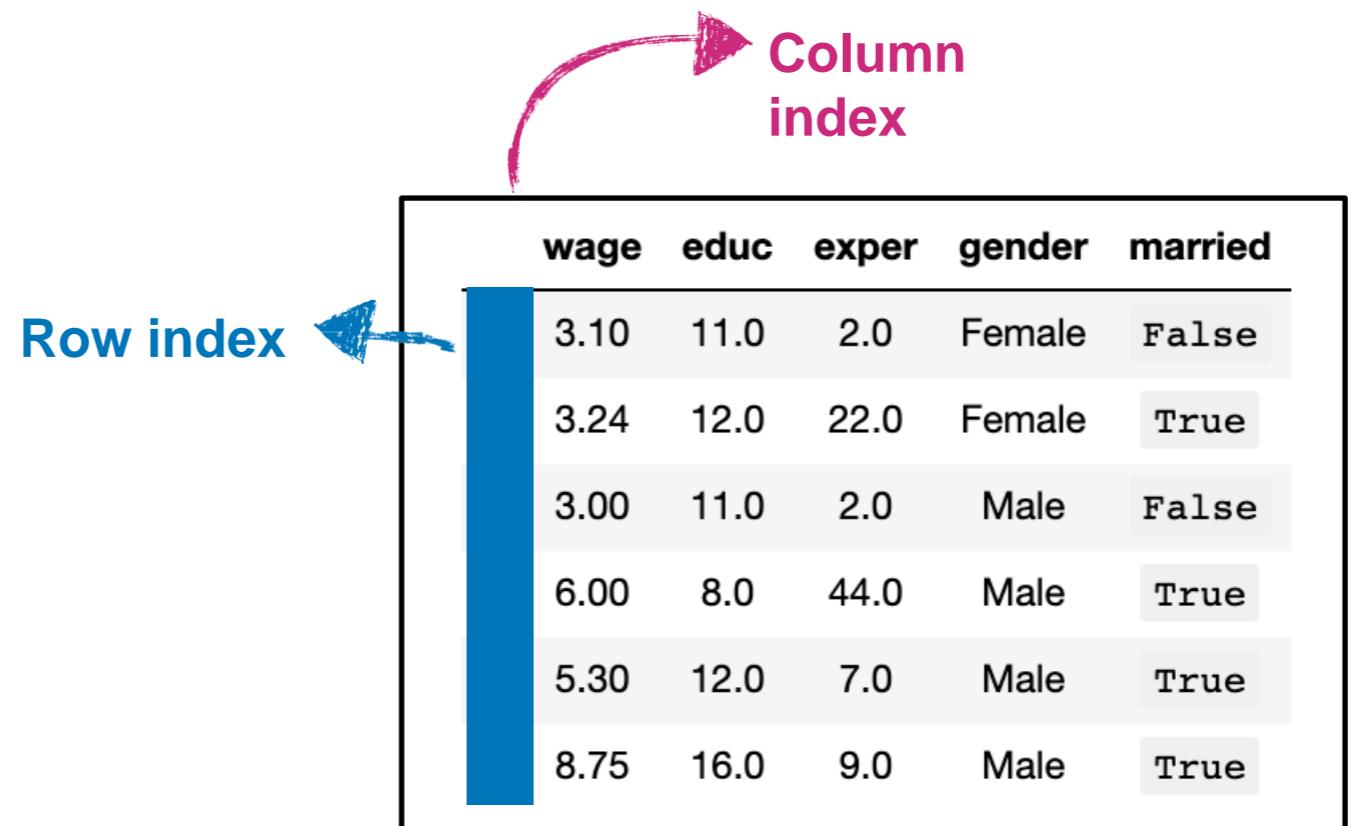
- Data representation
    - ▶ Variables (fields/attributes) as column labels
    - ▶ Observations (cases/records) as rows

A horizontal bar chart with five bars of equal length. The x-axis labels are 'wage', 'educ', 'exper', 'gender', and 'married'. The bars are colored blue, yellow, red, green, and pink from top to bottom.

# Datasets

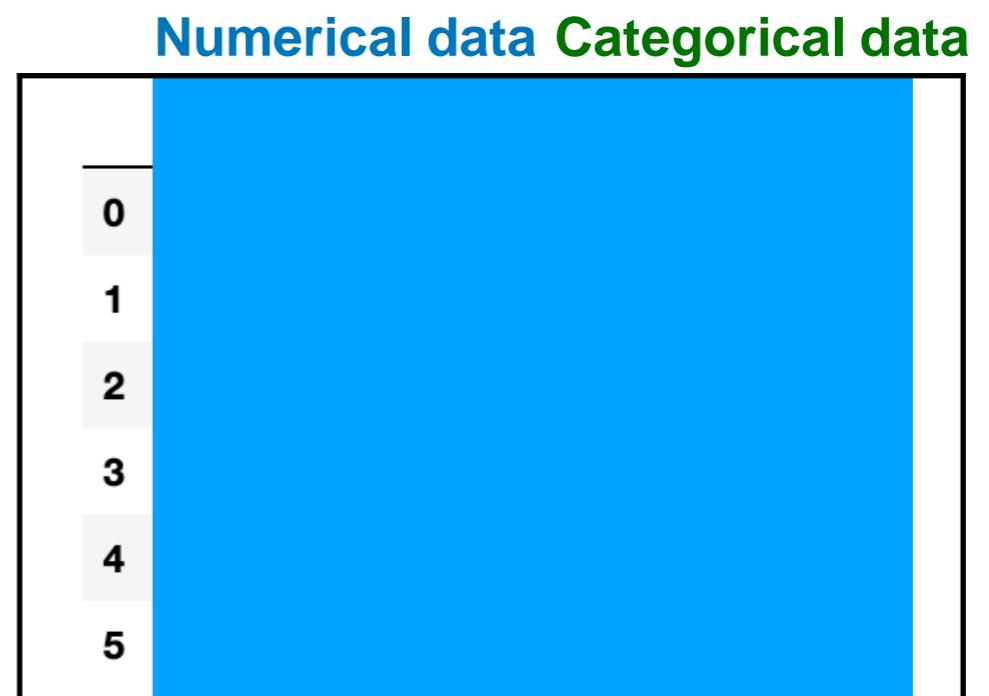
- Data representation

- ▶ Column index  
(label)
- ▶ Row index (label)



	wage	educ	exper	gender	married
1	3.10	11.0	2.0	Female	False
2	3.24	12.0	22.0	Female	True
3	3.00	11.0	2.0	Male	False
4	6.00	8.0	44.0	Male	True
5	5.30	12.0	7.0	Male	True
6	8.75	16.0	9.0	Male	True

- Types of variables
  - ▶ Numerical data
  - ▶ Categorical data



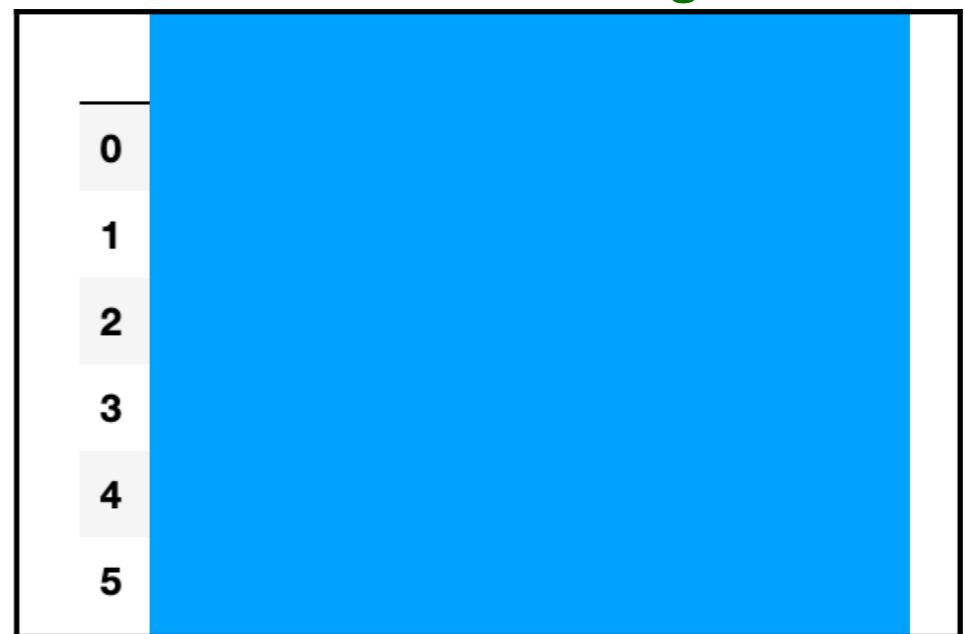
- Types of variables
  - ▶ Numerical (quantitative) data
  - ▶ Categorical (qualitative) data

```
print(int(True))  
print(int(False))
```

1  
0



Numerical data Categorical data



# Pandas for Data Analysis

---

- Introduction to Pandas
  - ▶ Labels for variables (or observations)
  - ▶ Heterogenous data
  - ▶ Possible missing values

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

---

- Introduction to Pandas

```
import pandas as pd  
pd.<function name>
```

# Pandas for Data Analysis

- The **Series** data structure
  - ▶ A one-dimensional array of indexed data
  - ▶ Create **DataFrames** object

```
1 wage = pd.Series([3.10, 3.24, 3.00, 6.00, 5.30, 8.75])  
2 wage
```

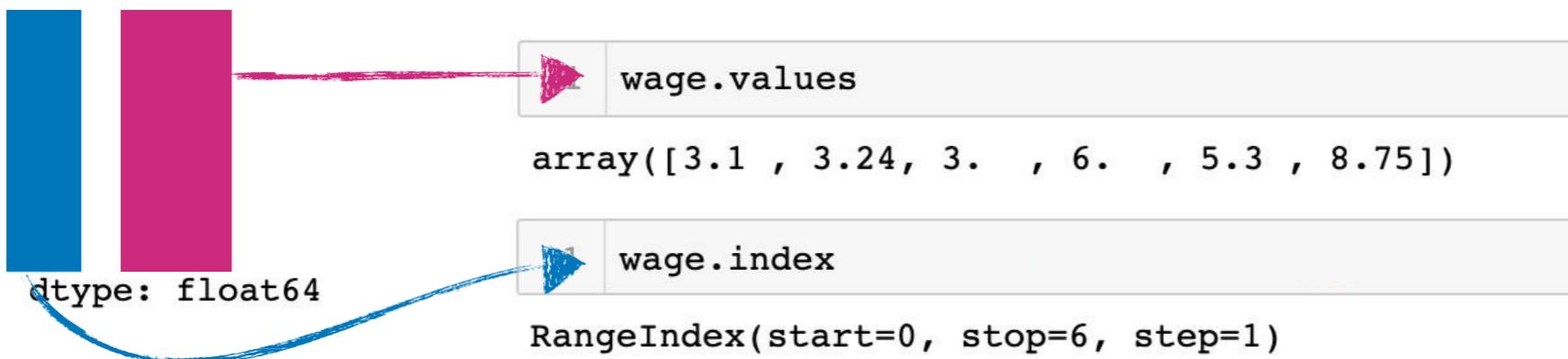
```
0    3.10  
1    3.24  
2    3.00  
3    6.00  
4    5.30  
5    8.75  
dtype: float64
```

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **Series** data structure
  - ▶ A one-dimensional array of indexed data
  - ▶ Create **Series** object

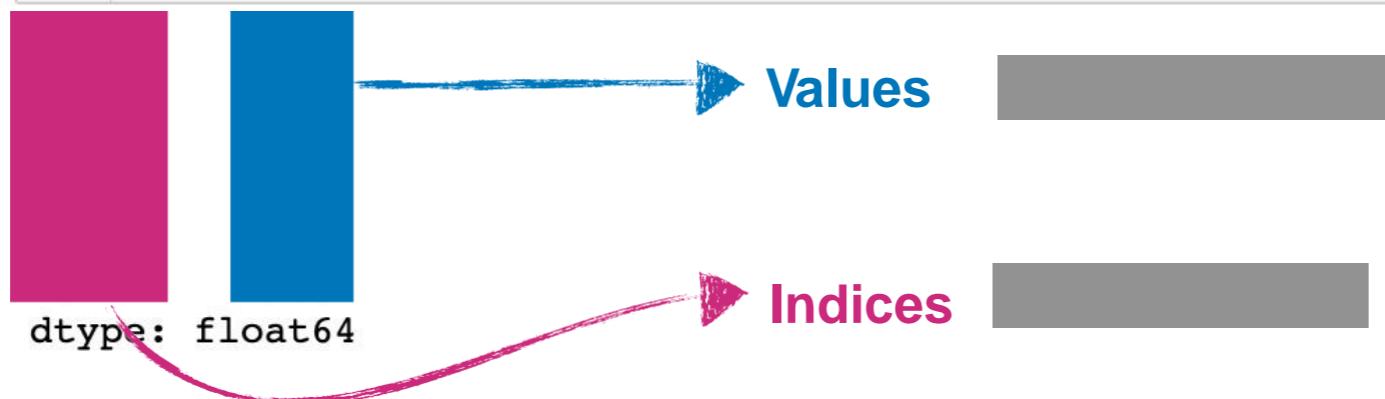
```
1 wage = pd.Series([3.10, 3.24, 3.00, 6.00, 5.30, 8.75])  
2 wage
```



