

LẬP TRÌNH PYTHON CHO PHÂN TÍCH DỮ LIỆU VÀ HỌC MÁY

Bài 4: Xử lý dữ liệu số với NumPy - 01

AI Academy Vietnam

Nội dung bài 4

1. Một số thư viện quan trọng của Python

2. Giới thiệu thư viện Numpy

3. Cách tạo vector, ma trận (matrix)

➤ Tạo mảng 1D, 2D, 3D | Tạo mảng với các hàm có sẵn | Tạo mảng từ file dữ liệu

4. Thao tác cơ bản với mảng

➤ Quan sát thuộc tính | Chuyển đổi kiểu dữ liệu | Truy cập phần tử

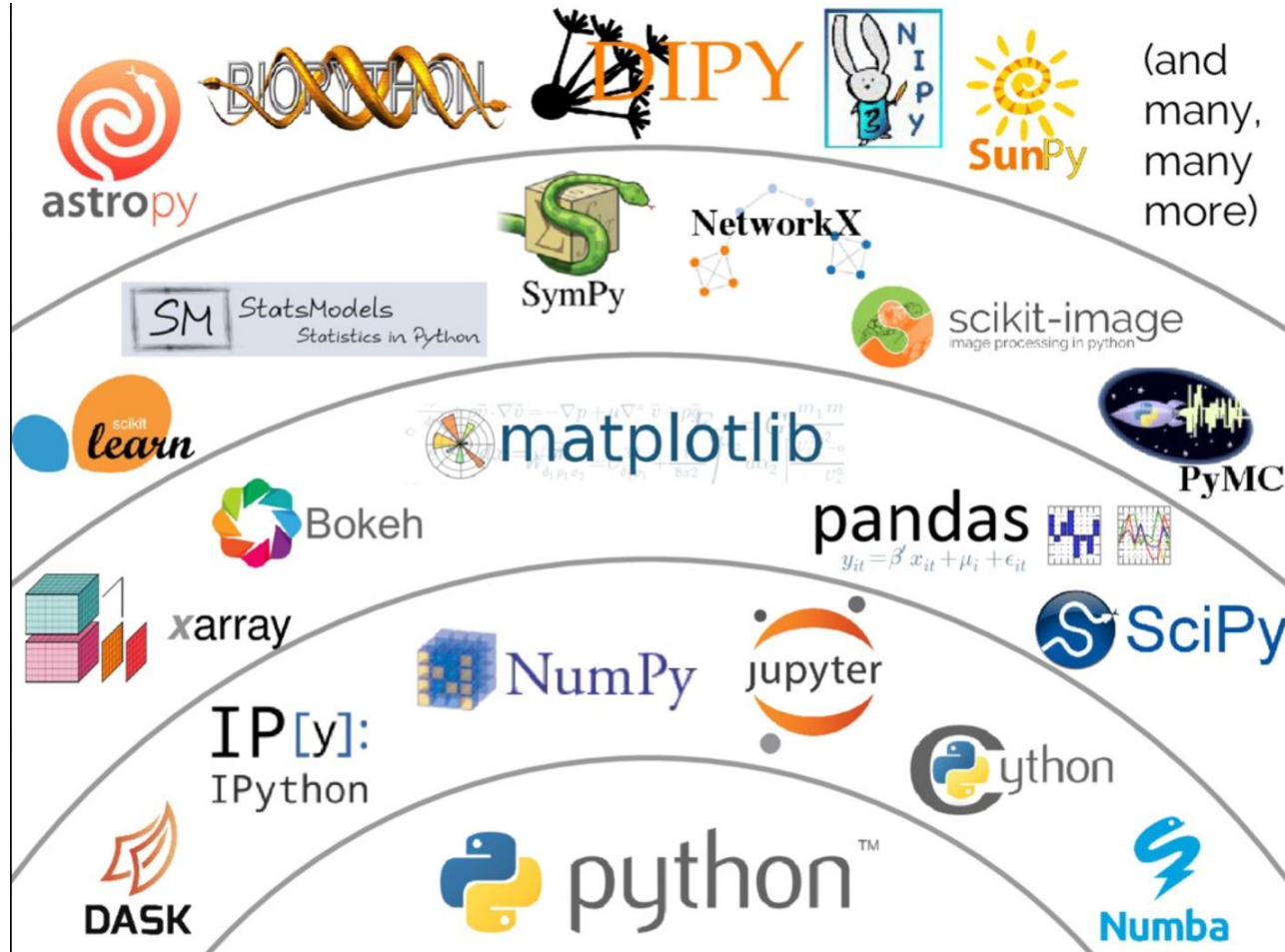
5. Tính toán các đặc trưng thống kê

➤ Min, Max, Mean, Median, Mode, Range, Std, Corrcoef

1. Một số thư viện quan trọng của Python

Thư viện Python

- **Python** có hệ thống thư viện rất phong phú, hỗ trợ nhiều lĩnh vực khác nhau.



- Do đó, tùy thuộc vào lĩnh vực nghiên cứu cụ thể, để lựa chọn và sử dụng các thư viện cho phù hợp.

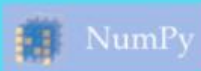
Thư viện Python

Top 5 Python Libraries for Data Science



Report from Cloud Academy suggests that the top technical skill in demand for data engineers is python. 67 percent of job posts mentioned python.

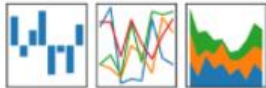
1.) NUMPY



Through NumPy, you can use it as an efficient multi-dimensional container of generic data. It also contains sophisticated (broadcasting) functions, tools for integrating C/C++ and Fortran code, useful linear algebra, Fourier transform, and random number capabilities.

2.) PANDAS

Pandas provides high-performance, easy-to-use data structures and data analysis tools for python. You can store and manage data from tables by performing manipulation over rows and columns.



3.) SCIKIT-LEARN



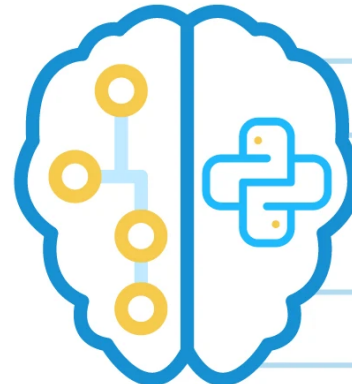
Scikit-Learn is a powerful library for machine learning in Python. It contains simple and efficient tools for data mining and data analysis. It is built on NumPy, SciPy and matplotlib.

Python libraries for Data Analysis

There are many interesting libraries that have made Python popular with Data Scientists:

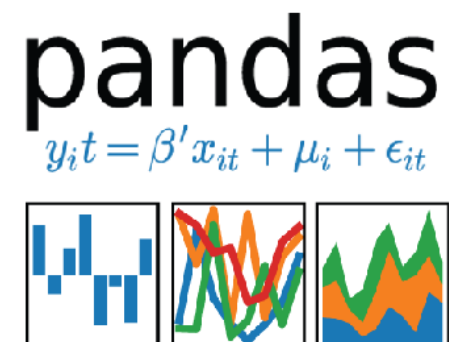
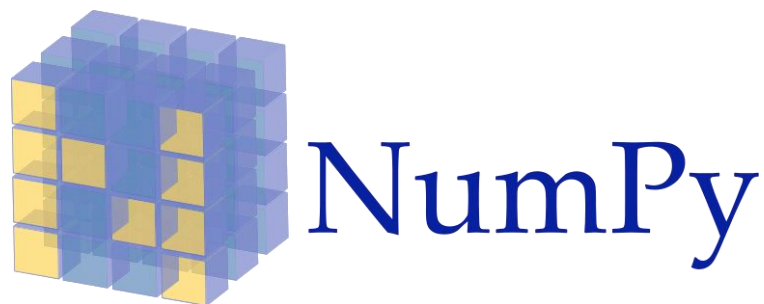


Python Libraries for Machine Learning

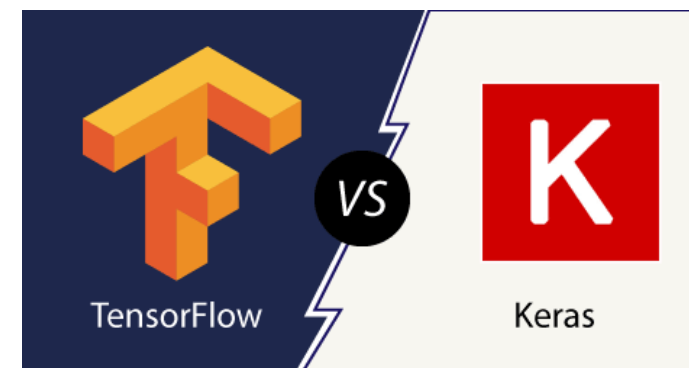
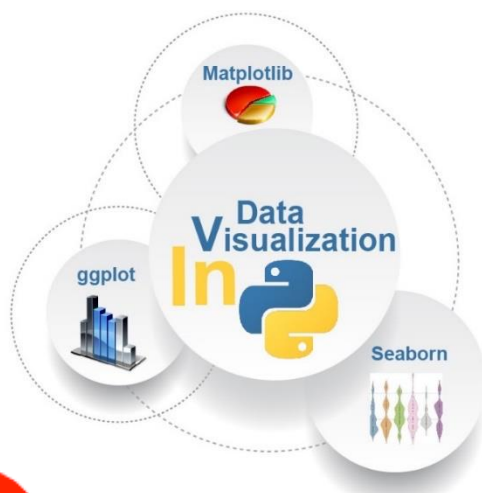


- | | |
|----------------|--------------|
| ● NumPy | ● Matplotlib |
| ● Pandas | ● Seaborn |
| ● Scikit-learn | ● OpenCV |
| ● Keras | ● spaCy |
| ● TensorFlow | ● SciPy |