# Pandas for Data Analysis

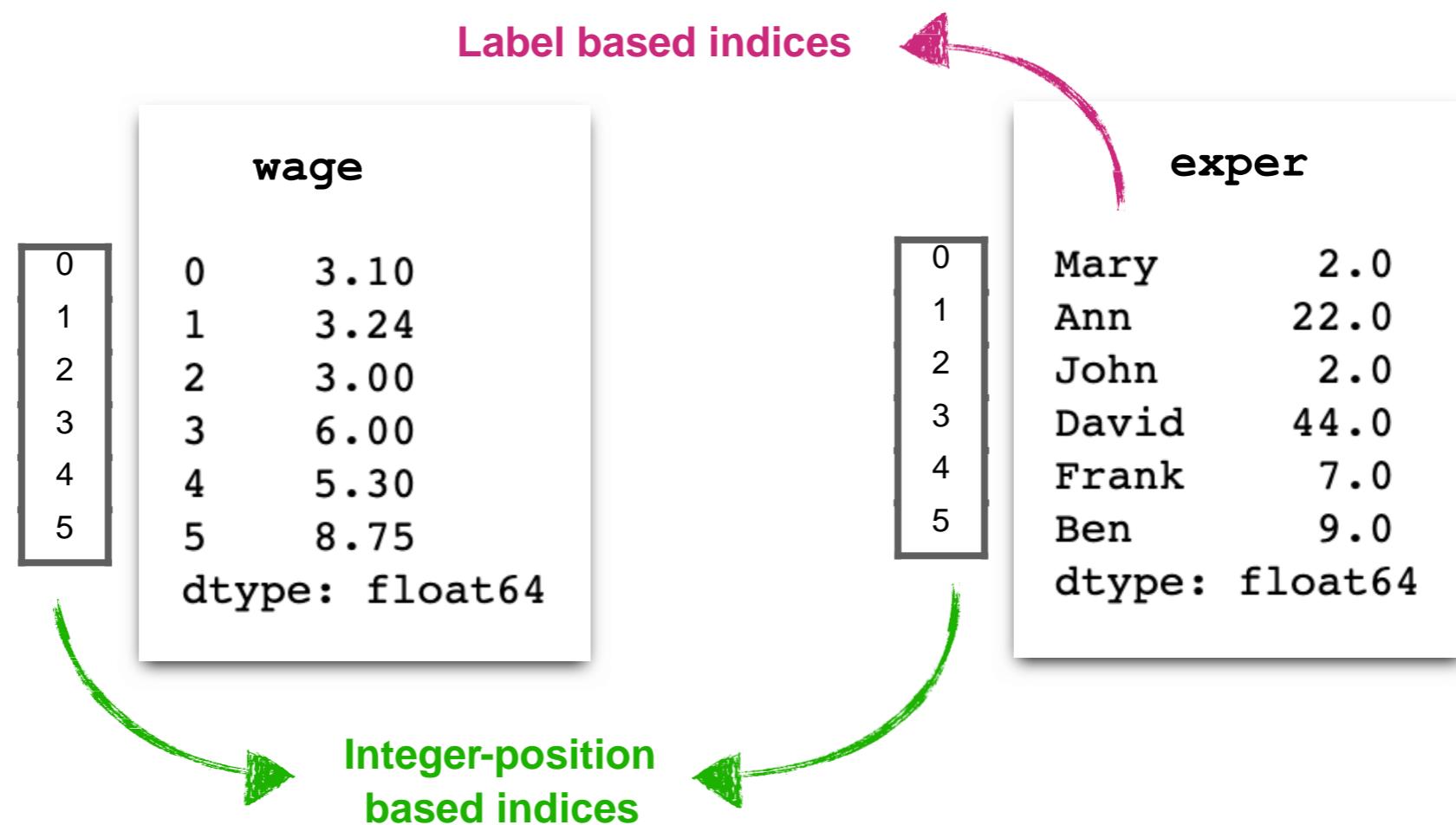
- The **Series** data structure
  - ▶ A one-dimensional array of indexed data
  - ▶ Create **pandas.Series** object

```
1 index = ['Mary', 'Ann', 'John', 'David', 'Frank', 'Ben']
2 exper = pd.Series([2.0, 22.0, 2.0, 44.0, 7.0, 9.0],
3                   index=index) # Specify series indices
4
5 exper
```



# Pandas for Data Analysis

- The **Series** data structure
  - Indexing and slicing of `pandas.Series` objects via brackets



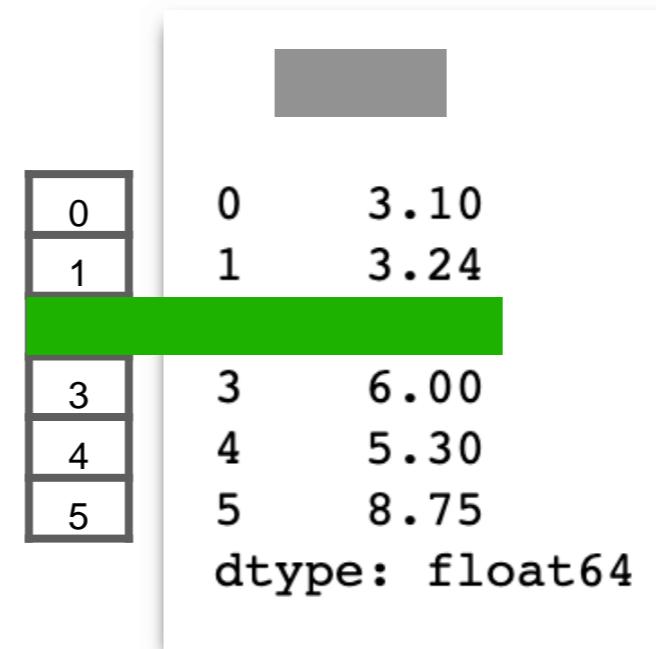
# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets
- ✓ Label based indices are the same as integer-position based indices

```
1 print(wage[2], '\n')
2 print(wage[:3], '\n')
3 print(wage[-2:])
```

```
3.0
0    3.10
1    3.24
2    3.00
dtype: float64

4    5.30
5    8.75
dtype: float64
```



0	0	3.10
1	1	3.24
2	3	6.00
3	4	5.30
4	5	8.75

dtype: float64

**Integer-position based indices**

# Pandas for Data Analysis

- The **Series** data structure
  - Indexing and slicing of **Series** objects via brackets
- ✓ Label based indices are the same as integer-position based indices

```
1 print(wage[2], '\n')
2 print(wage[1], '\n')
3 print(wage[-2:])
```

3.0



dtype: float64

4 5.30

5 8.75

dtype: float64



3	6.00
4	5.30
5	8.75
dtype: float64	

Integer-position  
based indices

# Pandas for Data Analysis

- The **Series** data structure
  - Indexing and slicing of **Series** objects via brackets
  - ✓ Label based indices are the same as integer-position based indices

```
1 print(wage[2], '\n')
2 print(wage[:3], '\n')
3 print(wage)
```

```
3.0
0    3.10
1    3.24
2    3.00
dtype: float64
```

```
dtype: float64
```

-6	0
-5	1
-4	2
-3	3

```
0    3.10
1    3.24
2    3.00
3    6.00
dtype: float64
```

**Integer-position  
based indices**

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets
  - ✓ Label based indices different from integer-position based indices

```
1 print(exper['John'])  
2 print(exper['Ann'])
```

22.0

**exper**

0	Mary	2.0
1	Ann	22.0
2	David	44.0
3	Frank	7.0
4	Ben	9.0
5		

**Label based indices**

dtype: float64

# Pandas for Data Analysis

- The **Series** data structure
  - Indexing and slicing of **Series** objects via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper['John'])  
2 print(exper[2])
```

exper

0	Mary 2.0
1	John 2.0
2	David 44.0
3	Frank 7.0
4	Ben 9.0
5	

Label based indices

dtype: float64

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper['John'], '\n')
2 print(exper[:'Ann'], '\n')
3 print(exper[['John', 'Ann']])
```

  
dtype: float64

```
Mary      2.0
Ann      22.0
dtype: float64

John      2.0
Ann      22.0
dtype: float64
```

0  
1  
2  
3  
4  
5

**exper**

0	Mary	2.0
1	Ann	22.0
2		
3		
4		
5		

  
dtype: float64

**Label based indices**

# Pandas for Data Analysis

- The **Series** data structure
  - Indexing and slicing of **Series** objects via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper['John'], '\n')
2 print(exper['David'], '\n')
3 print(exper[['John', 'Ann']])
```

```
John      2.0
David    44.0
Frank      7.0
Ben       9.0
dtype: float64
```

```
John      2.0
Ann     22.0
dtype: float64
```

0  
1  
2  
3  
4  
5

**exper**

0	
1	
2	
3	
4	
5	

```
John      2.0
David    44.0
Frank      7.0
Ben       9.0
dtype: float64
```

**Label based indices**

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper['John'], '\n')
2 print(exper['David'], '\n')
3 print(exper[['John', 'Ann']])
```

```
John      2.0
David    44.0
Frank      7.0
Ben       9.0
dtype: float64
```

```
John      2.0
Ann     22.0
dtype: float64
```

exper	
0	
1	
2	

John 2.0  
David 44.0

**Label based indices**

## Notes

In the slicing expressions for label based indices, the item indexed by stop are included in the selection. This is different from the list style slicing syntax we learned before.

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper['John'], '\n')
2 print(exper[:'Ann'], '\n')
3 print(exper)
```

```
John      2.0
David    44.0
Frank      7.0
Ben       9.0
dtype: float64
```

```
Mary      2.0
Ann      22.0
dtype: float64
```

```
dtype: float64
```

0  
1  
2  
3  
4  
5

**exper**

0	Mary      2.0
1	David    44.0
2	Frank      7.0
3	Ben       9.0
4	
5	

Label based indices

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper[3], '\n')
2 print(exper[:2], '\n')
3 print(exper[-3:], '\n')
4 print(exper[[2, 0, 1]])
```

```
Mary      2.0
Ann      22.0
dtype: float64
```

```
David     44.0
Frank      7.0
Ben       9.0
dtype: float64
```

```
John      2.0
Mary      2.0
Ann      22.0
dtype: float64
```

exper		
0	Mary	2.0
1	Ann	22.0
2	John	2.0
3		
4	Frank	7.0
5	Ben	9.0

dtype: float64

Integer-position  
based indices

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper[3], '\n')
2 print([4], '\n')
3 print(exper[-3:], '\n')
4 print(exper[[2, 0, 1]])
```

44.0

```
[4]
dtype: float64
```

```
David    44.0
Frank    7.0
Ben      9.0
dtype: float64
```

```
John    2.0
Mary    2.0
Ann    22.0
dtype: float64
```

exper		
2	John	2.0
3	David	44.0
4	Frank	7.0
5	Ben	9.0

dtype: float64

Integer-position  
based indices

# Pandas for Data Analysis

- The **data structure** is:
  - Indexing and slicing of **objects** via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper[3], '\n')
2 print(exper[:2], '\n')
3 print(exper[2], '\n')
4 print(exper[[2, 0, 1]])
```

```
44.0
Mary    2.0
Ann    22.0
dtype: float64
```

```
John    2.0
Mary    2.0
Ann    22.0
dtype: float64
```

```
John    2.0
Mary    2.0
Ann    22.0
dtype: float64
```

exper		
Mary	0	2.0
Ann	1	22.0
John	2	2.0

dtype: float64

Integer-position  
based indices

# Pandas for Data Analysis

- The **Series** data structure
  - Indexing and slicing of **Series** objects via brackets

✓ Label based indices different from integer-position based indices

```
1 print(exper[3], '\n')
2 print(exper[:2], '\n')
3 print(exper[-3:], '\n')
4 print(exper)
```

44.0

Mary 2.0

Ann 22.0

dtype: float64

David 44.0

Frank 7.0

Ben 9.0

dtype: float64

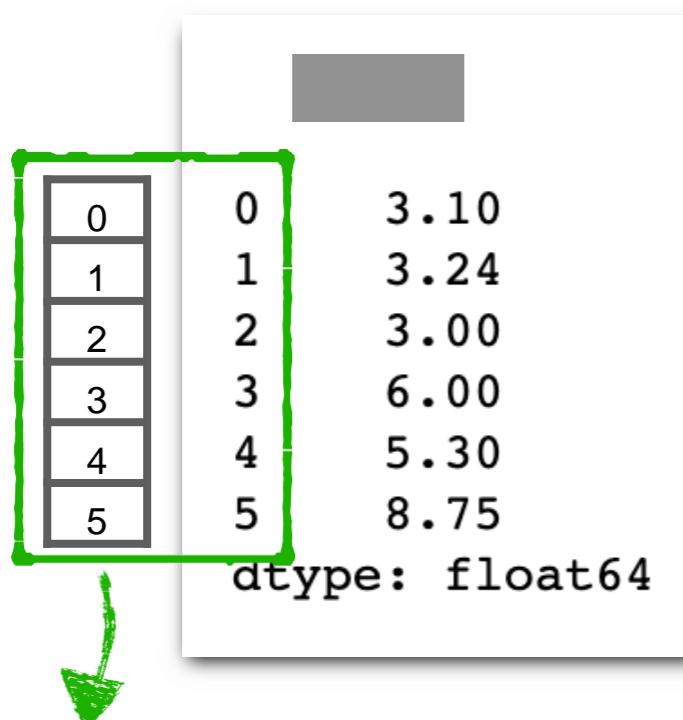
dtype: float64

exper		
1		
2		
3	David	44.0
4	Frank	7.0
5	Ben	9.0
		dtype: float64

Integer-position  
based indices

# Pandas for Data Analysis

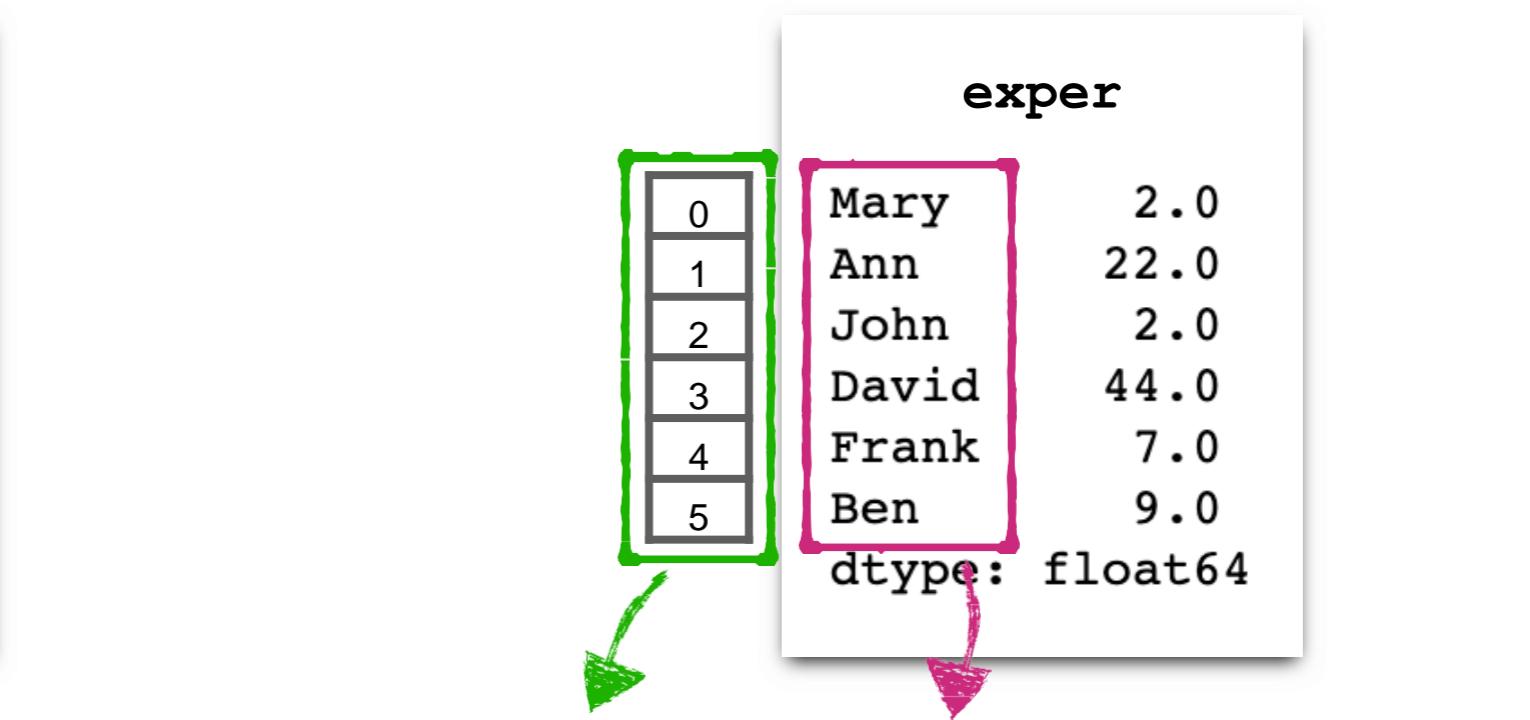
- The **data structure** is:
  - Indexing and slicing of **objects** via brackets



0	3.10
1	3.24
2	3.00
3	6.00
4	5.30
5	8.75

`wage`  
dtype: float64

```
1 print(wage[2], '\n')
2 print(wage[:3], '\n')
3 print(wage[-2:])
```



0	Mary	2.0
1	Ann	22.0
2	John	2.0
3	David	44.0
4	Frank	7.0
5	Ben	9.0

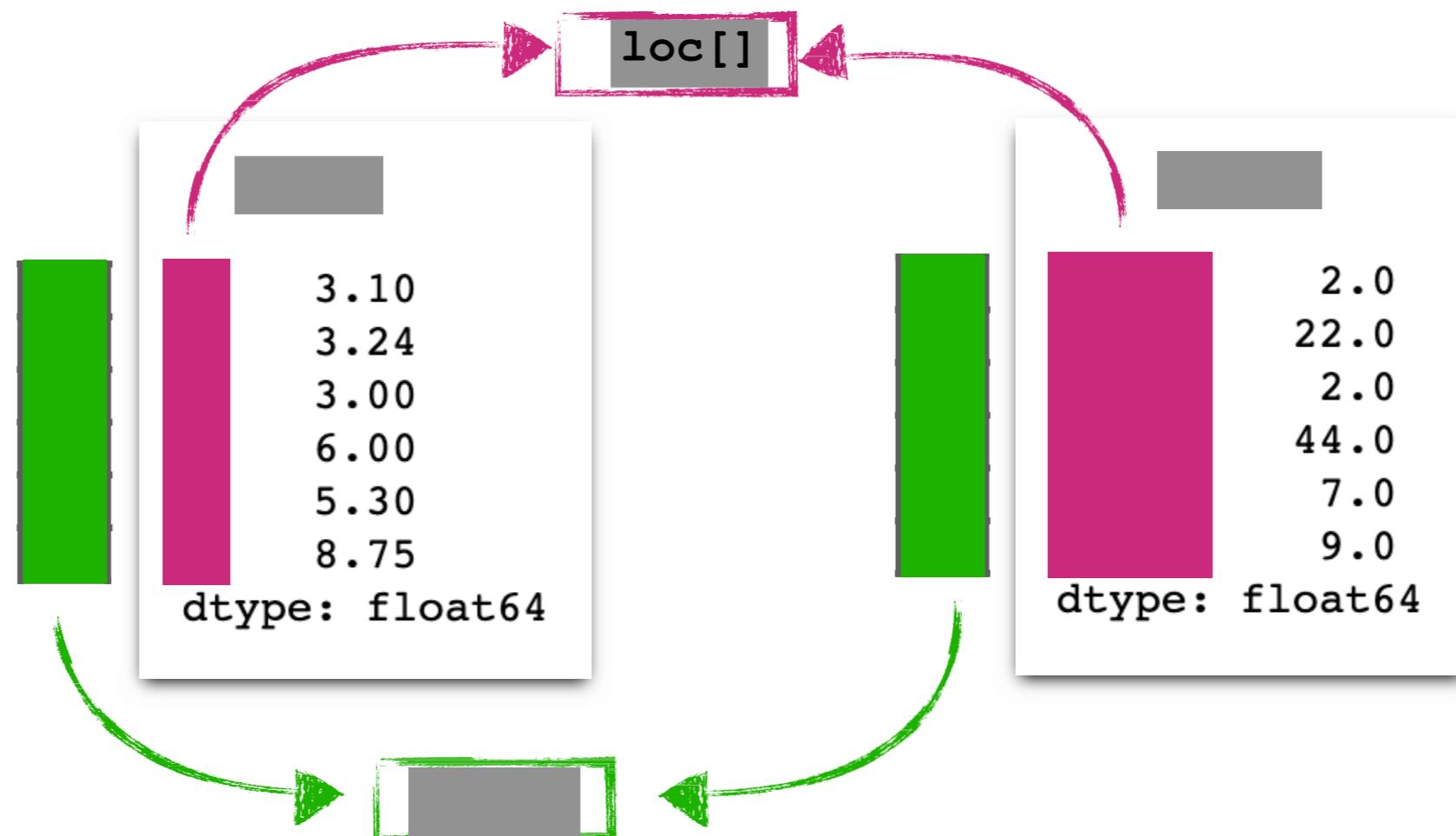
`exper`  
dtype: float64

```
1 print(exper[3], '\n')
2 print(exper[:2], '\n')
3 print(exper[-3:], '\n')
4 print(exper[[2, 0, 1]])
```

```
1 print(exper['John'], '\n')
2 print(exper[:'Ann'], '\n')
3 print(exper[['John', 'Ann']])
```

# Pandas for Data Analysis

- The **Series** data structure
  - Indexers: **loc[]** and **iloc[]**



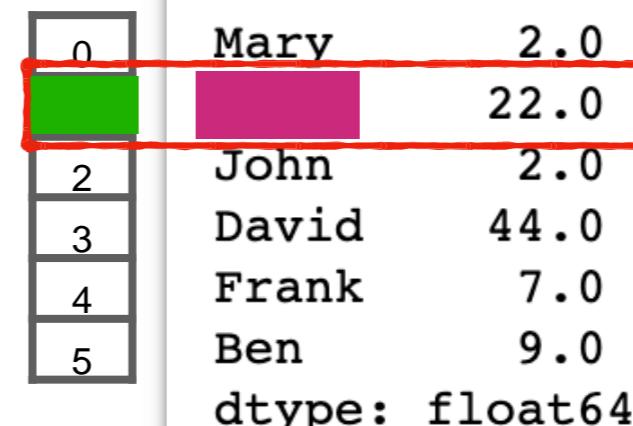
# Pandas for Data Analysis

- The **Series** data structure
- Indexers: **loc** and **iloc**

```
1 print(exper.loc['Ann'], '\n')
2 print(exper.loc['Frank'], '\n')
3 print(exper.loc[['Mary', 'John']], '\n')
4 print(exper.loc[::2])
```

```
1 print(exper['Age'], '\n')
2 print(exper.iloc[4:], '\n')
3 print(exper.iloc[[0, 2]], '\n')
4 print(exper.iloc[::2])
```

22.0



0	Mary	2.0
1	John	22.0
2	David	2.0
3	Frank	44.0
4	Ben	7.0
5		

dtype: float64

# Pandas for Data Analysis

- The **Series** data structure
- Indexers: **loc** and **iloc**

```
1 print(exper.loc['Ann'], '\n')
2 print(exper.loc['John'], '\n')
3 print(exper.loc[['Mary', 'John']], '\n')
4 print(exper.loc[::2])
```

```
1 print(exper.iloc[1], '\n')
2 print(exper.iloc[2], '\n')
3 print(exper.iloc[[0, 2]], '\n')
4 print(exper.iloc[::2])
```

```
Frank    7.0
Ben      9.0
dtype: float64
```

0	Mary	2.0
1	Ann	22.0
2	John	2.0
3	David	44.0
		7.0
		9.0

dtype: float64

# Pandas for Data Analysis

- The **Series** data structure
- Indexers: **loc** and **iloc**

```
1 print(exper.loc['Ann'], '\n')
2 print(exper.loc['Frank':], '\n')
3 print(exper.loc[1:4], '\n')
4 print(exper.loc[::2])
```

```
1 print(exper.iloc[1], '\n')
2 print(exper.iloc[4:], '\n')
3 print(exper.iloc[1:4], '\n')
4 print(exper.iloc[::2])
```

```
Mary    2.0
John    2.0
dtype: float64
```

1	Ann	22.0
2	David	44.0
3	Frank	7.0
4	Ben	9.0
5		

2.0  
2.0  
dtype: float64

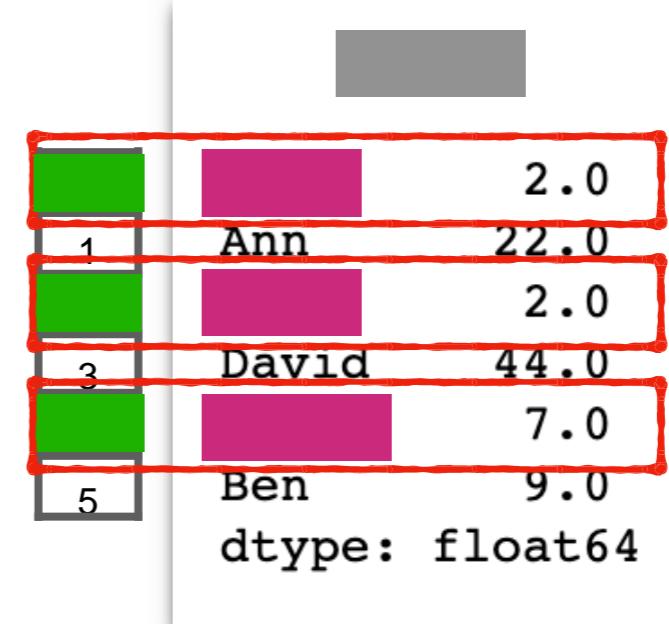
# Pandas for Data Analysis

- The **Series** data structure
- Indexers: **loc** and **iloc**

```
1 print(exper.loc['Ann'], '\n')
2 print(exper.loc['Frank':], '\n')
3 print(exper.loc[['Mary', 'John']], '\n')
4 print(exper.loc[0:2], '\n')
```

```
1 print(exper.iloc[1], '\n')
2 print(exper.iloc[4:], '\n')
3 print(exper.iloc[[0, 2]], '\n')
4 print(exper.iloc[0:3], '\n')
```

```
Mary      2.0
John      2.0
Frank     7.0
dtype: float64
```



```
0      2.0
1      Ann
2     22.0
3      2.0
4     David
5     44.0
6      7.0
7     Ben
8      9.0
dtype: float64
```

# Pandas for Data Analysis

- The **Series** data structure
- Indexers: **loc** and **iloc**

```
1 print(wage.loc[:2], '\n')
2 print(wage.iloc[:2])
```

Name: wage, dtype: float64

0 3.10  
1 3.24

Name: wage, dtype: float64

is included

0
1
2
3
4
5

6.00  
5.30  
8.75  
dtype: float64

# Pandas for Data Analysis

- The **Series** data structure
- Indexers: **loc** and **iloc**

```
1 print(wage.loc[:2], '\n')
2 print(wage.iloc[3:5])
```

```
0    3.10
1    3.24
2    3.00
Name: wage, dtype: float64
```

```
Name: wage, dtype: float64
```

**is excluded**

0	3.10
1	3.24
2	3.00
3	6.00
4	5.30
5	8.75

dtype: float64

# Pandas for Data Analysis

---

- The **Series** data structure
  - Indexers: `loc[]` and `iloc[]`
    - ✓ Label based indices are expressed by `loc[]`, with the `True` index included in the selection
    - ✓ Integer-position based indices are expressed by `iloc[]` with the `False` index excluded from the selection

## Coding Style

One guiding principle of Python code is that "explicit is better than implicit." The explicit nature of `loc[]` and `iloc[]` make them very useful in maintaining clean and readable code; especially in the case of integer indexes.