Thư viện Python



	BandName	WavelengthMax	WavelengthMin
0	CoastalAerosol	450	430
1	Blue	510	450
2	Green	590	530
3	Red	670	640
4	NearInfrared	880	850
5	ShortWaveInfrared_1	1650	1570
6	ShortWaveInfrared_2	2290	2110
7	Cirrus	1380	1360



Natural Language
Analyses with NLTK

2. Thư viện NumPy

Thư viện Numpy

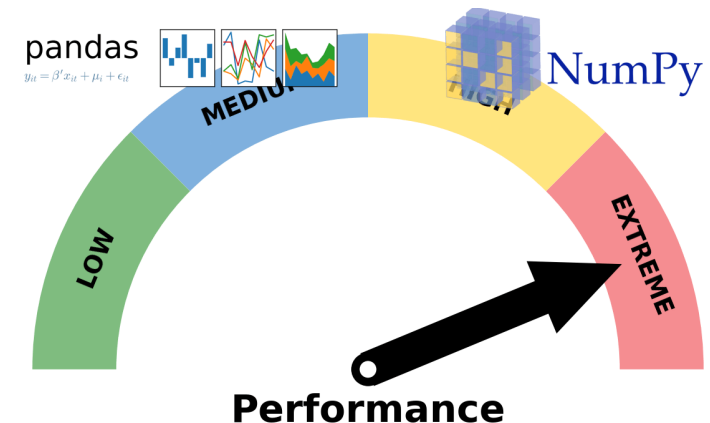
- **Numpy** (Numeric Python): là một thư viện toán học phổ biến và mạnh mẽ của Python.
- Cho phép làm việc hiệu quả với ma trận và mảng, đặc biệt là dữ liệu ma trận và mảng lớn với tốc độ xử lý nhanh hơn nhiều lần khi chỉ sử dụng “core Python” đơn thuần.
- Ngoài ra, Python cũng hỗ trợ một thư viện khác để mở rộng thêm các tính năng của Numpy là Scipy với ưu thế về các phép hồi quy hay biến đổi Fourier...
- Tham khảo thêm tại: <http://www.numpy.org/>

```
1 big_array = np.random.rand(1000000)
2 %timeit sum(big_array)
3 %timeit np.sum(big_array)
```

10 loops, best of 3: 171 ms per loop
1000 loops, best of 3: 380 µs per loop

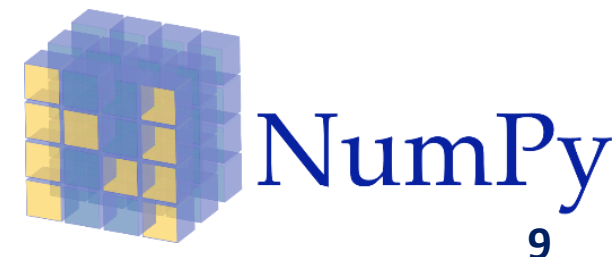
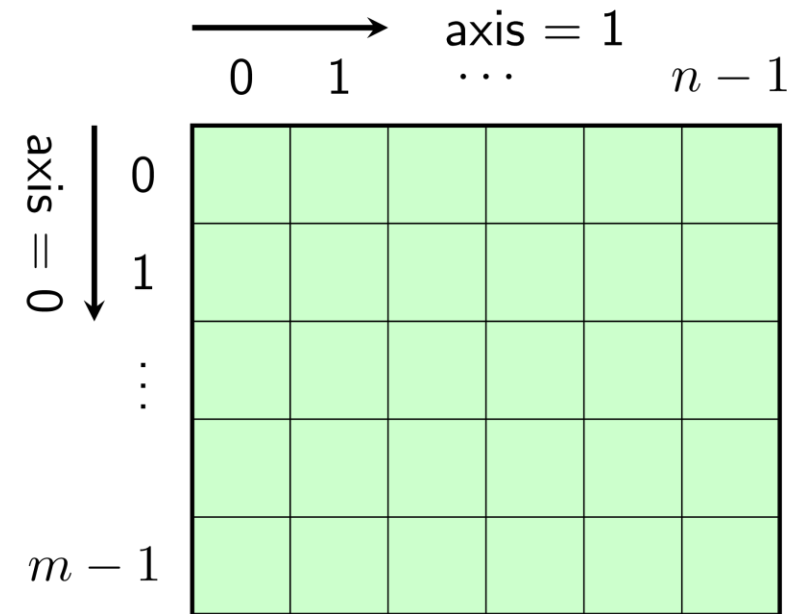
```
1 %timeit min(big_array)
2 %timeit np.min(big_array)
```

10 loops, best of 3: 103 ms per loop
1000 loops, best of 3: 432 µs per loop



Thư viện Numpy

- **Đối tượng chính của NumPy** là các mảng đa chiều đồng nhất:
 - Kiểu dữ liệu của các phần tử con trong mảng phải **giống nhau**
 - Mảng có thể có 1 chiều hoặc nhiều chiều
 - Các chiều được đánh số từ 0 trở đi
 - Số chiều được gọi là hạng (**rank**)
 - Có đến 24 kiểu số khác nhau.
 - Kiểu ndarray là lớp chính xử lý dữ liệu mảng nhiều chiều.
 - Có rất nhiều hàm và phương thức xử lý mảng



Thư viện Numpy

1	5	18	23
---	---	----	----

Vector (1D array)
Dimension = 1
(1 index required)

3	12	66
7	9	34
23	45	11

Matrix (2D array)
Dimension = 2
(2 indexes required)

3	12	66
7	9	34
23	45	11

3D array (3rd order Tensor)
Dimension = 3
(3 indexes required)

3	12	66
7	9	34
23	45	11

 ...

3	12	66
7	9	34
23	45	11

ND array
Dimension = N
(N indexes required)

3. Khởi tạo mảng

Khởi tạo mảng

Khởi tạo mảng 1 chiều – 1D (Vector)

```
1  #Khởi tạo mảng 1 chiều với thư viện Numpy
2  import numpy as np
3
4  #Tạo mảng 1 chiều (1D) - row
5  a = np.array([1, 2, 5, 7, 0, 8])
6
7  print(a)
8  print("Loại dữ liệu của biến a:", type(a))
9  print("Kiểu dữ liệu của phần tử trong mảng a:", a.dtype)
10 print("Kích thước của mảng a:", a.shape)
11 print("Số phần tử của mảng a:", a.size)
12 print("Số chiều của mảng a:", a.ndim)
```

[1 2 5 7 0 8]

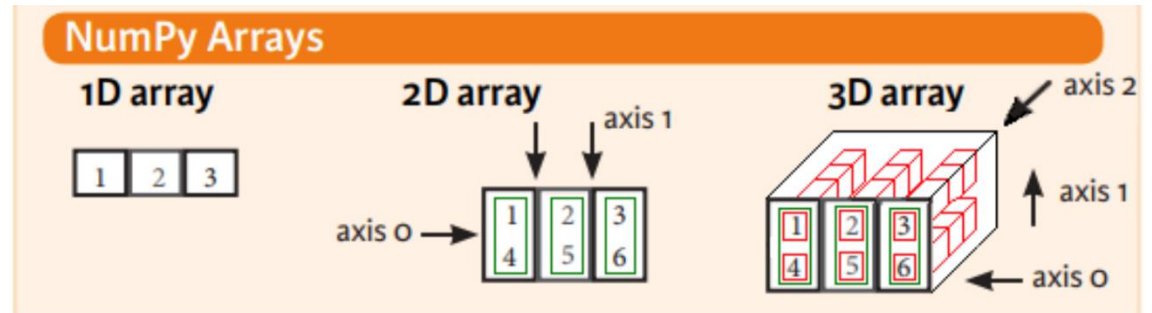
Loại dữ liệu của biến a: <class 'numpy.ndarray'>

Kiểu dữ liệu của phần tử trong mảng a: int32

Kích thước của mảng a: (6,)

Số phần tử của mảng a: 6

Số chiều của mảng a: 1

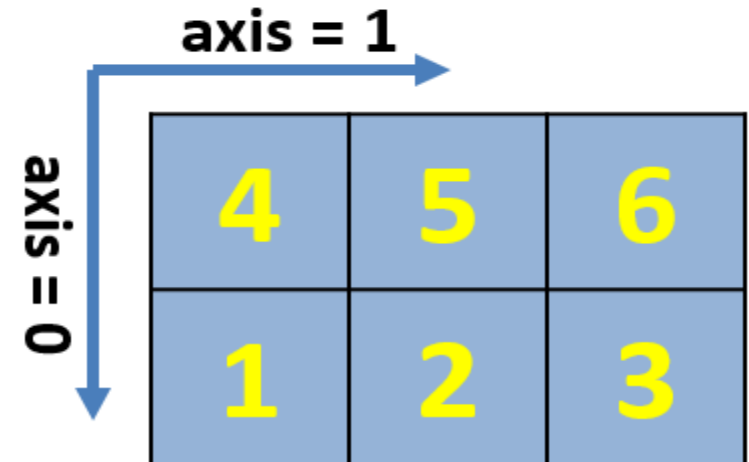


Khởi tạo mảng (2)

- Khởi tạo mảng 2 chiều – 2D (Matrix)

```
1 #Gọi thư viện numpy
2 import numpy as np
3
4 #Tạo mảng 2 chiều (2D - Ma trận)
5 b = np.array([(4, 5, 6.0),(1, 2, 3.5)])
6
7 print(b)
8 print("Loại dữ liệu của biến b:", type(b))
9 print("Kiểu dữ liệu của phần tử trong mảng b:", b.dtype)
10 print("Kích thước của mảng b:", b.shape)
11 print("Số phần tử của mảng b:", b.size)
12 print("Số chiều của mảng b:", b.ndim)
```

```
[[4.  5.  6. ]
 [1.  2.  3.5]]
Loại dữ liệu của biến b: <class 'numpy.ndarray'>
Kiểu dữ liệu của phần tử trong mảng b: float64
Kích thước của mảng b: (2, 3)
Số phần tử của mảng b: 6
Số chiều của mảng b: 2
```