# Pandas for Data Analysis

- The `pandas.DataFrame` data structure
  - ▶ Create `pandas.DataFrame` object

```
1 data_dict = {'wage': [3.10, 3.24, 3.00, 6.00, 5.30, 8.75],  
2             'educ': [11.0, 12.0, 11.0, 8.0, 12.0, 16.0],  
3             'exper': [2.0, 22.0, 2.0, 44.0, 7.0, 9.0],  
4             'gender': ['Female', 'Female', 'Male', 'Male', 'Male', 'Male'],  
5             'married': [False, True, False, True, True, True]}  
6  
7 data_frame = pd.DataFrame(data_dict)  
8 data_frame
```

```
1 type(data_frame)
```



```
1 type(data_frame)
```

```
pandas.core.frame.DataFrame
```

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**
  - ▶ Create **object**

```
1 data_dict = {'wage': [3.10, 3.24, 3.00, 6.00, 5.30, 8.75],  
2             'educ': [11.0, 12.0, 11.0, 8.0, 12.0, 16.0],  
3             'exper': [2.0, 22.0, 2.0, 44.0, 7.0, 9.0],  
4             'gender': ['Female', 'Female', 'Male', 'Male', 'Male', 'Male'],  
5             'married': [False, True, False, True, True, True]}
```

```
6  
7 data_frame = pd.DataFrame(data_dict)  
8 data_frame
```

Constructor

```
1 type(data_dict)
```

dict

```
1 type(
```

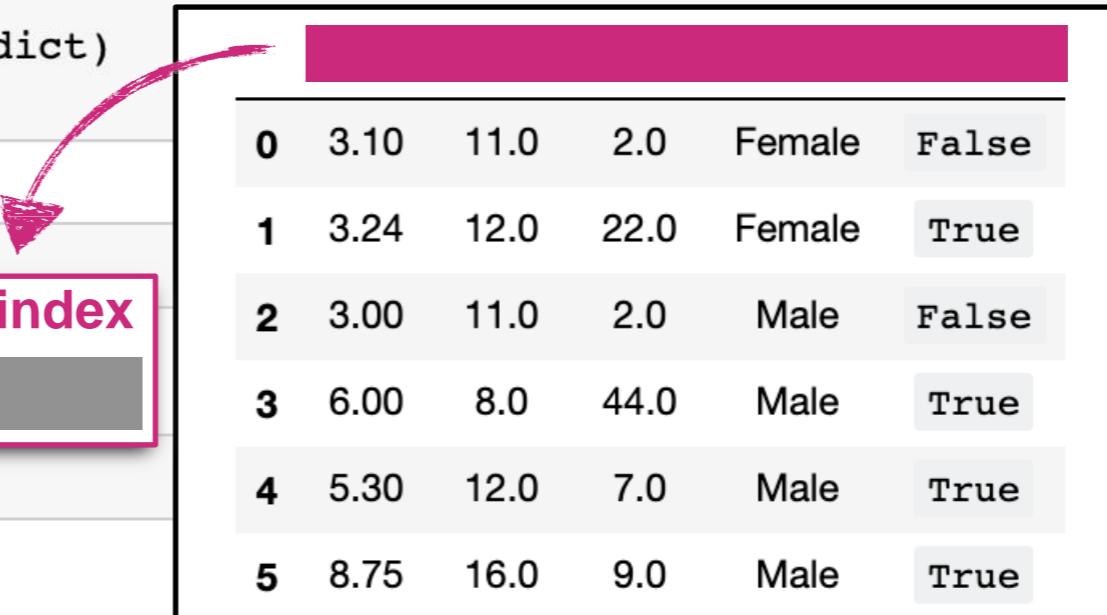
	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **Dataframe** data structure
- ▶ Create **Dataframe** object

```
1 data_dict = {Column index: [3.10, 3.24, 3.00, 6.00, 5.30, 8.75],  
2 Column index: [11.0, 12.0, 11.0, 8.0, 12.0, 16.0],  
3 Column index: [2.0, 22.0, 2.0, 44.0, 7.0, 9.0],  
4 Column index: ['Female', 'Female', 'Male', 'Male', 'Male', 'Male'],  
5 Column index: [False, True, False, True, True, True]}  
6  
7 data_frame = pd.DataFrame(data_dict)  
8 data_frame
```

```
1 type(data_dict)  
dict  
1 type(data_frame)  
pandas.core.frame.DataFrame
```



	0	1	2	3	4	5
<b>0</b>	3.10	11.0	2.0	Female	False	
<b>1</b>	3.24	12.0	22.0	Female	True	
<b>2</b>	3.00	11.0	2.0	Male	False	
<b>3</b>	6.00	8.0	44.0	Male	True	
<b>4</b>	5.30	12.0	7.0	Male	True	
<b>5</b>	8.75	16.0	9.0	Male	True	

# Pandas for Data Analysis

- The **Dataframe** data structure
  - ▶ Create **Dataframe** object

```
1 data_dict = {'wage': [3.10, 3.24, 3.00, 6.00, 5.30, 8.75],  
2             'educ': [11.0, 12.0, 11.0, 8.0, 12.0, 16.0],  
3             'exper': [2.0, 22.0, 2.0, 44.0, 7.0, 9.0],  
4             'gender': ['Female', 'Female', 'Male', 'Male', 'Male', 'Male'],  
5             'married': [False, True, False, True, True, True]}  
6  
7 data_frame = pd.DataFrame(data_dict)  
8 data_frame
```

```
1 type(data_dict)  
dict  
1 type(data_frame)  
pandas.core.frame.DataFrame
```

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True
" data-bbox="340 580 880 880"/>

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **Dataframe** data structure
  - ▶ Create **Dataframe** object

```
1 data_dict = {'wage': [11.0, 12.0, 11.0, 8.0, 12.0, 16.0],  
2             'educ': [11.0, 12.0, 11.0, 8.0, 12.0, 16.0],  
3             'exper': [2.0, 22.0, 2.0, 44.0, 7.0, 9.0],  
4             'gender': ['Female', 'Female', 'Male', 'Male', 'Male', 'Male'],  
5             'married': [False, True, False, True, True, True]})  
6  
7 data_frame = pd.DataFrame(data_dict)  
8 data_frame
```

```
1 type(data_dict)
```

dict

```
1 type(data_frame)
```

pandas.core.frame.DataFrame

	wage	educ	exper	gender	married
0		11.0	2.0	Female	False
1		12.0	22.0	Female	True
2		11.0	2.0	Male	False
3		8.0	44.0	Male	True
4		12.0	7.0	Male	True
5		16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**

- ▶ Create **pandas.DataFrame** object

```
1 index = ['Mary', 'Ann', 'John', 'David', 'Frank', 'Ben']
2 data_frame_new = pd.DataFrame(data_dict, index=index)
3
4 data_frame_new
```



	wage	educ	exper	gender	married
<b>Mary</b>	3.10	11.0	2.0	Female	False
<b>Ann</b>	3.24	12.0	22.0	Female	True
<b>John</b>	3.00	11.0	2.0	Male	False
<b>David</b>	6.00	8.0	44.0	Male	True
<b>Frank</b>	5.30	12.0	7.0	Male	True
<b>Ben</b>	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

---

- The  data structure
  - Accessing columns and rows via brackets 
  - ✓ Columns: a single column index or a list of column indices

# Pandas for Data Analysis

- The **data structure**

e

- Accessing columns and rows via brackets

✓ Columns: a single column index or a list of column indices

```
1 wage = data_frame['wage']  
2 wage
```

```
0    3.10  
1    3.24  
2    3.00  
3    6.00  
4    5.30  
5    8.75  
Name: wage, dtype: float64
```

```
1 type(wage)
```

pandas.core.series.Series

		educ	exper	gender	married
Mary		11.0	2.0	Female	False
Ann		12.0	22.0	Female	True
John		11.0	2.0	Male	False
David		8.0	44.0	Male	True
Frank		12.0	7.0	Male	True
Ben		16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**

e

- Accessing columns and rows via brackets

✓ Columns: a single column index or

```
1 skills = data_frame[['educ', 'exper']]  
2 skills
```

	educ	exper
0	11.0	2.0
1	12.0	22.0
2	11.0	2.0
3	8.0	44.0
4	12.0	7.0
5	16.0	9.0

	wage	gender	married
Mary	3.10	Female	False
Ann	3.24	Female	True
John	3.00	Male	False
David	6.00	Male	True
Frank	5.30	Male	True
Ben	8.75	Male	True

# Pandas for Data Analysis

---

- The **pandas.DataFrame** data structure
  - Accessing columns and rows via brackets `[]`
    - ✓ Rows: a slicing expression of the label based or integer-position based indices

# Pandas for Data Analysis

- The **data frame** data structure
    - ▶ Accessing columns and rows via brackets
- ✓ Rows: a slicing expression of the label based or integer-position

```
1 john = data_frame_new[2:3]
2 john
```

	wage	educ	exper	gender	married
John	3.0	11.0	2.0	Male	False
Mary	3.10	11.0	2.0	Female	False
Ann	3.24	12.0	22.0	Female	True
David	6.00	8.0	44.0	Male	True
Frank	5.30	12.0	7.0	Male	True
Ben	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**

e

- Accessing columns and rows via brackets

✓ Rows: **of the label based or integer-position**

```
1 david = data_frame_new[ 'David':'Frank' ]  
2 david
```

	wage	educ	exper	gender	married
<b>David</b>	6.0	8.0	44.0	Male	True
<b>Frank</b>	5.3	12.0	7.0	Male	True

	wage	educ	exper	gender	married
<b>Mary</b>	3.10	11.0	2.0	Female	False
<b>Ann</b>	3.24	12.0	22.0	Female	True
<b>John</b>	3.00	11.0	2.0	Male	False
<b>Ben</b>	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**
- Accessing columns and rows via brackets `[]`

## Question

Given a data table, how to access the highlighted subset?

- A. `1 data[1][[2, 3]]`
- B. `1 data[2][[1, 2]]`
- C. `1 data[2:3][[1, 2]]`
- D. `1 data[1:2][[1, 2]]`
- E. `1 data[1:2][[2, 3]]`
- F. `1 data[2:3][[1, 2]]`

	1	2	3	4
1	1	2	3	4
2	5			8
3	9	10	11	12
4	13	14	15	16

# Pandas for Data Analysis

---

- The  data structure
  - Accessing columns and rows via brackets 
  - ✓ Columns: a single column index or a list of column indices
  - ✓ Rows: a slicing expression of the label based or integer-position based indices

# Pandas for Data Analysis

---

- The **dataframe** data structure
  - Indexers: **loc** and **iloc**
- 

## **.iloc selections - position based selection**

`data.iloc[<row selection>, <column selection>]`

*Integer list of rows: [0,1,2]  
Slice of rows: [4:7]  
Single values: 1*

*Integer list of columns: [0,1,2]  
Slice of columns: [4:7]  
Single column selections: 1*

---

## **loc selections - label based selection**

`data.loc[<row selection>, <column selection>]`

*Index/Label value: 'john'  
List of labels: ['john', 'sarah']  
Logical/Boolean index: `data['age'] == 10`*

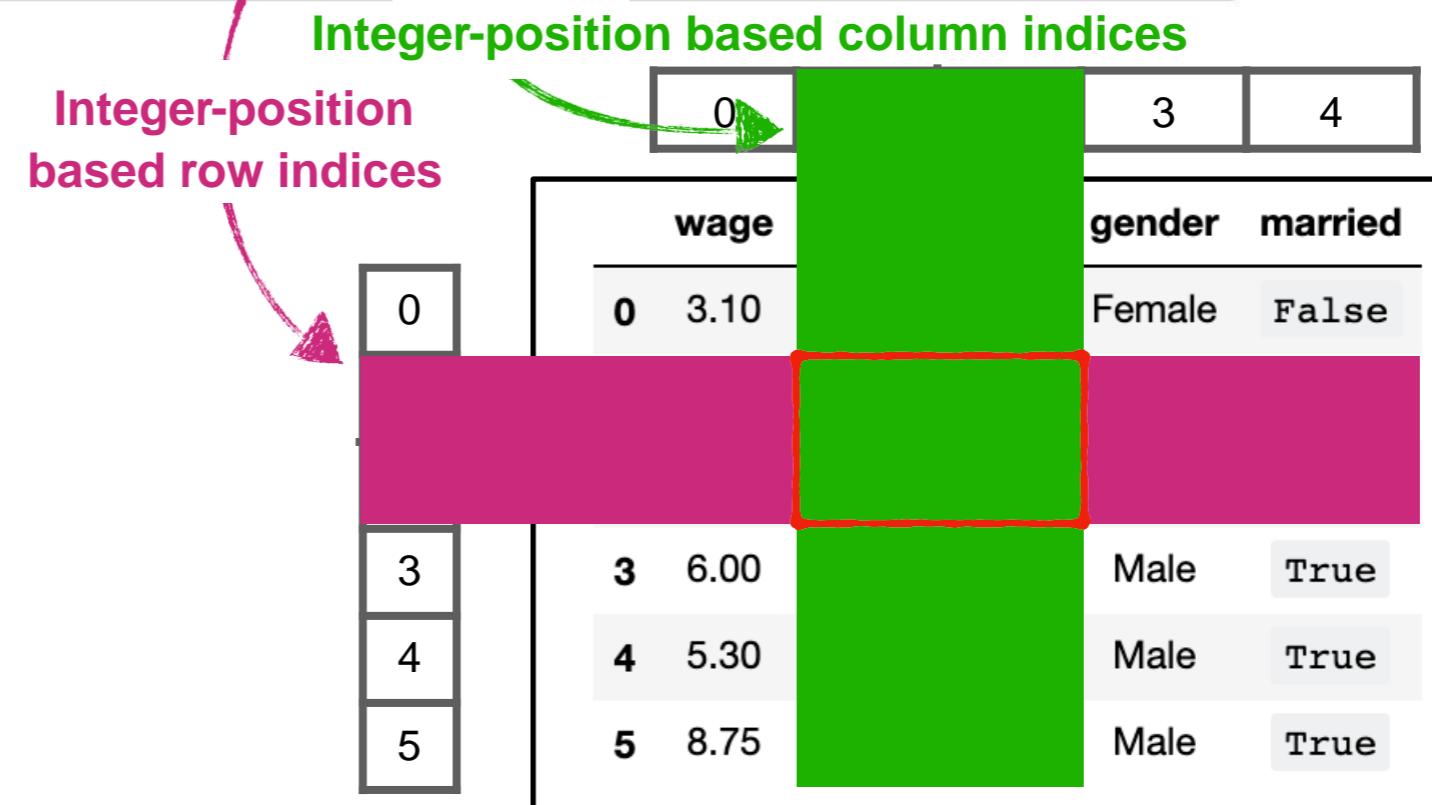
*Named column: 'first\_name'  
List of column names: ['first\_name', 'age']  
Slice of columns: 'first\_name':'address'*

# Pandas for Data Analysis

- The **data frame** data structure
- Indexers: **loc** and **iloc**

```
6 data_subset = data_frame.iloc[1:3, 1:3]
7 data_subset
```

	educ	exper
1	12.0	22.0
2	11.0	2.0



# Pandas for Data Analysis

- The **data frame** data structure
- Indexers: **loc** and **iloc**

```
6 data_new_subset = data_frame_new.iloc[1:3, 1:3]  
7 data_new_subset
```

	educ	exper
Ann	12.0	22.0
John	11.0	2.0

The diagram illustrates a Pandas DataFrame with the following structure:

	0	1	2	3	4
0	Ann	John	Mary	Female	False
1	12.0	2.0	3.10	Male	True
2	22.0		6.00	Male	True
3			5.30	Male	True
4			8.75		
5					
6					
7					

Annotations explain the indexing:

- Integer-position based row indices:** Points to the index row (0, 1, 2, 3, 4, 5, 6, 7) on the left of the DataFrame.
- Integer-position based column indices:** Points to the index column (0, 1, 2, 3, 4) at the top of the DataFrame.