Khởi tạo mảng (3)

- Khởi tạo mảng 3 chiều – 3D

```
1 import numpy as np
2
3 c = np.array([[ (2,4,0,6), (4,7,5,6)],
4               [ (0,3,2,1), (9,4,5,6)],
5               [ (5,8,6,4), (1,4,6,8) ]]) #mảng 3 chiều (3D)
6
7 print(c)
8 print("Phần tử đầu tiên của mảng c:",c[0,0,0])
9 print("Kiểu dữ liệu của phần tử trong mảng c:",c.dtype)
10 print("Kích thước của mảng c:",c.shape)
11 print("Số phần tử của mảng c:",c.size)
12 print("Số chiều của mảng c:",c.ndim)
```

```
[[[2 4 0 6]
  [4 7 5 6]]
```

```
[[0 3 2 1]
 [9 4 5 6]]
```

```
[[5 8 6 4]
 [1 4 6 8]]]
```

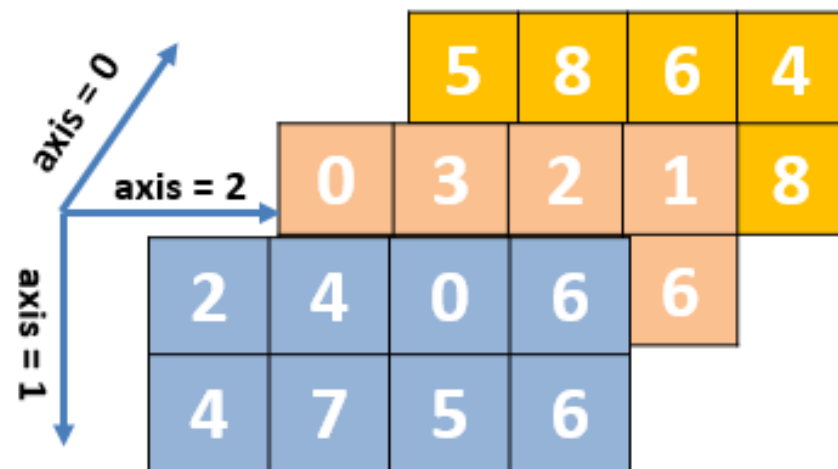
Phần tử đầu tiên của mảng c: 2

Kiểu dữ liệu của phần tử trong mảng c: int32

Kích thước của mảng c: (3, 2, 4)

Số phần tử của mảng c: 24

Số chiều của mảng c: 3



Khởi tạo mảng

Với các hàm có sẵn của NumPy

Khởi tạo mảng (4)

- Khởi tạo mảng với các hàm sẵn có của Numpy

Initial Placeholders

```
>>> np.zeros((3,4))
>>> np.ones((2,3,4),dtype=np.int16)
>>> d = np.arange(10,25,5)

>>> np.linspace(0,2,9)

>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2X2 identity matrix
Create an array with random values
Create an empty array

Khởi tạo mảng (5)

- Vd 1: Tạo ma trận 0/1 kích thước m x n

```
1 # Phương thức zeros: Tạo ma trận 0 kích thước 5 hàng x 3 cột
2 import numpy as np
3
4 array_zeros = np.zeros((5, 3))
5
6 print(array_zeros)
7 print("Kiểu dữ liệu trong mảng array_zeros:", array_zeros.dtype)
8 print("Kích thước của mảng array_zeros:", array_zeros.shape)
9 print("Số phần tử của mảng array_zeros:", array_zeros.size)
10 print("Số chiều của mảng array_zeros:", array_zeros.ndim)
```

```
[[0. 0. 0.]
 [0. 0. 0.]
 [0. 0. 0.]
 [0. 0. 0.]
 [0. 0. 0.]]
```

Kiểu dữ liệu trong mảng array_zeros: float64

Kích thước của mảng array_zeros: (5, 3)

Số phần tử của mảng array_zeros: 15

Số chiều của mảng array_zeros: 2

```
1 # Phương thức ones: Tạo ma trận 1 kích thước 3 hàng x 5 cột
2 import numpy as np
3
4 array_one = np.ones((3, 5), dtype=np.int)
5
6 print(array_one)
7 print("Kiểu dữ liệu trong mảng array_one:", array_one.dtype)
8 print("Kích thước của mảng array_one:", array_one.shape)
9 print("Số phần tử của mảng array_one:", array_one.size)
10 print("Số chiều của mảng array_one:", array_one.ndim)
```

```
[[1 1 1 1 1]
 [1 1 1 1 1]
 [1 1 1 1 1]]
```

Kiểu dữ liệu trong mảng array_one: int32

Kích thước của mảng array_one: (3, 5)

Số phần tử của mảng array_one: 15

Số chiều của mảng array_one: 2

Khởi tạo mảng (6)

- Vd 2: Tạo ma trận đơn vị cấp n

```
1 #Phương thức eye: Tạo ma trận đơn vị cấp 5
2 import numpy as np
3 array_eye = np.eye(5)
4
5 print(array_eye)
6 print("Kiểu dữ liệu của phần tử trong mảng array_eye:", array_eye.dtype)
7 print("Kích thước của mảng array_eye:", array_eye.shape)
8 print("Số phần tử của mảng array_eye:", array_eye.size)
9 print("Số chiều của mảng array_eye:", array_eye.ndim)
```

```
[[1. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0.]
 [0. 0. 1. 0. 0.]
 [0. 0. 0. 1. 0.]
 [0. 0. 0. 0. 1.]]
```

Kiểu dữ liệu của phần tử trong mảng array_eye: float64

Kích thước của mảng array_eye: (5, 5)

Số phần tử của mảng array_eye: 25

Số chiều của mảng array_eye: 2

Khởi tạo mảng (7)

- Vd 3: Tạo ma trận với các phần tử ngẫu nhiên trong khoảng (0,1)

```
1 #Phương thức random: Tạo một ma trận (7x5) các phần tử ngẫu nhiên [0,1]
2 import numpy as np
3 array_random = np.random.random((7,5))
4
5 print(array_random)
6 print("Kiểu dữ liệu của phần tử trong mảng array_random:", array_random.dtype)
7 print("Kích thước của mảng array_random:", array_random.shape)
8 print("Số phần tử của mảng array_random:", array_random.size)
9 print("Số chiều của mảng array_random:", array_random.ndim)
```