Specific cells are highlighted in green and red:

- Cells at index 1, column 2 (John) and index 1, column 3 (3.10) are highlighted in green.
- A red box highlights the cells at index 1, column 2 (John) and index 1, column 3 (3.10).

# Pandas for Data Analysis

- The **data structure**
  - ▶ Indexers: **Label based row indices** and **Label based column indices**

```
6 data_subset = data_frame.loc[1:2, 'educ':'exper']
7 data_subset
```

	educ	exper
1	12.0	22.0
2	11.0	2.0

**Label based  
row indices**

**Label based column indices**

	wage	gender	married
0	3.10	Female	False
1	5.10	Male	True
2	5.10	Male	True
3	6.00	Male	True
4	5.30	Male	True
5	8.75	Male	True

# Pandas for Data Analysis

- The **data structure**
  - ▶ Indexers: **Label based row indices** and **Label based column indices**

```
6 data_subset = data_frame.loc[1:2, 'educ':'exper']
7 data_subset
```

	educ	exper
1	12.0	22.0
2	11.0	2.0

**Label based row indices**

## Notes

In the slicing expressions for label based indices (`loc[]`), the item indexed by stop are included in the selection. This is different from the list style slicing syntax we learned before.

**Label based column indices**

	wage	gender	married
0	3.10	Female	False
1	9.98	Male	True
2	4.54	Male	True
3	6.00	Male	True
4	5.30	Male	True
5	8.75	Male	True

# Pandas for Data Analysis

- The **data frame** data structure

- Indexers: **loc** and **iloc**

- ✓ Access all rows or columns via **loc**

```
1 data_frame.iloc[:, 2]  
0    2.0  
1   22.0  
2    2.0  
3   44.0  
4    7.0  
5    9.0  
Name: exper, dtype: float64
```

All rows  
Column indices

0	1	2	3	4
wage	educ	gender	married	
0				
1				
2				
3				
4				
5				

# Pandas for Data Analysis

- The **data frame** data structure
  - Indexers: **loc** and **iloc**

✓ Access all rows or columns via **loc**

The diagram illustrates the use of the `loc` indexer to select rows from a DataFrame. On the left, a code snippet shows two lines of Python code:

```
1 rows = data_frame_new.loc['John':'Ben', :]  
2 rows
```

A green arrow points from the colon in the slice to the text "All columns" in green. On the right, a diagram shows a DataFrame with 5 rows and 6 columns. The rows are labeled "Mary", "Ann", "John", "David", and "Frank". The columns are labeled "wage", "educ", "exper", "gender", and "married". The "gender" and "married" columns are highlighted in pink, and the "wage", "educ", and "exper" columns are highlighted in green. The "All indices" label is positioned above the "gender" column.

	wage	educ	exper	gender	married
John	3.00	11.0	2.0	Male	False
David	6.00	8.0	44.0	Male	True
Frank	5.30	12.0	7.0	Male	True
Ben	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**
  - ▶ Indexers:  **and**

✓ Access all rows or columns via

1 rows = data\_frame\_new.loc[ ]  
2 rows

All indices

	wage	educ	exper	gender	married
John	3.00	11.0	2.0	Male	False
David	6.00	8.0	44.0	Male	True
Frank	5.30	12.0	7.0	Male	True
Ben	8.75	16.0	9.0	Male	True

All columns as default

# Pandas for Data Analysis

- The **data structure**

- ▶ Indexers:  **and**

- ✓ Data types of indices

```
1 data_frame_new.loc['John']
```

```
wage          3
educ         11
exper         2
gender      Male
married     False
Name: John, dtype: object
```

```
1 data_frame_new.loc[['John']]
```

	wage	educ	exper	gender	married
John	3.0	11.0	2.0	Male	False

	wage	educ	exper	gender	married
Mary	3.10	11.0	2.0	Female	False
Ann	3.24	12.0	22.0	Female	True
David	6.00	8.0	44.0	Male	True
Frank	5.30	12.0	7.0	Male	True
Ben	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data structure**
  - ▶ Indexers: **iloc** and **loc**
  - ✓ Data types of indices

```
1 data_frame.iloc[:, 2]  
0    2.0  
1   22.0  
2    2.0  
3   44.0  
4    7.0  
5    9.0  
Name: exper, dtype: float64
```

0	1	3	4	
	wage	educ	gender	married
0	3.10	11.0	Female	False
1	3.24	12.0	Female	True
2	3.00	11.0	Male	False
3	6.00	8.0	Male	True
4	5.30	12.0	Male	True
5	8.75	16.0	Male	True

# Pandas for Data Analysis

- The **data frame** data structure

- ▶ Indexers: **loc** and **iloc**

- ✓ Data types of indices

```
1 data_frame.iloc[:, 2:3]
```

	exper
0	2.0
1	22.0
2	2.0
3	44.0
4	7.0
5	9.0

	0	1	3	4
	wage	educ	gender	married
0	3.10	11.0	Female	False
1	3.24	12.0	Female	True
2	3.00	11.0	Male	False
3	6.00	8.0	Male	True
4	5.30	12.0	Male	True
5	8.75	16.0	Male	True

# Pandas for Data Analysis

---

- The **pandas** library provides a **DataFrame** data structure
  - ▶ Modify a subset of a `pandas.DataFrame` object
    - ✓ Similar to modifying a part of a built-in compound data object (list/dictionary)

# Pandas for Data Analysis

- The **data frame** data structure
  - ▶ Modify a subset of a **data frame** object

```
1 data_frame.loc[2:3, 'educ'] = 9.0
2 data_frame
```

	wage	exper	gender	married
0	3.10	2.0	Female	False
1	3.24	22.0	Female	True
2		2.0	Male	False
3		44.0	Male	True
4	5.30	7.0	Male	True
5	8.75	9.0	Male	True



	wage	exper	gender	married	
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	9.0	2.0	Male	False
3	6.00	9.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

- The **data frame** data structure
  - Modify a subset of a **data frame** object

```
1 data_frame.iloc[2, 4] = 1.0  
2 data_frame
```

The diagram illustrates the modification of a data frame using the `iloc` selector. On the left, the original data frame is shown with a pink row highlighted at index 2. On the right, the modified data frame is shown with the value `1.0` at index 2, column 4, highlighted with a red box. The data frame contains columns `wage`, `educ`, `exper`, `gender`, and `married`.

	0	3	4	
	wage	gender	married	
0	3.10	Female	False	
1	3.24	Female	True	
2	3.00	Male	True	
3	6.00	Male	True	
4	5.30	Male	True	
5	8.75	Male	True	

	0	1	3	4	
	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	1.0	1.0	Male	False
3	6.00	9.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

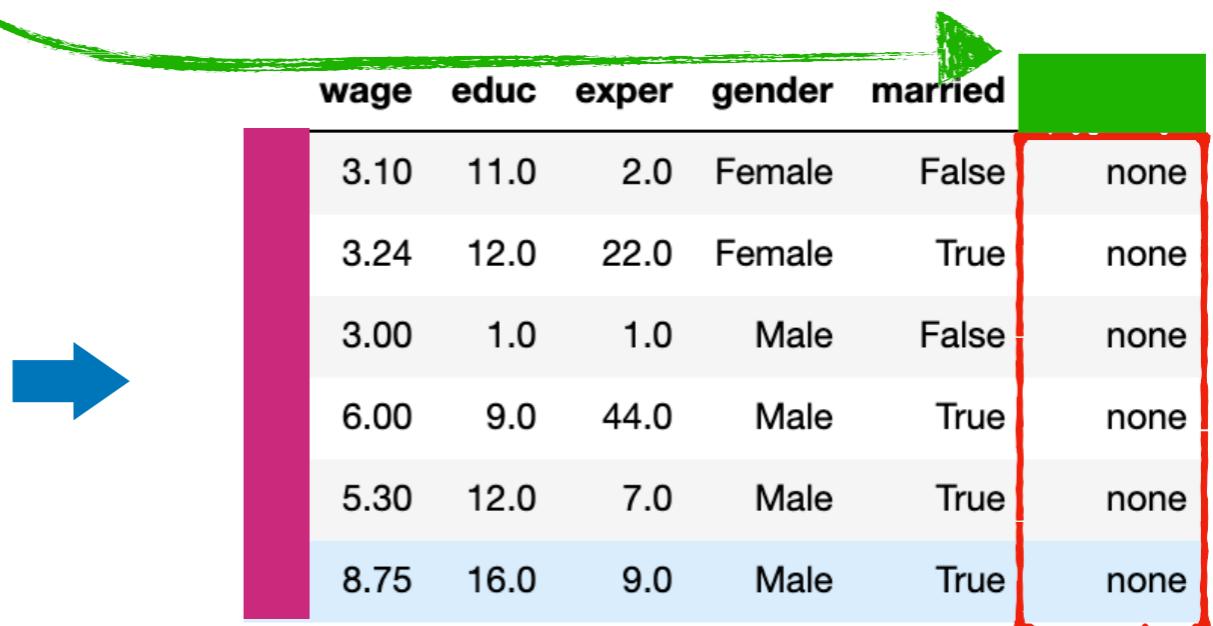
# Pandas for Data Analysis

- The **pandas** data structure
  - Modify a subset of a `pandas.DataFrame` object

```
1 data_frame.loc[1, 'married'] = 'none'  
2 data_frame
```

New label added as a new column

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	9.0	2.0	Male	False
3	6.00	9.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True



	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	1.0	1.0	Male	False
3	6.00	9.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

# Pandas for Data Analysis

---

- The **pandas** data structure
  - ▶ Boolean **pandas.Series** and boolean indexing

A subset for  
female workers?

	wage	educ	exper	gender	married	remarks
0						
1						
2						
3						
4						
5						

# Pandas for Data Analysis

- The **data structure**
  - ▶ Boolean **and boolean indexing**

```
1 is_female = data_frame['gender'] == 'Female'  
2 is_female  
  
0    True  
1    True  
2   False  
3   False  
4   False  
5   False  
  
Name: gender, dtype: bool
```

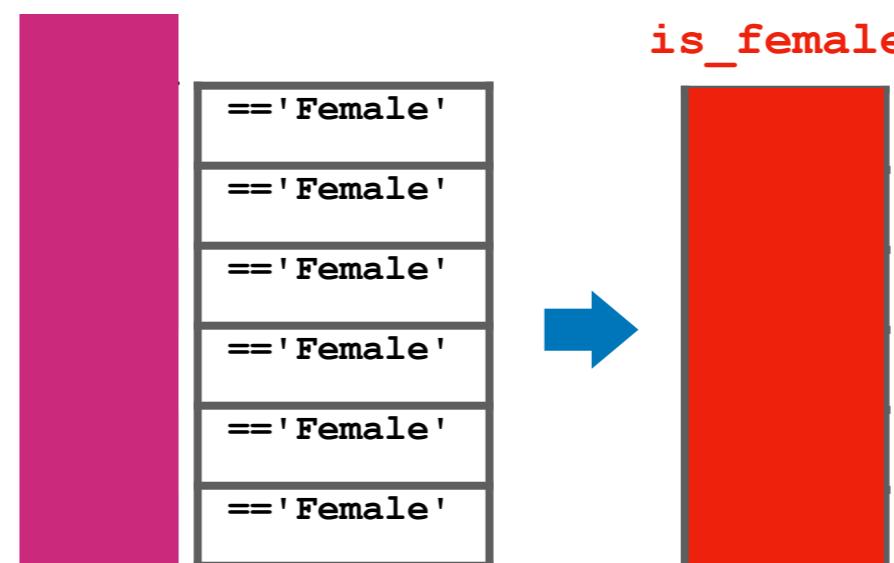
	wage	educ	exper	gender	married	remarks
0	3.10	11.0	2.0	Female	False	none
1	3.24	12.0	22.0	Female	True	none
2	3.00	1.0	1.0	Male	False	none
3	6.00	9.0	44.0	Male	True	none
4	5.30	12.0	7.0	Male	True	none
5	8.75	16.0	9.0	Male	True	none

# Pandas for Data Analysis

- The **Series** data structure
  - ▶ Boolean **Series** and boolean indexing

```
1 is_female = data_frame['gender'] == 'Female'  
2 is_female
```

```
0    True  
1    True  
2   False  
3   False  
4   False  
5   False  
Name: gender, dtype: bool
```

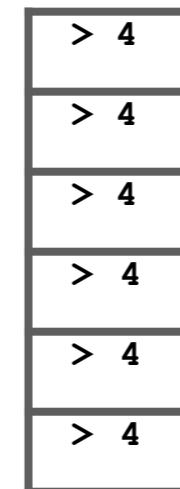


# Pandas for Data Analysis

- The **Boolean** data structure
  - ▶ Boolean **\_** and boolean indexing

```
1 is_high_wage = data_frame[ 'wage' ] > 4
2 is_high_wage
```

```
0    False
1    False
2    False
3    True
4    True
5    True
Name: wage, dtype: bool
```



**is\_high\_wage**



# Pandas for Data Analysis

- The **data structure**
  - ▶ Boolean **and boolean indexing**

```
3 is_wife = (data_frame['gender']=='Female') & (data_frame['married'])  
4 is_wife
```

```
0    False  
1    True  
2   False  
3   False  
4   False  
5   False  
dtype: bool
```



# Pandas for Data Analysis

- The **data structure**
  - ▶ Boolean **and boolean indexing**

```
1 is_female = data_frame['gender'] == 'Female'  
2 female = data_frame.loc[is_female]  
3 female
```

is_female	wage	educ	exper	gender	married	remarks
False	2 3.00	1.0	1.0	Male	False	none
False	3 6.00	9.0	44.0	Male	True	none
False	4 5.30	12.0	7.0	Male	True	none
False	5 8.75	16.0	9.0	Male	True	none

# Pandas for Data Analysis

- The **data structure**
  - ▶ Boolean **and boolean indexing**

```
1 female = data_frame['gender'] == 'Female'  
2 female = data_frame.loc[female]  
3 female
```

	wage	educ	exper	gender	married	remarks
0	3.10	11.0	2.0	Female	False	none
1	3.24	12.0	22.0	Female	True	none

	wage	educ	exper	gender	married	remarks
0	3.10	11.0	2.0	Female	False	none
1	3.24	12.0	22.0	Female	True	none
2	3.00	1.0	1.0	Male	False	none
3	6.00	9.0	44.0	Male	True	none
4	5.30	12.0	7.0	Male	True	none
5	8.75	16.0	9.0	Male	True	none

# Pandas for Data Analysis

- The **data structure**
  - ▶ Boolean **and boolean indexing**

```
1 is_wife = (data_frame['gender']=='Female') & (data_frame['married'])
2 data_frame.loc[is_wife, 'remarks'] = 'Wife'
3 data_frame
```

is_wife	wage	educ	exper	gender	married	
FALSE	0 3.10	11.0	2.0	Female	False	
FALSE	2 3.00	1.0	1.0	Male	False	none
FALSE	3 6.00	9.0	44.0	Male	True	none
FALSE	4 5.30	12.0	7.0	Male	True	none
FALSE	5 8.75	16.0	9.0	Male	True	none

# Pandas for Data Analysis

- The **data structure**
  - ▶ Boolean **and boolean indexing**

```
1 is_wife = (data_frame['gender'] == 'Female') & (data_frame['married'])
2 data_frame.loc[is_wife] = 'Wife'
3 data_frame
```

is_wife	wage	educ	exper	gender	married	
FALSE	0 3.10	11.0	2.0	Female	False	
FALSE	2 3.00	1.0	1.0	Male	False	none
FALSE	3 6.00	9.0	44.0	Male	True	none
FALSE	4 5.30	12.0	7.0	Male	True	none
FALSE	5 8.75	16.0	9.0	Male	True	none

# Pandas for Data Analysis

---

- The **pandas.DataFrame** data structure
  - Boolean **pandas.Series** and boolean indexing

## Question

Please change values of the column “remarks” by yourself according to the following rule:

- The value for married male is changed to the string 'Husband';
- The value for unmarried male or female is changed to the string 'Single'.

# Roadmap

---

- 1 Array Data Structure
- 2 Data Visualization
- 3 Pandas
- 4 **Descriptive Analytics in Python**

# Descriptive Analytics in Python

---

“You can, for example, never foretell what any one man will do, but you can say with precision what an average number will be up to. Individuals vary, but percentages remain constant”

—*Sherlock Holmes, A Sign of Four*



# Descriptive Analytics in Python

---

- Read Data from files

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
...	...	...	...	...	...
521	15.00	16.0	14.0	Female	True
522	2.27	10.0	2.0	Female	False
523	4.67	15.0	13.0	Male	True
524	11.56	16.0	5.0	Male	True
525	3.50	14.0	5.0	Female	False

526 rows × 5 columns

# Descriptive Analytics in Python

- Read Data from files

```
1 data = pd.read_csv('wage.csv') # Read data from a file "wage.csv"  
2 data.head(6) # Return the first six rows of data
```

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
5	8.75	16.0	9.0	Male	True

	wage	educ	exper	gender	married
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
...	...	...	...	...	...
521	15.00	16.0	14.0	Female	True
522	2.27	10.0	2.0	Female	False
523	4.67	15.0	13.0	Male	True
524	11.56	16.0	5.0	Male	True
525	3.50	14.0	5.0	Female	False

526 rows × 5 columns

# Descriptive Analytics in Python

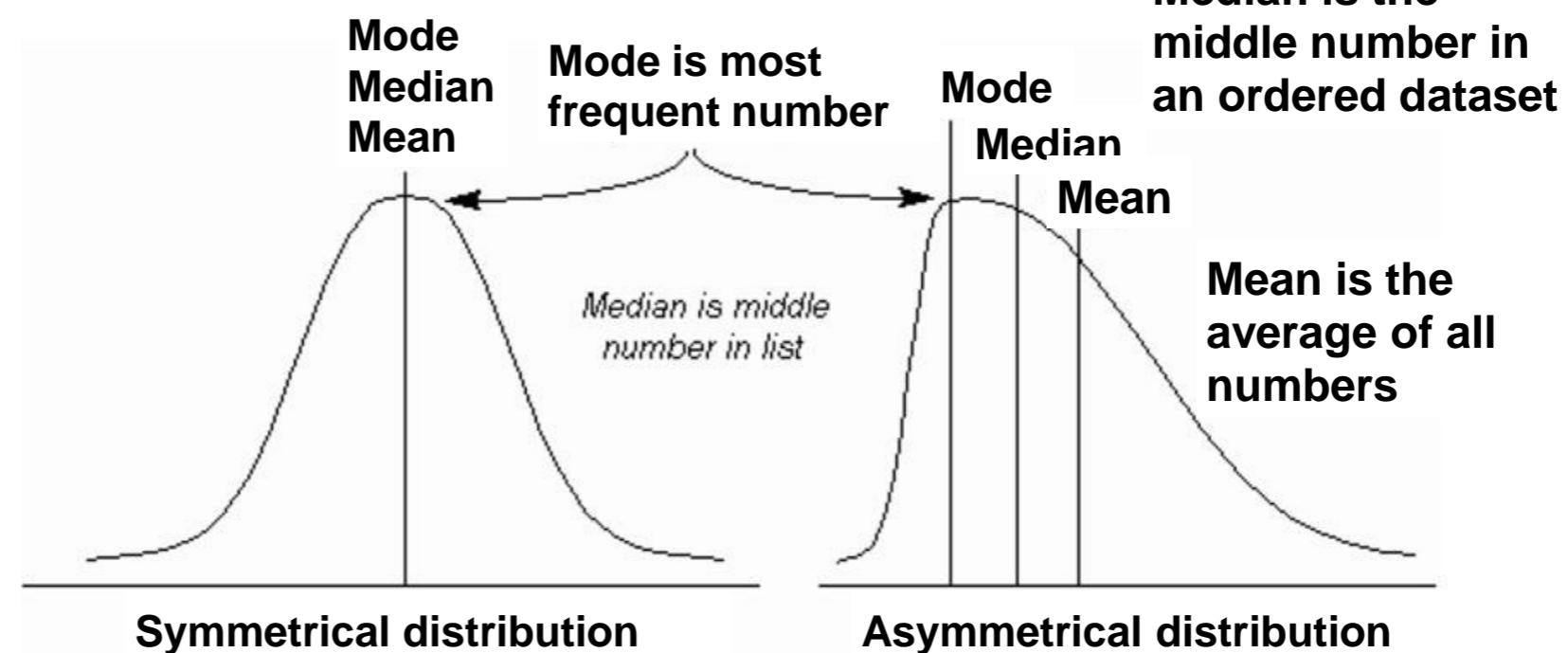
- Descriptive measures

- Measures of center

- ✓ Mean

- ✓ Median

- ✓ Mode



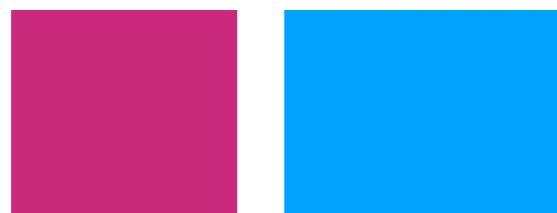
# Descriptive Analytics in Python

- Descriptive Measures

- Measures of center

- ✓ Mean

```
1 print(data.mean(), '\n')      # Mean value
2 print(type(data.mean())))    # Show the data type
```



dtype: float64

				gender	
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
...	...	...	...	...	...
521	15.00	16.0	14.0	Female	True
522	2.27	10.0	2.0	Female	False
523	4.67	15.0	13.0	Male	True
524	11.56	16.0	5.0	Male	True
525	3.50	14.0	5.0	Female	False

# Descriptive Analytics in Python

- Descriptive Measures

- Measures of center

- ✓ Mean

```
1 print(data.mean(), '\n')      # Mean value of each column
2 print(type(data.mean()))     # Show the data type of the results
```

```
wage      5.896103
educ      12.562738
exper     17.017110
dtype: float64
```

```
<class 'pandas.core.series.Series'>
```



$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Total case number

Number of married cases  
 $x_i = \begin{cases} 1, & \text{if True} \\ 0, & \text{if False} \end{cases}$

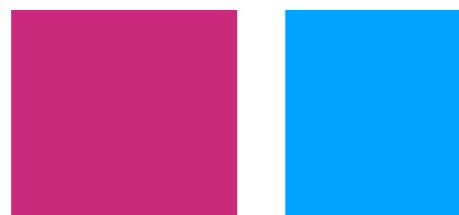
# Descriptive Analytics in Python

- Descriptive Measures

- Measures of center

- ✓ Median

```
1 print(data.median(), '\n')      # Median of all columns
2 print()                         # Show the data
```

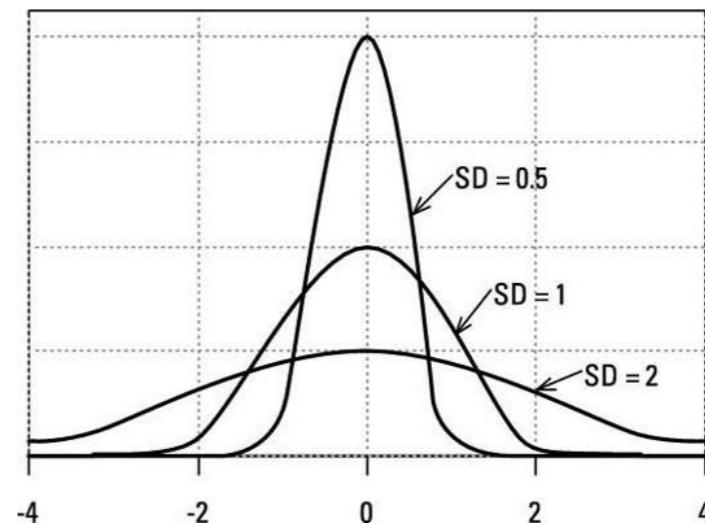


dtype: float64

				gender	
0	3.10	11.0	2.0	Female	False
1	3.24	12.0	22.0	Female	True
2	3.00	11.0	2.0	Male	False
3	6.00	8.0	44.0	Male	True
4	5.30	12.0	7.0	Male	True
...	...	...	...	...	...
521	15.00	16.0	14.0	Female	True
522	2.27	10.0	2.0	Female	False
523	4.67	15.0	13.0	Male	True
524	11.56	16.0	5.0	Male	True
525	3.50	14.0	5.0	Female	False

# Descriptive Analytics in Python

- Descriptive Measures
  - Measures of variation
    - ✓ Standard deviation
    - ✓ Variance



```
1 data.std()          # Sample Standard deviation of each column

wage      3.693086
educ      2.769022
exper     13.572160
married   0.488580
dtype: float64
```

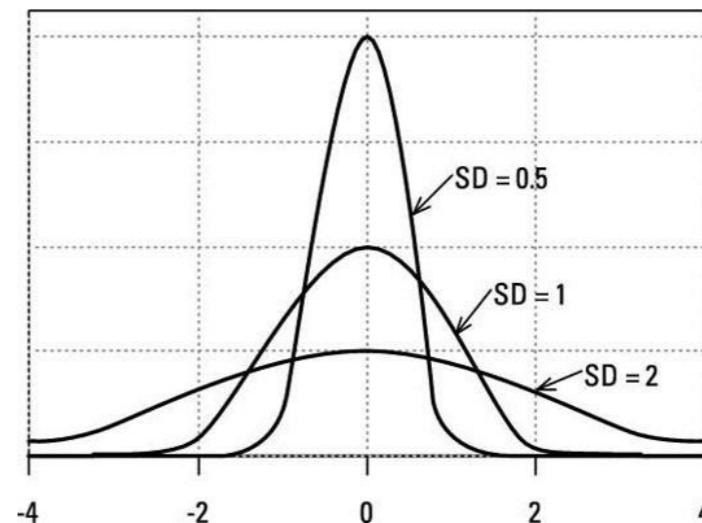
# Descriptive Analytics in Python

- Descriptive Measures

- Measures of variation

- ✓ Standard deviation

- ✓ Variance



```
1 data.var()          # Sample Variance of each column
```

```
wage      13.638884
educ      7.667485
exper     184.203516
married    0.238711
dtype: float64
```

# Descriptive Analytics in Python

---

- Descriptive Measures

- Extreme points

- ✓ Maximum**

```
1 data.max()          # Maximum value of each column
```

```
wage      24.98
educ      18
exper      51
gender    Male
married    True
dtype: object
```

- ✓ Minimum**

```
1 data.min()          # Minimum value of each column
```

```
wage      0.53
educ      0
exper      1
gender    Female
married   False
dtype: object
```

# Descriptive Analytics in Python

---

- Descriptive Measures
  - Counts of categorical values

```
1 data['gender'].value_counts()
```

```
Male      274
Female    252
Name: gender, dtype: int64
```

Actual count of each category (as default)

```
1 data['gender'].value_counts(normalize=True)
```

```
Male      0.520913
Female    0.479087
Name: gender, dtype: float64
```

Proportion of each category

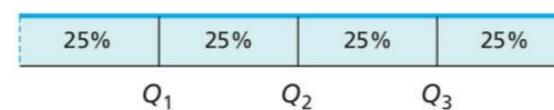
# Descriptive Analytics in Python

- Descriptive Measures

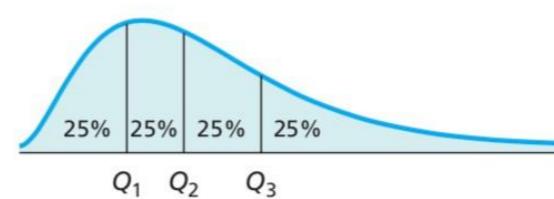
- Method `describe()`

```
1 wage_summary = data.describe()  
2 wage_summary
```