```
[[0.57738653 0.38330643 0.84085595 0.88920867 0.11759141]
 [0.13200344 0.40891213 0.46518628 0.81332657 0.62117097]
 [0.0255157  0.10842881 0.36001561 0.06382023 0.40403947]
 [0.37338483 0.35678386 0.38280971 0.97395415 0.8950108 ]
 [0.86054013 0.93742679 0.13039088 0.599897   0.0071806 ]
 [0.83131566 0.35010989 0.47524521 0.56107776 0.13418245]
 [0.39089618 0.40229334 0.73431802 0.53481456 0.45046071]]
```

Kiểu dữ liệu của phần tử trong mảng array_random: float64

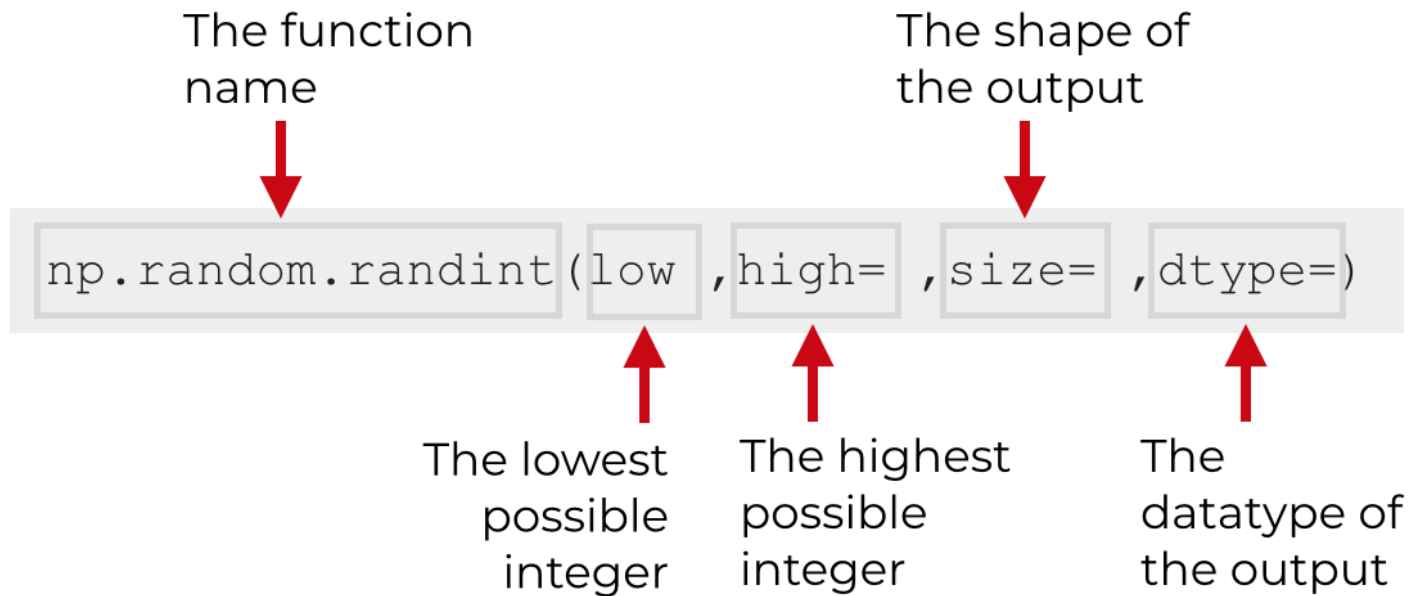
Kích thước của mảng array_random: (7, 5)

Số phần tử của mảng array_random: 35

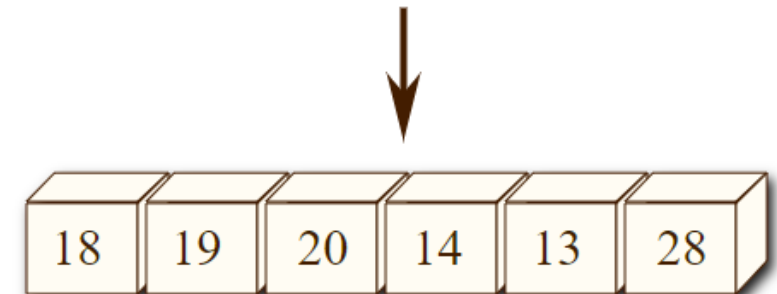
Số chiều của mảng array_random: 2

Khởi tạo mảng (8)

- Vd 4: Tạo vector, ma trận với các phần tử là số nguyên ngẫu nhiên trong khoảng (low,high)



`np.random.randint(low=10, high=30, size=6)`



Khởi tạo mảng (9)

- Vd 5: Tạo vector với các tham số thiết lập

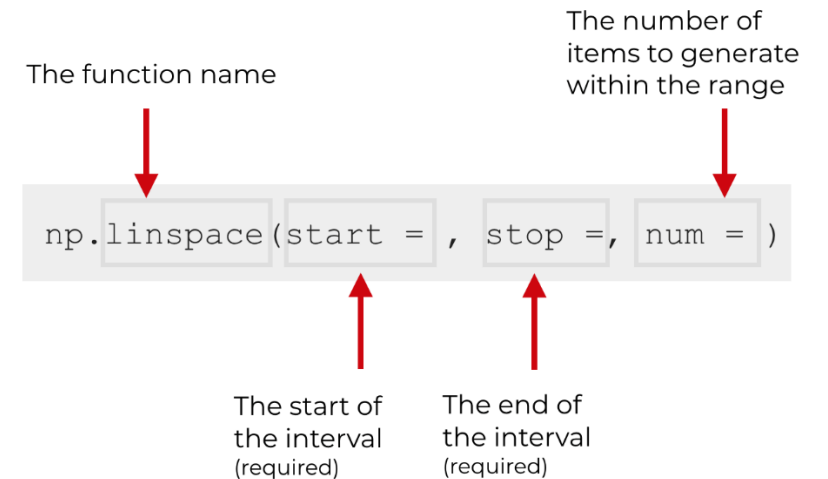
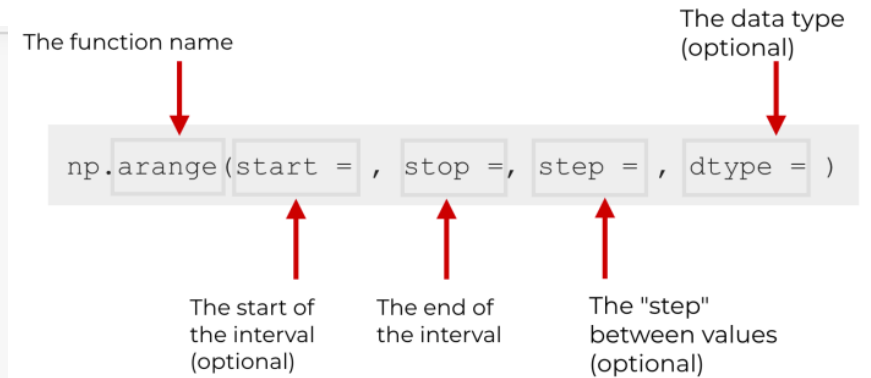
```
1 #Phương thức arange(a, b, steps):
2 #Tạo vector:
3 # Phần tử đầu tiên = a,
4 # kết thúc <b,
5 # mỗi phần tử cách nhau một khoảng = steps
6 d = np.arange(1, 15, 2)
7 print('Vector d:', d)
8 print("Số phần tử của vector d:", d.size)
9
10 print('-----')
11 #Phương thức linspace(a, b, num)
12 #Tạo vector:
13 #Phần tử đầu tiên = a,
14 #Phần tử kết thúc = b,
15 #Số phần tử của ma trận = num
16 f = np.linspace(1,15,11)
17 print('Vector f:', f)
18 print("Số phần tử của vector f:", f.size)
19
```

Vector d: [1 3 5 7 9 11 13]

Số phần tử của vector d: 7

Vector f: [1. 2.4 3.8 5.2 6.6 8. 9.4 10.8 12.2 13.6 15.]

Số phần tử của vector f: 11



Khởi tạo mạng

từ file dữ liệu .txt

Khởi tạo mảng (10)

- Bảng điểm của lớp 2A** (bao gồm 20 học sinh, tương ứng với 20 cột, của 10 môn học, tương ứng với 10 hàng) lưu trong file Diem_2A.txt

Diem_2A.txt

2,	4,	3,	7,	5,	6,	5,	6,	8,	9,	3,	6,	1,	9,	8,	7,	3,	3,	9,	5
3,	5,	3,	10,	9,	1,	9,	8,	3,	1,	6,	0,	7,	10,	8,	5,	2,	7,	7,	1
1,	10,	4,	9,	6,	9,	0,	2,	3,	1,	8,	6,	8,	4,	2,	9,	2,	9,	5,	0
6,	3,	0,	8,	3,	7,	7,	2,	6,	8,	7,	3,	4,	1,	5,	9,	1,	0,	2,	10
4,	3,	6,	7,	4,	5,	2,	6,	9,	4,	3,	9,	9,	4,	5,	7,	2,	10,	9,	4
2,	3,	8,	10,	4,	5,	9,	5,	4,	7,	10,	1,	8,	4,	3,	9,	6,	3,	6,	7
9,	9,	1,	10,	9,	9,	5,	9,	6,	3,	9,	5,	1,	10,	7,	10,	2,	8,	8,	1
8,	8,	7,	8,	6,	7,	7,	8,	6,	7,	8,	6,	7,	6,	8,	8,	7,	6,	8,	8
6,	7,	8,	9,	10,	9,	2,	2,	6,	1,	10,	9,	6,	3,	9,	5,	9,	8,	1,	1
7,	8,	7,	8,	6,	10,	10,	6,	8,	10,	8,	9,	8,	8,	5,	10,	8,	7,	8,	7

Điểm của môn học i
(10 môn học)

Điểm của học sinh j
(20 học sinh)

Khởi tạo mảng (11)

- Đọc dữ liệu từ file txt vào biến mảng.

```
#Đọc dữ liệu từ file Diem_2A.txt
diem_2a = np.loadtxt('Data/Diem_2A.txt',
                    delimiter=',',
                    dtype=np.int8)

print(diem_2a)
print("Kiểu dữ liệu của phần tử trong mảng diem_2a:", diem_2a.dtype)
print("Kích thước của mảng diem_2a:", diem_2a.shape)
print("Số phần tử của mảng diem_2a:", diem_2a.size)
print("Số chiều của mảng diem_2a:", diem_2a.ndim)
```

```
[[ 2  4  3  7  5  6  5  6  8  9  3  6  1  9  8  7  3  3  9  5]
 [ 3  5  3 10  9  1  9  8  3  1  6  0  7 10  8  5  2  7  7  1]
 [ 1 10  4  9  6  9  0  2  3  1  8  6  8  4  2  9  2  9  5  0]
 [ 6  3  0  8  3  7  7  2  6  8  7  3  4  1  5  9  1  0  2 10]
 [ 4  3  6  7  4  5  2  6  9  4  3  9  9  4  5  7  2 10  9  4]
 [ 2  3  8 10  4  5  9  5  4  7 10  1  8  4  3  9  6  3  6  7]
 [ 9  9  1 10  9  9  5  9  6  3  9  5  1 10  7 10  2  8  8  1]
 [ 8  8  7  8  6  7  7  8  6  7  8  6  7  6  8  8  7  6  8  8]
 [ 6  7  8  9 10  9  2  2  6  1 10  9  6  3  9  5  9  8  1  1]
 [ 7  8  7  8  6 10 10  6  8 10  8  9  8  8  5 10  8  7  8  7]]
```

Kiểu dữ liệu của phần tử trong mảng diem_2a: int8

Kích thước của mảng diem_2a: (10, 20)

Số phần tử của mảng diem_2a: 200

Số chiều của mảng diem_2a: 2

4. Các thao tác cơ bản

4.1 Quan sát mảng

#a.shape: Cho biết kích thước của mảng a:

```
print('kích thước của mảng diem_2a:', diem_2a.shape)
```

kích thước của mảng diem_2a: (10, 20)

#a.ndim: Cho biết Số chiều của mảng a:

```
print('Số chiều của mảng diem_2a:', diem_2a.ndim)
```

Số chiều của mảng diem_2a: 2

#a.size: Cho biết số phần tử của mảng a:

```
print('Số phần tử của mảng diem_2a: ', diem_2a.size)
```

Số phần tử của mảng diem_2a: 200

#a.dtype: Cho biết kiểu dữ liệu của các phần tử trong mảng a

```
print('Kiểu dữ liệu của các phần tử trong mảng diem_2a:', diem_2a.dtype)
```

Kiểu dữ liệu của các phần tử trong mảng diem_2a: int8

4.2 Chuyển đổi kiểu dữ liệu