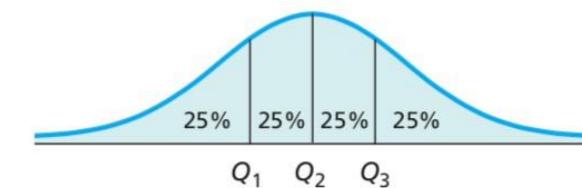
	wage	educ	exper
count	526.000000	526.000000	526.000000
mean	5.896103	12.562738	17.01711
std	3.693086	2.769022	13.57216
min	0.530000	0.000000	1.00000
max	24.980000	18.000000	51.00000



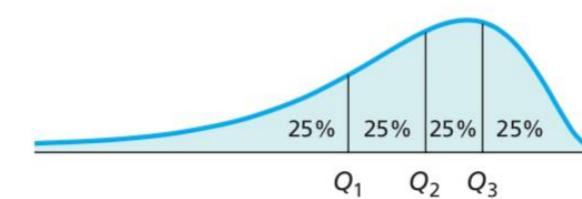
(a) Uniform



(c) Right skewed



(b) Bell shaped



(d) Left skewed

**Q1**      **Quartiles: values**

**Q2**      **that divide a dataset**

**Q3**      **into quarters:**

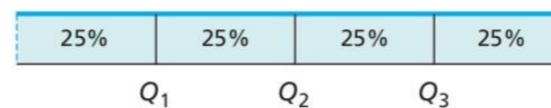
# Descriptive Analytics in Python

- Descriptive Measures

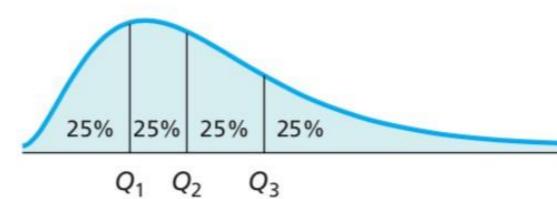
- Method

```
1 wage_summary = data.describe()  
2 wage_summary
```

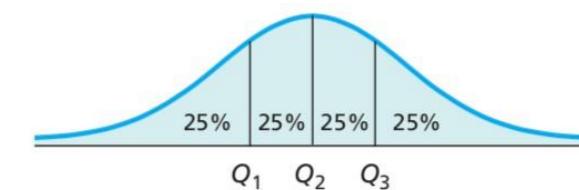
	wage	educ	exper
<b>count</b>	526.000000	526.000000	526.000000
<b>mean</b>	5.896103	12.562738	17.01711
<b>std</b>	3.693086	2.769022	13.57216
<b>min</b>	0.530000	0.000000	1.00000
<b>max</b>	24.980000	18.000000	51.00000



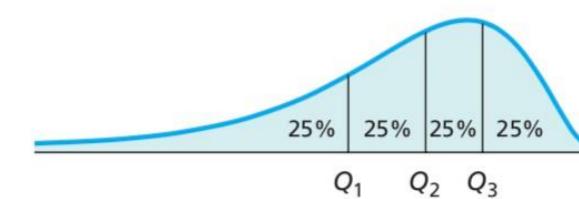
(a) Uniform



(c) Right skewed



(b) Bell shaped

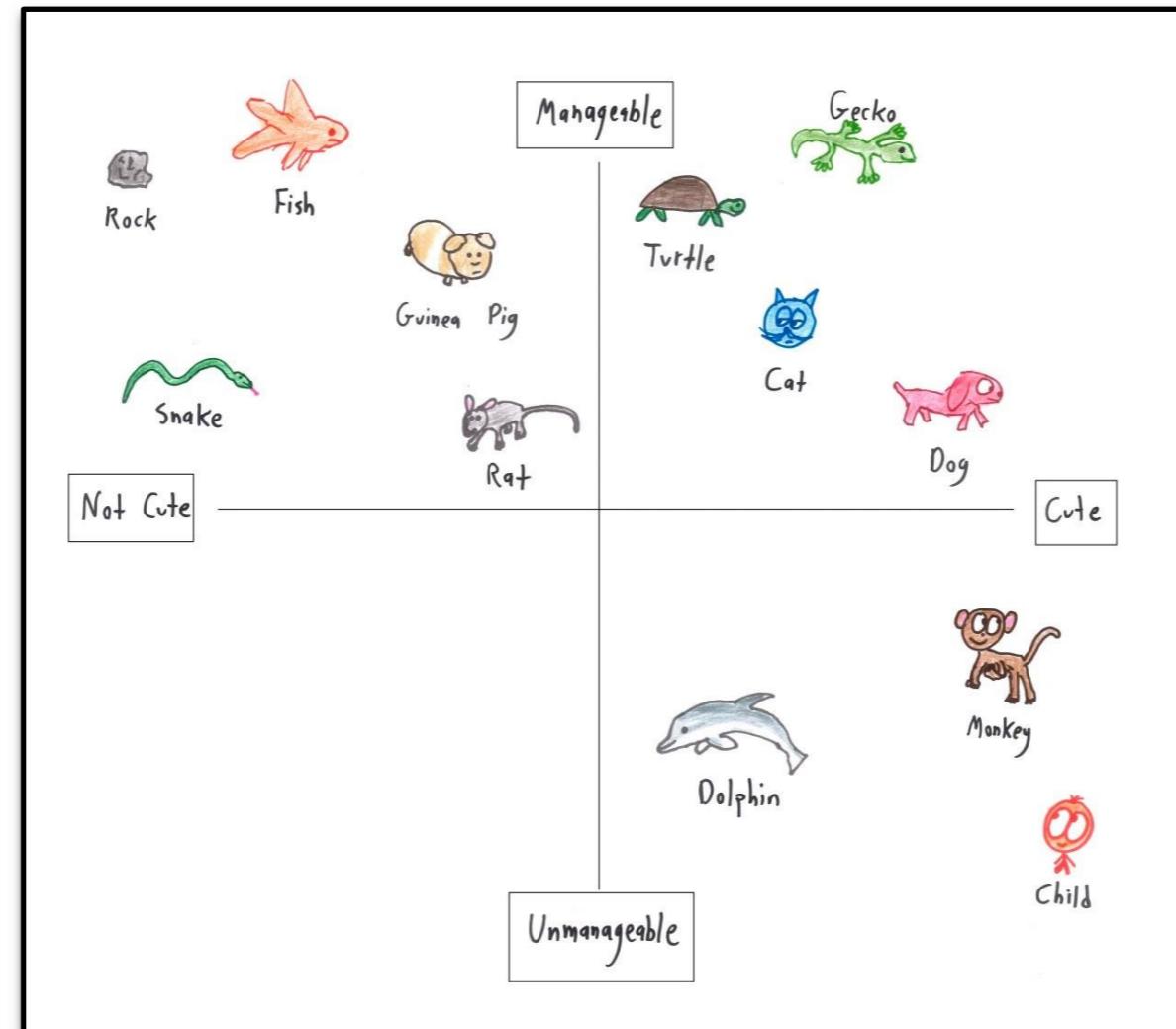


(d) Left skewed

The value **Q<sub>3</sub> - Q<sub>1</sub>** is called the interquartile range (IQR).