```
1 #a.astype(kiểu mới): Chuyển đổi kiểu dữ liệu của các phần tử
2 a_float = np.linspace(0,15,11)
3 print(a_float)
4 print('Kiểu Dữ liệu: ', a_float.dtype)
5 print('-----')
6 #Chuyển từ kiểu float --> int
7 a_int = a_float.astype(np.int16)
8 print(a_int)
9 print('Dữ liệu sau khi chuyển: ', a_int.dtype)
```

```
[ 0.  1.5  3.  4.5  6.  7.5  9. 10.5 12. 13.5 15. ]
Kiểu Dữ liệu: float64
```

```
-----
[ 0  1  3  4  6  7  9 10 12 13 15]
Dữ liệu sau khi chuyển: int16
```

```
1 #Chuyển từ kiểu float --> int
2 a_str = a_int.astype(np.str)
3 print(a_str)
4 print('Dữ liệu sau khi chuyển: ', a_str.dtype)
5 print('-----')
6 #Chuyển từ kiểu float --> boolean
7 a_bool = a_int.astype(np.bool)
8 print(a_bool)
9 print('Dữ liệu sau khi chuyển: ', a_bool.dtype)
```

```
['0' '1' '3' '4' '6' '7' '9' '10' '12' '13' '15']
Dữ liệu sau khi chuyển: <U6
```

```
-----
[False  True  True  True  True  True  True  True  True  True  True]
Dữ liệu sau khi chuyển: bool
```

Chuyển đổi kiểu dữ liệu (2)



NumPy dtypes

Basic Type	Available NumPy types	Comments
Boolean	<code>bool</code>	Elements are 1 byte in size
Integer	<code>int8</code> , <code>int16</code> , <code>int32</code> , <code>int64</code> , <code>int128</code> , <code>int</code>	<code>int</code> defaults to the size of <code>int</code> in C for the platform
Unsigned Integer	<code>uint8</code> , <code>uint16</code> , <code>uint32</code> , <code>uint64</code> , <code>uint128</code> , <code>uint</code>	<code>uint</code> defaults to the size of unsigned <code>int</code> in C for the platform
Float	<code>float32</code> , <code>float64</code> , <code>float</code> , <code>longfloat</code> ,	Float is always a double precision floating point value (64 bits). <code>longfloat</code> represents large precision floats. Its size is platform dependent.
Complex	<code>complex64</code> , <code>complex128</code> , <code>complex</code>	The real and complex elements of a <code>complex64</code> are each represented by a single precision (32 bit) value for a total size of 64 bits.
Strings	<code>str</code> , <code>unicode</code>	Unicode is always UTF32 (UCS4)
Object	<code>object</code>	Represent items in array as Python objects.
Records	<code>void</code>	Used for arbitrary data structures in record arrays.

Chuyển đổi kiểu dữ liệu (3)

Data type	Description
<code>bool_</code>	Boolean (True or False) stored as a byte
<code>int_</code>	Default integer type (same as C <code>long</code> ; normally either <code>int64</code> or <code>int32</code>)
<code>intc</code>	Identical to C <code>int</code> (normally <code>int32</code> or <code>int64</code>)
<code>intp</code>	Integer used for indexing (same as C <code>ssize_t</code> ; normally either <code>int32</code> or <code>int64</code>)
<code>int8</code>	Byte (-128 to 127)
<code>int16</code>	Integer (-32768 to 32767)
<code>int32</code>	Integer (-2147483648 to 2147483647)
<code>int64</code>	Integer (-9223372036854775808 to 9223372036854775807)
<code>uint8</code>	Unsigned integer (0 to 255)
<code>uint16</code>	Unsigned integer (0 to 65535)
<code>uint32</code>	Unsigned integer (0 to 4294967295)
<code>uint64</code>	Unsigned integer (0 to 18446744073709551615)
<code>float_</code>	Shorthand for <code>float64</code> .
<code>float16</code>	Half precision float: sign bit, 5 bits exponent, 10 bits mantissa
<code>float32</code>	Single precision float: sign bit, 8 bits exponent, 23 bits mantissa
<code>float64</code>	Double precision float: sign bit, 11 bits exponent, 52 bits mantissa
<code>complex_</code>	Shorthand for <code>complex128</code> .
<code>complex64</code>	Complex number, represented by two 32-bit floats
<code>complex128</code>	Complex number, represented by two 64-bit floats

4.3 Truy cập tới các phần tử

- Truy cập tới phần tử trong một vector (1D)

```
1 #Truy cập tới một phần tử của Vector: a[index]
2 #Note: index phần tử đầu tiên 0
3 #      : index phần tử cuối cùng -1
4 a = np.array([3, 5, 3, 10, 9, 1, 9, 8, 3, 1])
5
6 print('các phần tử của Vector a:\n', a)
7 print('-----')
8 print('phần tử đầu tiên:', a[0])
9 print('phần tử thứ 3:', a[3])
10 print('phần tử cuối cùng:', a[-1])
```

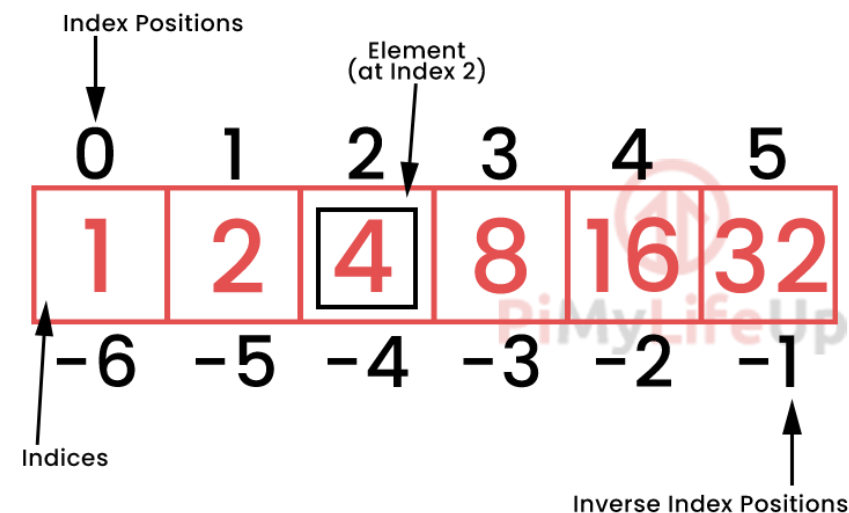
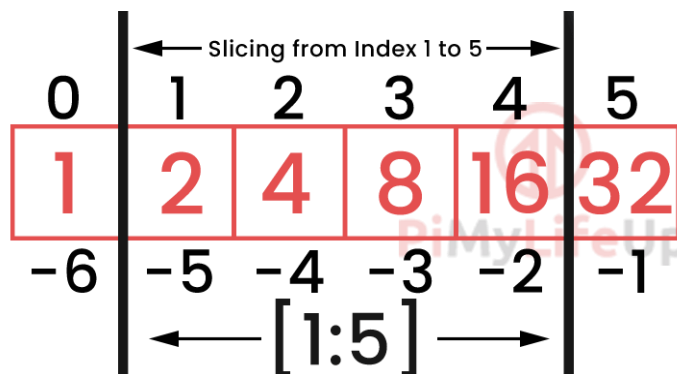
các phần tử của Vector a:

[3 5 3 10 9 1 9 8 3 1]

phần tử đầu tiên: 3

phần tử thứ 3: 10

phần tử cuối cùng: 1



```
1 #Truy cập tới nhiều phần tử của Vector: a[index1:index2]
2 print('các phần tử của Vector a:\n', a)
3 print('-----')
4 print('3 Phần tử đầu tiên:', a[:3])
5 print('Từ phần tử thứ 5 tới hết:', a[5:])
6 print('Từ phần tử 2 đến phần tử <6 của vector:', a[2:6])
```

các phần tử của Vector a:

[3 5 3 10 9 1 9 8 3 1]

3 Phần tử đầu tiên: [3 5 3]

Từ phần tử thứ 5 tới hết: [1 9 8 3 1]

Từ phần tử 2 đến phần tử <6 của vector: [3 10 9 1]

Truy cập tới các phần tử (2)

- Truy cập tới các phần tử trong một ma trận (2D)

```
1 #Truy cập tới 1 phần tử của ma trận (2D): a[index_row, index_col]
2 print('Điểm môn học đầu tiên, của học sinh đầu tiên:',diem_2a[0,0])
3 print('Điểm môn học thứ 1, của học sinh thứ 3:',diem_2a[1,3])
4 print('Điểm môn cuối cùng, của học sinh cuối cùng:',diem_2a[-1,-1])
5 print('-----')
6 print('Bảng điểm lớp 2A:\n',diem_2a)
```

Điểm môn học đầu tiên, của học sinh đầu tiên: 2

Điểm môn học thứ 1, của học sinh thứ 3: 10

Điểm môn cuối cùng, của học sinh cuối cùng: 7

	0	1	2	
0	(0,0)	(0,1)	(0,2)	← Column Index
1	(1,0)	(1,1)	(1,2)	
2	(2,0)	(2,1)	(2,2)	

↑ Row Index

		0	1	...	$n - 1$
axis = 0 ↓	0				
	1				
	...				
	$m - 1$				

Truy cập tới các phần tử (3)

- Truy cập tới các phần tử trong một ma trận (2D)

```
1 #Truy cập tới nhiều phần tử trong ma trận: a[index_row1:index_row2,index_col1:index_col2]
2 #Lấy điểm tất cả các môn (tất cả các hàng) của học sinh 5:
3 diem_hs5 = diem_2a[:,5]
4 print("Điểm các môn của học sinh 5:",diem_hs5)
5
6 #Lấy điểm môn học cuối cùng của tất cả học sinh (tất cả các cột)
7 diem_mon = diem_2a[-1,:]
8 print("Điểm môn học cuối cùng của tất cả học sinh: \n",diem_mon)
9
10 #Lấy điểm 5 môn học đầu tiên của 10 học sinh đầu tiên
11 diem5_hs10 = diem_2a[:5,:10]
12 print("Bảng điểm 5 môn học đầu tiên của 10 học sinh đầu của lớp:\n",diem5_hs10)
```

Điểm các môn của học sinh 5: [6 1 9 7 5 5 9 7 9 10]

Điểm môn học cuối cùng của tất cả học sinh:

```
[ 7  8  7  8  6 10 10  6  8 10  8  9  8  8  5 10  8  7  8  7  9  9  8  7
 7  7 10  8  9  7]
```

Bảng điểm 5 môn học đầu tiên của 10 học sinh đầu của lớp:

```
[[ 2  4  3  7  5  6  5  6  8  9]
 [ 3  5  3 10  9  1  9  8  3  1]
 [ 1 10  4  9  6  9  0  2  3  1]
 [ 6  3  0  8  3  7  7  2  6  8]
 [ 4  3  6  7  4  5  2  6  9  4]]
```


THỰC HÀNH 1

Thực hành 1



Yêu cầu 1: Học viên tạo một ma trận vuông cấp n (n hàng x n cột), bao gồm các phần tử là những số nguyên ngẫu nhiên trong khoảng (0-100) như minh họa, với $n = 12$

```
[[61 21 68 72 84 90 81 80 79 80 49 53]
 [92 10  1 43 49 93 76  6  2  6 69  2]
 [12 79 88 10 37 55 37  6 59 75 77 64]
 [97 22 75 32 39 39 93 19 28 64 55 87]
 [88 25 88 58 11 96 58 14 88 16 22 64]
 [ 3  5 60 14 65 50 80 42  8 27 44 52]
 [84 38 54 27 86 13 67 77 77 12 66 40]
 [96  9 94 24 61 19  2 80 95 92 72 32]
 [49 21 78 92 35 92 84 86 85 62 64 29]
 [ 5 53 95  2 43 30 72 66 97 17  8 23]
 [16 27  1 71 19 22 90 81 12 93 14 53]
 [40 50 83 25 37 16 49 73 42 86 11 18]]
```

Kiểu dữ liệu của phần tử trong ma trận: int32

Kích thước của mảng ma trận: (12, 12)

Số phần tử của mảng ma trận: 144

Số chiều của mảng ma trận: 2

Thực hành 1



Yêu cầu 2: Sử dụng ma trận tạo được trong yêu cầu 1, Học viên tạo 2 vector như sau:

- **v_chinh:** bao gồm các phần tử nằm trên đường chéo chính của ma trận.
- **v_phu:** bao gồm các phần tử nằm trên đường chéo phụ của ma trận

```
[[61 21 68 72 84 90 81 80 79 80 49 53]
 [92 10 1 43 49 93 76 6 2 6 69 2]
 [12 79 88 10 37 55 37 6 59 75 77 64]
 [97 22 75 32 39 39 93 19 28 64 55 87]
 [88 25 88 58 11 96 58 14 88 16 22 64]
 [ 3  5 60 14 65 50 80 42  8 27 44 52]
 [84 38 54 27 86 13 67 77 77 12 66 40]
 [96  9 94 24 61 19  2 80 95 92 72 32]
 [49 21 78 92 35 92 84 86 85 62 64 29]
 [ 5 53 95  2 43 30 72 66 97 17  8 23]
 [16 27  1 71 19 22 90 81 12 93 14 53]
 [40 50 83 25 37 16 49 73 42 86 11 18]]
```

Vector các phần tử nằm trên đường chéo chính:

```
[61. 10. 88. 32. 11. 50. 67. 80. 85. 17. 14. 18.]
```

Vector các phần tử nằm trên đường chéo phụ:

```
[40. 27. 95. 92. 61. 13. 80. 14. 28. 75. 69. 53.]
```

Thực hành 1



Yêu cầu 3: Nhập vào số nguyên x bất kỳ trong khoảng (0-100), đếm xem có bao nhiêu phần tử trong ma trận sinh ra ở yêu cầu 1 có giá trị bằng, lớn hơn và nhỏ hơn giá trị x:

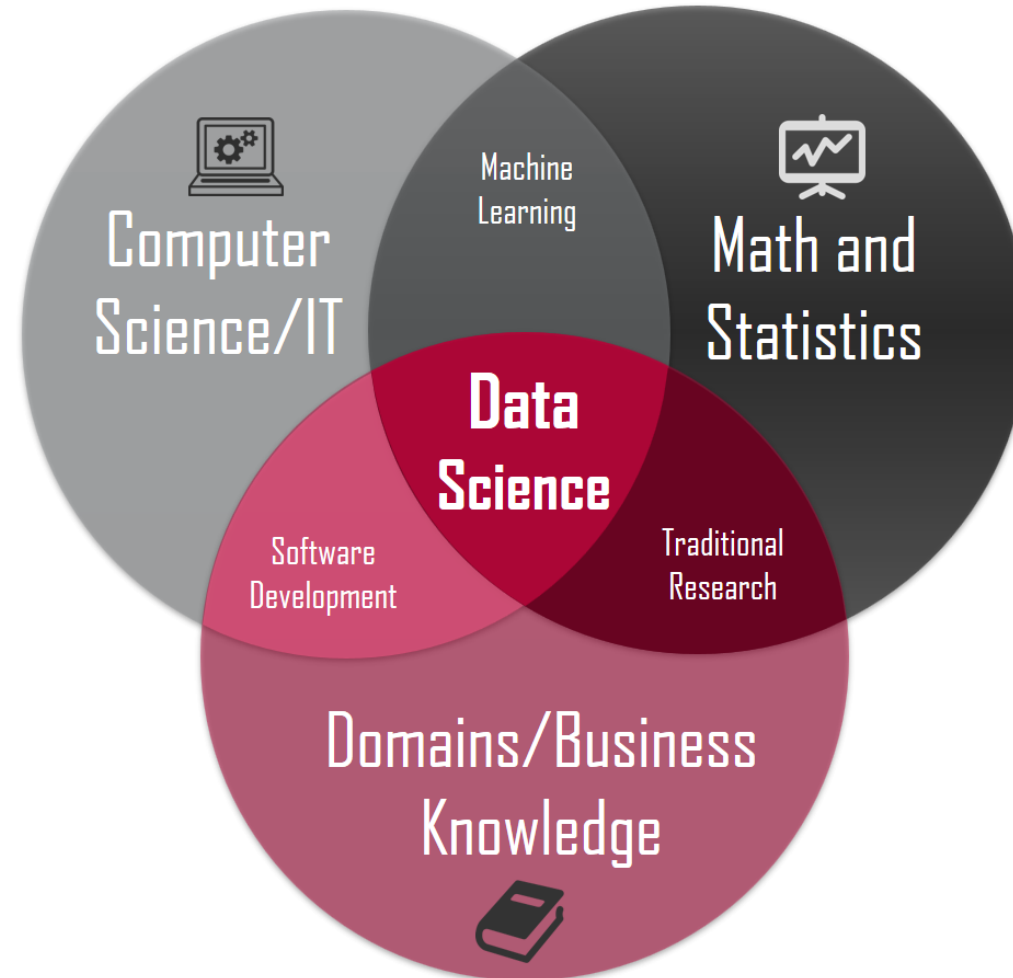
```
[[61 21 68 72 84 90 81 80 79 80 49 53]
 [92 10 1 43 49 93 76 6 2 6 69 2]
 [12 79 88 10 37 55 37 6 59 75 77 64]
 [97 22 75 32 39 39 93 19 28 64 55 87]
 [88 25 88 58 11 96 58 14 88 16 22 64]
 [3 5 60 14 65 50 80 42 8 27 44 52]
 [84 38 54 27 86 13 67 77 77 12 66 40]
 [96 9 94 24 61 19 2 80 95 92 72 32]
 [49 21 78 92 35 92 84 86 85 62 64 29]
 [5 53 95 2 43 30 72 66 97 17 8 23]
 [16 27 1 71 19 22 90 81 12 93 14 53]
 [40 50 83 25 37 16 49 73 42 86 11 18]]
```

Nhập vào giá trị x (0-100):88

-
1. Số phần tử có giá trị bằng x trong ma trận: 4
 2. Số phần tử nhỏ hơn giá trị x trong ma trận: 124
 3. Số phần tử lớn hơn giá trị x trong ma trận: 16

5. Tính toán các đặc trưng thống kê

5. Các đặc trưng thống kê



Toán học và thống kê có một vai trò rất quan trọng trong khoa học dữ liệu!

5.1 Max - Min

- **a.max():** Lấy giá trị lớn nhất của mảng a
- **b.min():** Lấy giá trị nhỏ nhất của mảng b

```
1 #Max - Min: Xác định giá trị lớn nhất, nhỏ nhất:  
2 #1) Hiển thị điểm cao nhất, thấp nhất của lớp 2A  
3 print('Điểm cao nhất của lớp:',diem_2a.max())  
4 print('Điểm thấp nhất của lớp:',diem_2a.min())
```

Điểm cao nhất của lớp: 10

Điểm thấp nhất của lớp: 0

```
1 #2) Liệt kê điểm cao nhất và thấp nhất theo môn học  
2 for i in range(0,diem_2a.shape[0]):  
3     print('Môn ', i,': Điểm Max: ', diem_2a[i,:].max(),  
4         '-- Điểm Min:',diem_2a[i,:].min())
```

Môn 0 : Điểm Max: 9 -- Điểm Min: 1

Môn 1 : Điểm Max: 10 -- Điểm Min: 0

Môn 2 : Điểm Max: 10 -- Điểm Min: 0