# Descriptive Analytics in Python

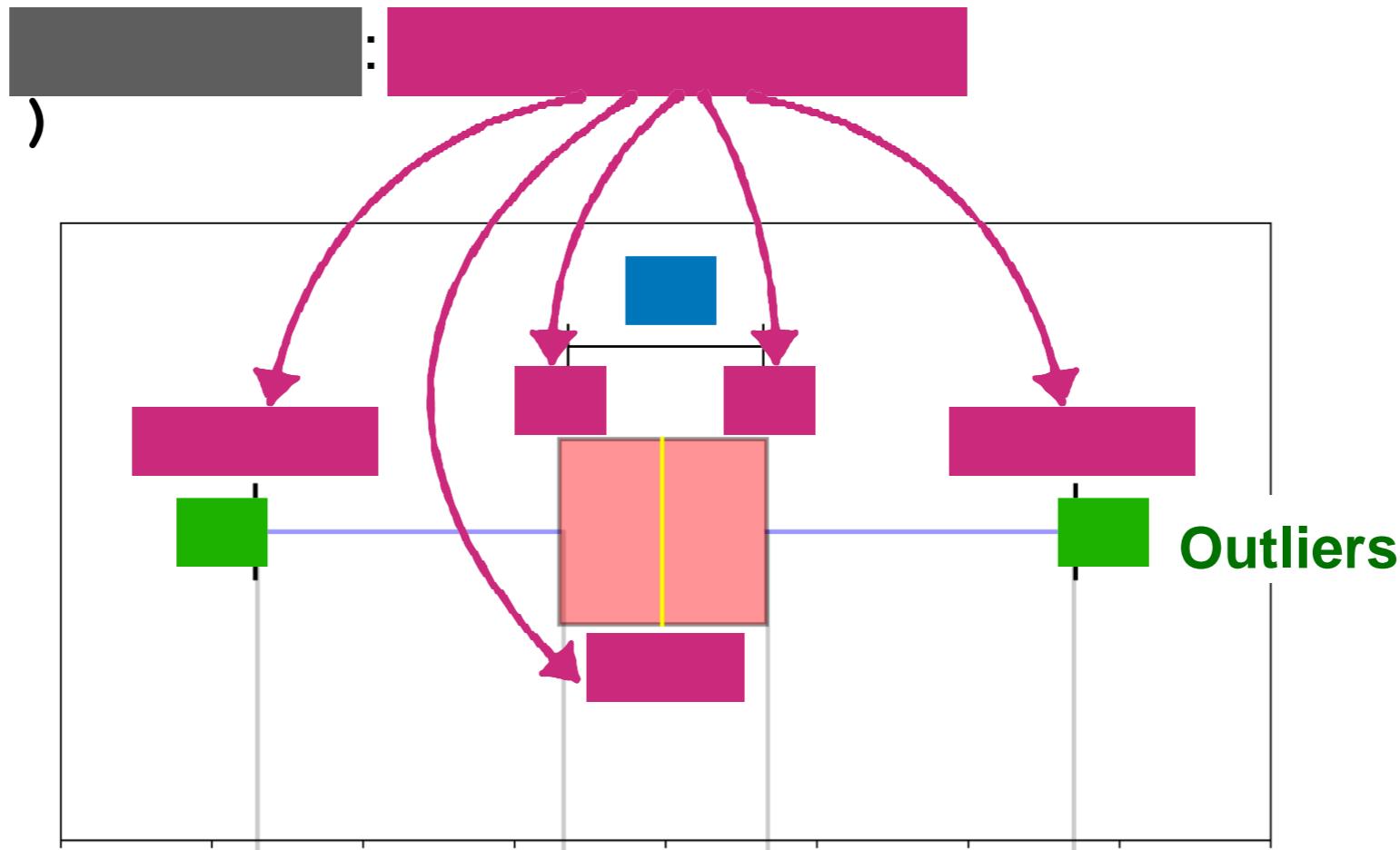
- Data Visualization



# Descriptive Analytics in Python

- Data Visualization
  - Distributions of variables

✓ Function



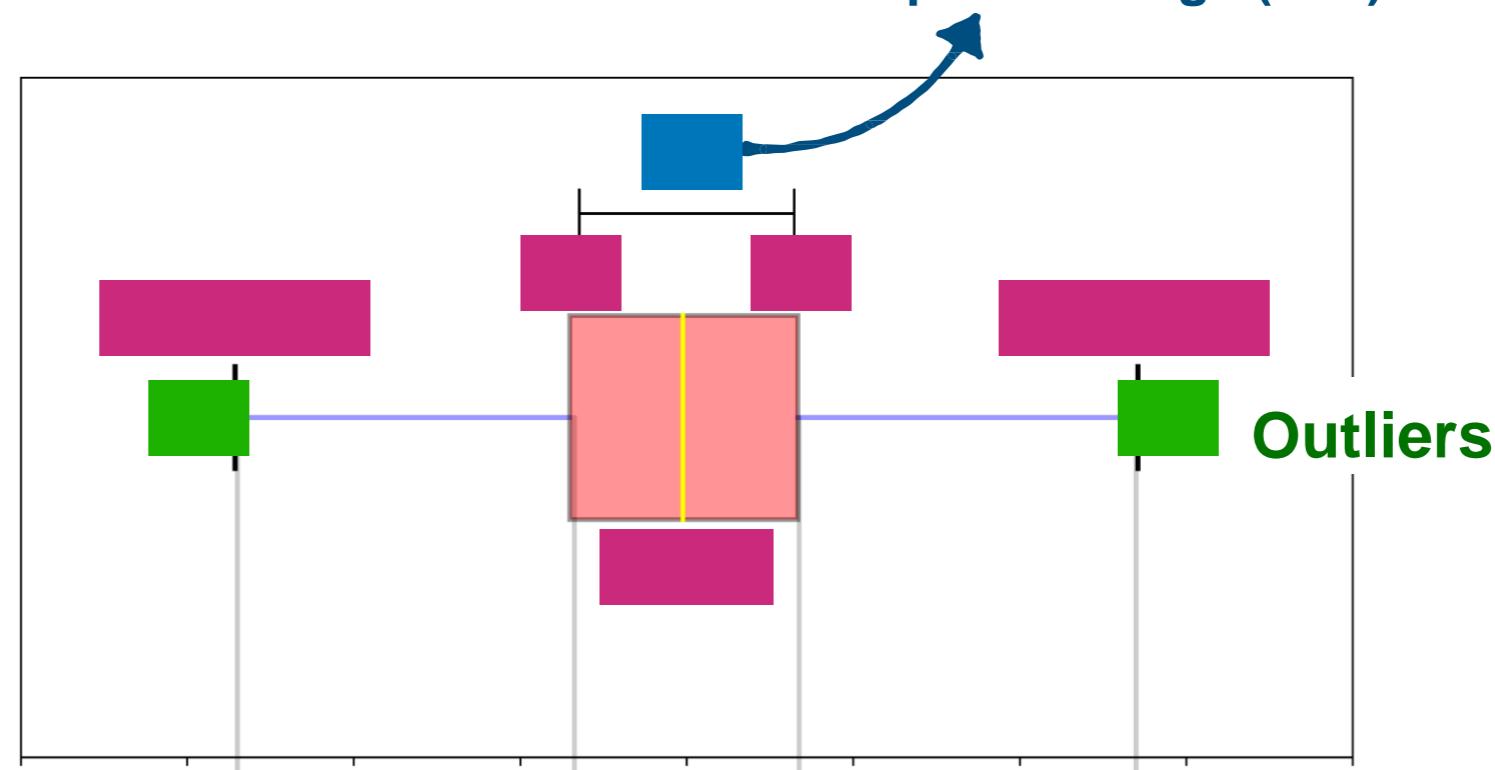
# Descriptive Analytics in Python

- Data Visualization
  - Distributions of variables

✓ Function

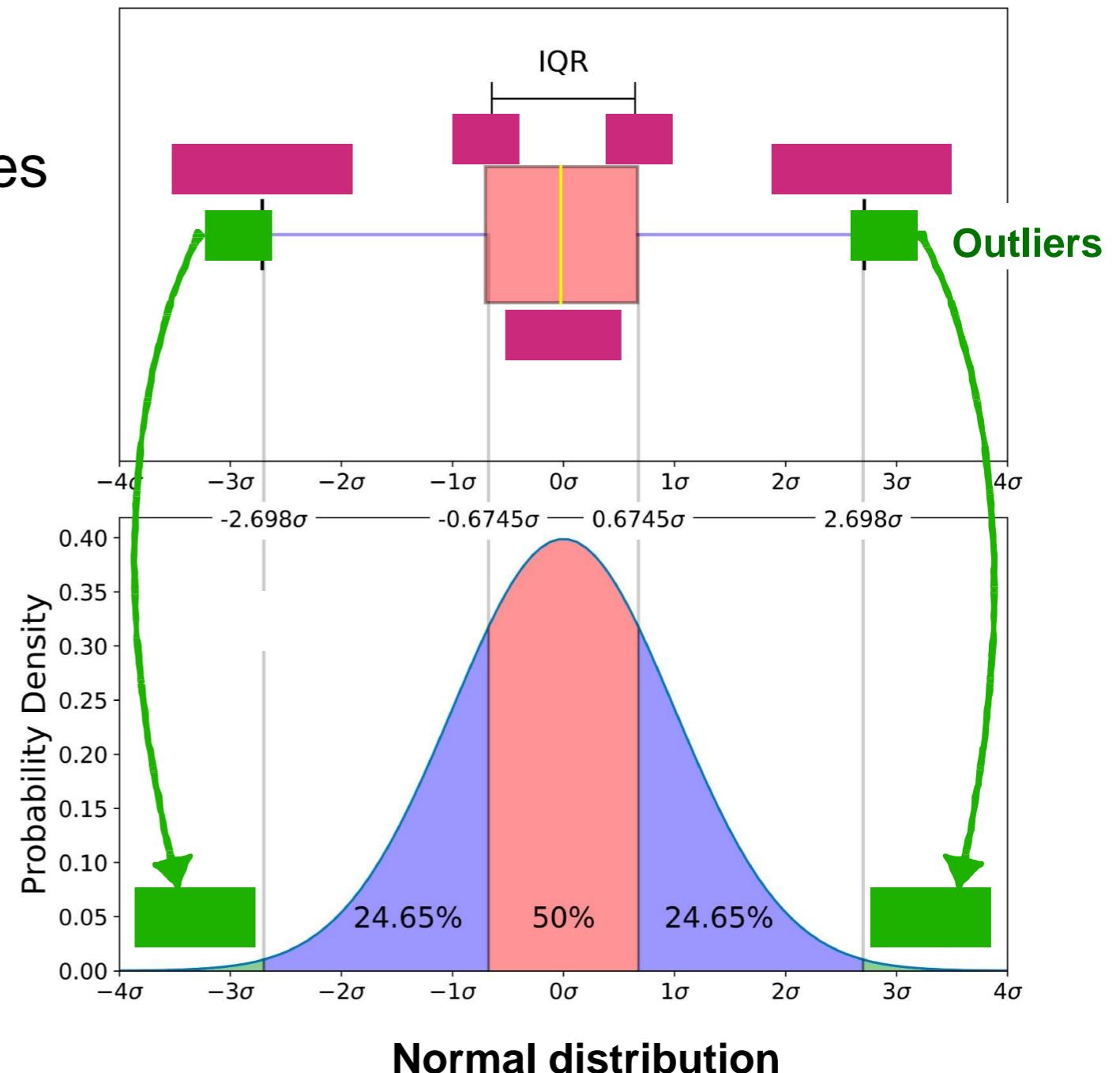
```
 )
```

The value  $Q3 - Q1$  is called  
the interquartile range (IQR).



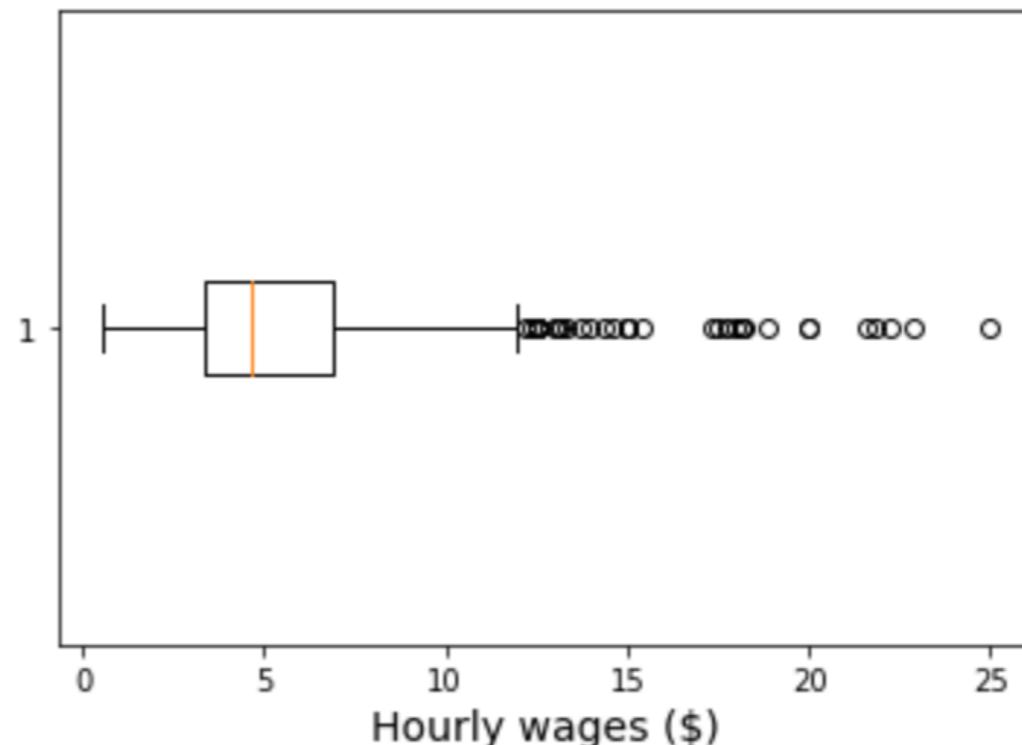
# Descriptive Analytics in Python

- Data Visualization
  - Distributions of variables
  - ✓ Function `boxplot()`



# Descriptive Analytics in Python

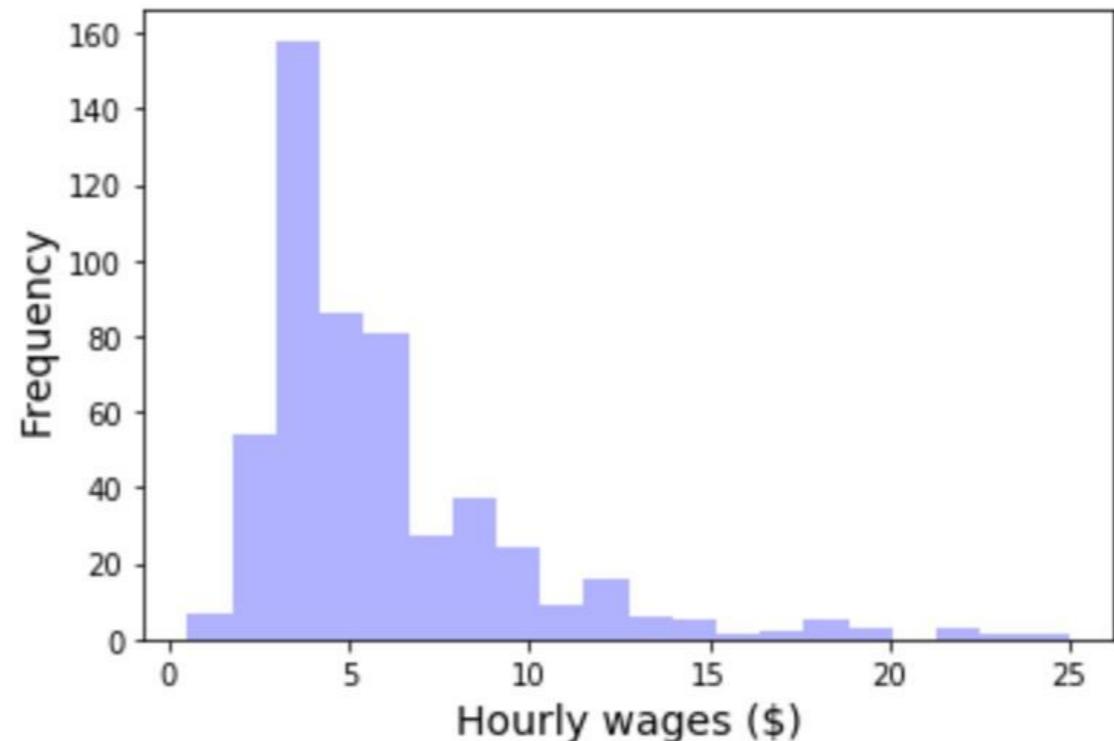
- Data Visualization
  - Distributions of variables
    - ✓ Package **matplotlib**
    - ✓ Function **boxplot()**



```
1 plt.boxplot(data['wage'],  
2               vert=False);           # Create a box plot of the wages  
3  
4 plt.xlabel('Hourly wages ($)',  
5            fontsize=14);           # Make the plot horizontal  
6  
7 plt.show()                         # Label of the x axis
```

# Descriptive Analytics in Python

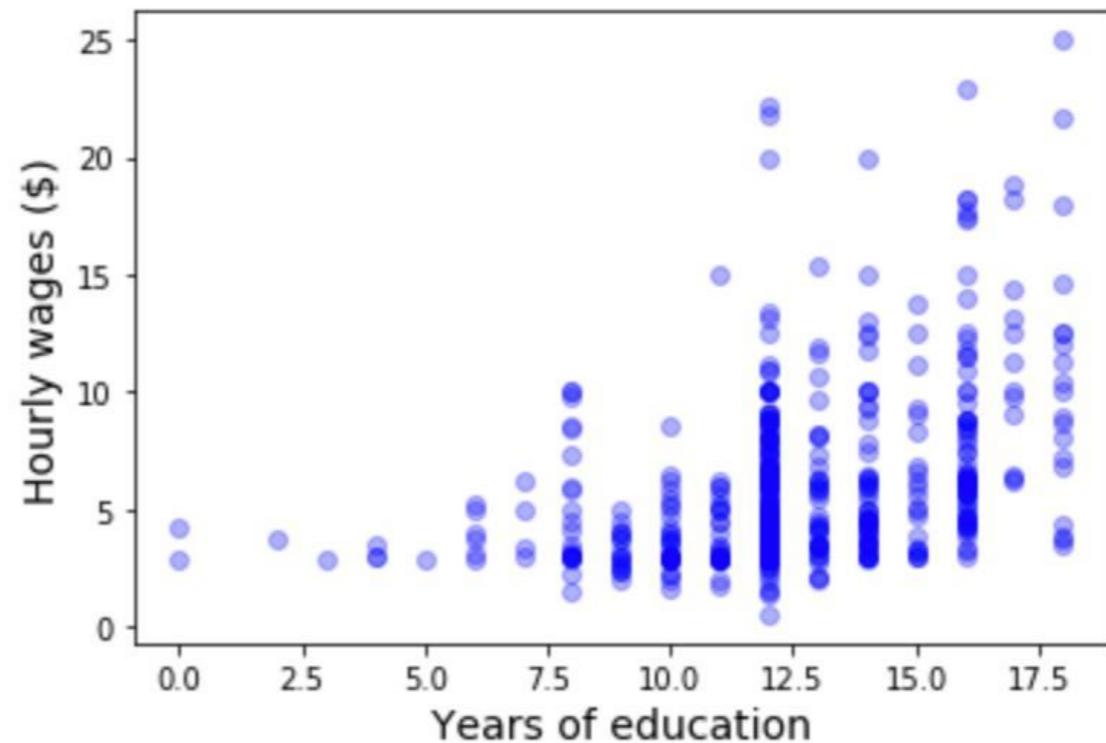
- Data Visualization
  - Distributions of variables
    - ✓ Package 
    - ✓ Function 



```
1 plt.hist(data['wage'], 20,      # Histogram of wages with 20 bins
2           color='b', alpha=0.3); # Color is blue, opacity is 0.3
3
4 plt.xlabel('Hourly wages ($)', fontsize=14)
5 plt.ylabel('Frequency', fontsize=14)
6
7 plt.show()
```

# Descriptive Analytics in Python

- Data Visualization
  - Relationship between variables



```
1 plt.scatter(data['educ'],          # x data is the years of education
2             data['wage'],          # y data is the hourly wages
3             c='b', alpha=0.3)      # Color is blue, opacity is 0.3
4
5 plt.xlabel('Years of education', fontsize=14)
6 plt.ylabel('Hourly wages ($)', fontsize=14)
7
8 plt.show()
```

# Descriptive Analytics in Python

- Data Visualization

- ▶ Relationship between variables

```
1 | data.corr()
```

	wage	educ	exper	married
wage	1.000000		0.112903	0.228817
educ		1.000000	-0.299542	0.068881
exper	0.112903	-0.299542	1.000000	0.316984
married	0.228817	0.068881	0.316984	1.000000