```
1 #3) Liệt kê điểm cao nhất và thấp nhất của mỗi học sinh  
2 for i in range(0,diem_2a.shape[1]):  
3     print('Học sinh ', i,': Điểm Max: ', diem_2a[:,i].max(),  
4         '-- Điểm Min:',diem_2a[:,i].min())
```

Học sinh 0 : Điểm Max: 9 -- Điểm Min: 1

Học sinh 1 : Điểm Max: 10 -- Điểm Min: 3



5.2 Sum

- **a.sum():**Tính tổng tất cả các phần tử của mảng a

```
1  #Sum:Tính tổng các phần tử trong mảng
2  print('Tổng tất các điểm trong của lớp 2A:', diem_2a.sum())
3  print('-----')
4
5  #Tính tổng điểm của từng học sinh:
6  for i in range(0,diem_2a.shape[1]):
7      print('Tổng điểm các môn của học sinh ', i, ' : ', diem_2a[:,i].sum())
```

Tổng tất các điểm trong của lớp 2A: 1731

Tổng điểm các môn của học sinh	0	:	48
Tổng điểm các môn của học sinh	1	:	60
Tổng điểm các môn của học sinh	2	:	47
Tổng điểm các môn của học sinh	3	:	86
Tổng điểm các môn của học sinh	4	:	62
Tổng điểm các môn của học sinh	5	:	68
Tổng điểm các môn của học sinh	6	:	56
Tổng điểm các môn của học sinh	7	:	54
Tổng điểm các môn của học sinh	8	:	59
Tổng điểm các môn của học sinh	9	:	51



5.3 Mean, Median, Mode, Range

Statistics – Mean, Median, Mode and Range

EZY MATHS

Mean

$$\text{Mean} = \frac{\text{Total of all values}}{\text{number of values}}$$

3, 3, 4, 5, 5, 8, 9, 15

$$\text{Mean} = \frac{52}{8} = 6.5$$

Collect it all together and share it out evenly

Using the mean to find the total amount

Mean \times Number of values

Ezytown FC have scored an average of 3.8 goals per game in their last 15 matches. How many goals have they scored?

$$3.8 \times 15 = 57 \text{ goals}$$

Median

Median = Middle value
(Numbers written in order)

3, 3, 4, 5, 5, 8, 9, 15

Median = 5

Finds the middle value

Use of formula to find location of median

$$\text{Location} = \frac{n + 1}{2}$$

The median of 45 values would be the 23rd number when written in order

$$\frac{45 + 1}{2} = 23$$

Mode

Mode = Most common value/item

3, 3, 4, 5, 5, 8, 9, 15

Mode = 3 and 5

Average usually used for qualitative data

Occurrence of no mode

If every value appears equally, there is no mode

1, 1, 3, 3, 7, 7

Each value appears twice so there is no mode

Range

Range = Largest - Smallest

3, 3, 4, 5, 5, 8, 9, 15

Range = 15 - 3 = 12

Reveals how close/far apart the values are

Interpreting measures of spread

The Smaller the range, the closer and more 'consistent' the values are.

The Larger the range, the more varied and more 'inconsistent' the values are.

Mean

```
1 # a.mean(): Giá trị trung bình của mảng a
2 print('Điểm trung bình của cả lớp 2A:', diem_2a.mean())
3 print('-----')
4 #Tính điểm trung bình của các học sinh trong lớp:
5 #CÁCH 1:
6 for i in range(0, diem_2a.shape[1]):
7     print('Điểm trung bình của học sinh ', i, ' : ', diem_2a[:,i].mean())
```

Điểm trung bình của cả lớp 2A: 5.77

Điểm trung bình của học sinh 0 : 4.8
Điểm trung bình của học sinh 1 : 6.0
Điểm trung bình của học sinh 2 : 4.7

```
1 #Tính điểm trung bình của các học sinh trong lớp:
2 #CÁCH 2:
3 mean_2a = diem_2a.mean(axis=0)
4 #axis = 0: theo hàng
5 #axis = 1: theo cột
6 for i in range(0, mean_2a.size):
7     print('Điểm trung bình của học sinh ', i, ' : ', mean_2a[i])
```

Điểm trung bình của học sinh 0 : 4.8
Điểm trung bình của học sinh 1 : 6.0
Điểm trung bình của học sinh 2 : 4.7

Mean

$$\text{Mean} = \frac{\text{Total of all values}}{\text{number of values}}$$

3, 3, 4, 5, 5, 8, 9, 15

$$\text{Mean} = \frac{52}{8} = 6.5$$

Collect it all together and
share it out evenly

Using the mean to find the
total amount

$\text{Mean} \times \text{Number of values}$

Ezytown FC have scored an
average of 3.8 goals per game
in their last 15 matches. How
many goals have they scored?

$$3.8 \times 15 = 57 \text{ goals}$$

Median

- `np.median(a)`: Tìm trung vị của mảng a

```
1 #median(): Giá trị trung vị trong một tập hợp các phần tử.  
2 #Trường hợp số phần tử trong mảng là lẻ  
3 a=diem_2a[1,:15]  
4  
5 print('Mảng a ban đầu: \n', a)  
6 print('Số phần tử trong mảng a: ', a.size)  
7 print('Mảng a đã sắp xếp: \n', np.sort(a))  
8 print('Giá trị trung bình mean:', np.mean(a))  
9 print('Giá trị trung vị median:', np.median(a))
```

Mảng a ban đầu:

[3 5 3 10 9 1 9 8 3 1 6 0 7 10 8]

Số phần tử trong mảng a: 15

Mảng a đã sắp xếp:

[0 1 1 3 3 3 5 6 7 8 8 9 9 10 10]

Giá trị trung bình mean: 5.533333333333333

Giá trị trung vị median: 6.0

Mảng a ban đầu:

[9 1 1 8 4 7 3 7 1 10]

Số phần tử trong mảng a: 10

Mảng a đã sắp xếp:

[1 1 1 3 4 7 7 8 9 10]

Giá trị trung bình mean: 5.1

Giá trị trung vị median: 5.5

Median

Median = Middle value
(Numbers written in order)

3, 3, 4, 5, 5, 8, 9, 15

Median = 5

Finds the middle value

Use of formula to find
location of median

$$Location = \frac{n + 1}{2}$$

The median of 45 values
would be the 23rd number
when written in order

$$\frac{45 + 1}{2} = 23$$

Mode

```
1 #C) Mode: là giá trị xuất hiện nhiều nhất trong tập hợp.
2 #Trong trường hợp không có giá trị nào được lặp lại thì không có Mode.
3 #Liệt kê điểm xuất hiện nhiều nhất theo từng môn học
4 from scipy import stats as sp #sử dụng thư viện scipy để dùng hàm mode
5
6 for i in range(0,diem_2a.shape[0]):
7     a = sp.mode(diem_2a[i,:])
8     print('Môn ', i, ': Điểm xuất hiện nhiều nhất: ', a[0],
9           ' số lần: ', a[1])
10 print(type(a))
```

Môn 0	: Điểm xuất hiện nhiều nhất:	[1]	số lần:	[6]
Môn 1	: Điểm xuất hiện nhiều nhất:	[1]	số lần:	[6]
Môn 2	: Điểm xuất hiện nhiều nhất:	[9]	số lần:	[8]
Môn 3	: Điểm xuất hiện nhiều nhất:	[6]	số lần:	[5]
Môn 4	: Điểm xuất hiện nhiều nhất:	[4]	số lần:	[6]
Môn 5	: Điểm xuất hiện nhiều nhất:	[5]	số lần:	[5]
Môn 6	: Điểm xuất hiện nhiều nhất:	[9]	số lần:	[8]
Môn 7	: Điểm xuất hiện nhiều nhất:	[8]	số lần:	[15]
Môn 8	: Điểm xuất hiện nhiều nhất:	[8]	số lần:	[7]
Môn 9	: Điểm xuất hiện nhiều nhất:	[8]	số lần:	[10]

<class 'scipy.stats.stats.ModeResult'>

Mode

Mode = Most common
value/item

3, 3, 4, 5, 5, 8, 9, 15

Mode = 3 and 5

Average usually used for
qualitative data

Occurrence of no mode

If **every** value appears
equally, there is **no mode**

1, 1, 3, 3, 7, 7

Each value appears twice so
there is no mode

Range

Trong thư viện numpy không có hàm tính range, ta có thể xác định giá trị range bằng cách tính thông qua max - min

```
1 #D) Range: là sự khác biệt, khoảng cách giữa phần tử dưới và phần tử trên,  
2 #giữa giá trị nhỏ nhất (Min) với giá trị lớn nhất (Max) trong tập hợp.  
3 #Xác định độ chênh điểm max - min của từng học sinh  
4  
5 for i in range(0,diem_2a.shape[1]):  
6     print('Độ chênh điểm của học sinh ', i, ' : ',  
7         diem_2a[:,i].max()-diem_2a[:,i].min())
```

```
Độ chênh điểm của học sinh 0 : 8  
Độ chênh điểm của học sinh 1 : 7  
Độ chênh điểm của học sinh 2 : 8  
Độ chênh điểm của học sinh 3 : 3  
Độ chênh điểm của học sinh 4 : 7  
Độ chênh điểm của học sinh 5 : 9  
Độ chênh điểm của học sinh 6 : 10  
Độ chênh điểm của học sinh 7 : 7  
Độ chênh điểm của học sinh 8 : 6  
Độ chênh điểm của học sinh 9 : 9  
Độ chênh điểm của học sinh 10 : 7
```

Range

Range = Largest - Smallest

3, 3, 4, 5, 5, 8, 9, 15

Range = 15 - 3 = 12

Reveals how close/far
apart the values are

Interpreting measures of
spread

The **Smaller** the range, the
closer and more 'consistent'
the values are.

The **Larger** the range, the
more varied and more
'inconsistent' the values are.

5.4 std

Độ lệch tiêu chuẩn (*standard deviation*) là đại lượng thường được sử dụng để phản ánh **mức độ phân tán** của một **biến số xung quanh số bình quân**.

```
1 #E) Std: Tính độ lệch chuẩn
2 a = np.array([10,1,1,9,12,1,9,12,10])
3 print('Phần tử của mảng a:',a)
4 print('Giá trị trung bình:',a.mean())
5 print('Độ lệch chuẩn:',a.std())
6
7 print('-----')
8 b = np.array([7,7,8,7,8,7,7,7,7])
9 print('Phần tử của mảng b:',b)
10 print('Giá trị trung bình:',b.mean())
11 print('Độ lệch chuẩn:',b.std())
```

Phần tử của mảng a: [10 1 1 9 12 1 9 12 10]
Giá trị trung bình: 7.222222222222222
Độ lệch chuẩn: 4.516089207311461

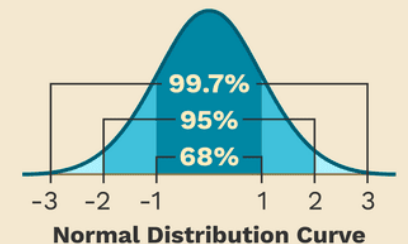
Phần tử của mảng b: [7 7 8 7 8 7 7 7 7]
Giá trị trung bình: 7.222222222222222
Độ lệch chuẩn: 0.41573970964154905

- Nếu độ lệch chuẩn bằng 0, suy ra các giá trị quan sát cũng chính là giá trị trung bình. Nói cách khác là không có sự biến thiên.
- Nếu độ lệch chuẩn càng lớn, suy ra sự biến thiên xung quanh giá trị trung bình càng lớn.

Calculating Standard Deviation

$$S_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

- n = The number of data points
- x_i = Each of the values of the data
- \bar{x} = The mean of x_i



THỰC HÀNH 2

Thực hành 2



Yêu cầu 1: Sử dụng dữ liệu bảng điểm của lớp 2A. Cho biết:

- 1. ĐTB của từng học sinh trong lớp.**
- 2. Học sinh có điểm TB cao nhất.**
- 3. Học sinh có điểm trung bình thấp nhất**

Điểm TB của từng học sinh trong lớp:

```
[4.8 6.  4.7 8.6 6.2 6.8 5.6 5.4 5.9 5.1 7.2 5.4 5.9 5.9 6.  7.9 4.2 6.1  
6.3 4.4 4.7 5.9 5.6 4.7 6.4 6.2 6.5 4.6 5.8 4.3]
```

Điểm TB cao nhất: 8.6

Của học sinh thứ: 3

Bảng điểm đầy đủ của học sinh: [7 10 9 8 7 10 10 8 9 8]

Điểm TB thấp nhất: 4.2

Của học sinh thứ: 16

Bảng điểm đầy đủ của học sinh: [3 2 2 1 2 6 2 7 9 8]

Thực hành 2



Yêu cầu 2: Sử dụng dữ liệu bảng điểm của lớp 2A. Cho biết:

- 1. ĐTB của từng môn học.**
- 2. Môn học có điểm TB cao nhất.**
- 3. Môn học có điểm trung bình thấp nhất**

Điểm TB của từng môn học:

```
[4.73 4.43 5.5 4.83 4.97 5.6 6.23 7.3 6.13 7.97]
```

Điểm TB cao nhất của môn học: 7.97

Môn học thứ [9]

Bảng điểm đầy đủ của môn học:

```
[[ 7 8 7 8 6 10 10 6 8 10 8 9 8 8 5 10 8 7 8 7 9 9 8 7
  7 7 10 8 9 7]]
```

Điểm TB thấp nhất của môn học: 4.43

Môn học thứ [1]

Bảng điểm đầy đủ của môn học:

```
[[ 3 5 3 10 9 1 9 8 3 1 6 0 7 10 8 5 2 7 7 1 1 6 1 6
  3 0 2 2 1 6]]
```

Thực hành 2



Yêu cầu 3: Sử dụng dữ liệu bảng điểm của lớp 2A. Cho biết:

- 1. Sinh viên có điểm đồng đều nhất tất cả các môn. Sinh viên có điểm các môn lệch nhất trong lớp.**
- 2. Môn học có điểm đồng đều nhất. Môn học có điểm chênh lệch nhất.**

Học sinh học đồng đều nhất: [3]

Bảng điểm của học sinh học đồng đều: `[[7 10 9 8 7 10 10 8 9 8]]`

Học sinh học lệch nhất: [19]

Bảng điểm của học sinh học lệch: `[[5 1 0 10 4 7 1 8 1 7]]`

Môn học có điểm đồng đều nhất: [7]

Bảng điểm của môn học đồng đều: `[[8 8 7 8 6 7 7 8 6 7 8 6 7 6 8 8 7 6 8 8 8 7 8 8 8 6 8 7 7 8]]`

Môn học có điểm lệch nhất: [2]

Bảng điểm của môn học lệch: `[[1 10 4 9 6 9 0 2 3 1 8 6 8 4 2 9 2 9 5 0 4 1 7 3
8 9 8 9 9 9]]`

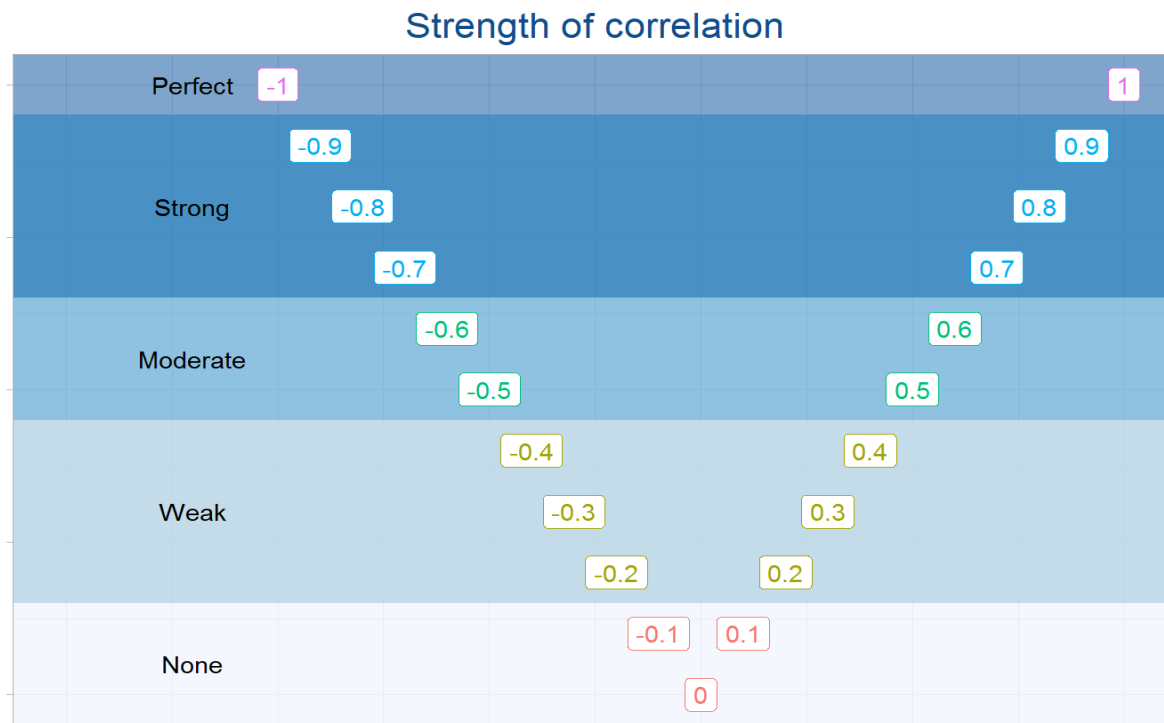
5.5 Hệ số tương quan

Hệ số tương quan đo lường mức độ quan hệ tuyến tính giữa hai biến.

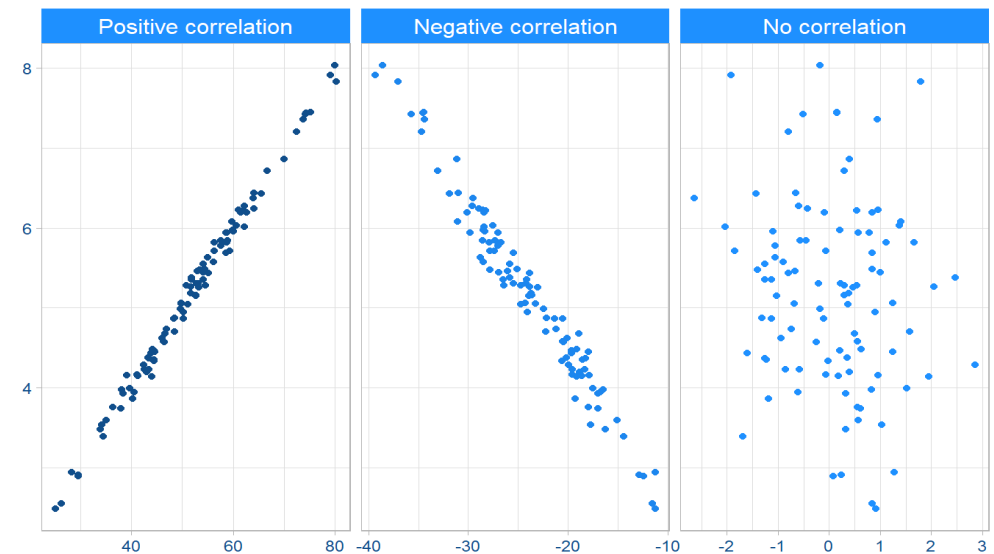
- Hệ số tương quan không có đơn vị
- Hệ số tương quan nằm trong khoảng $[-1,1]$

Correlation Coefficient Formula

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$



Author: NNB



Author: NNB

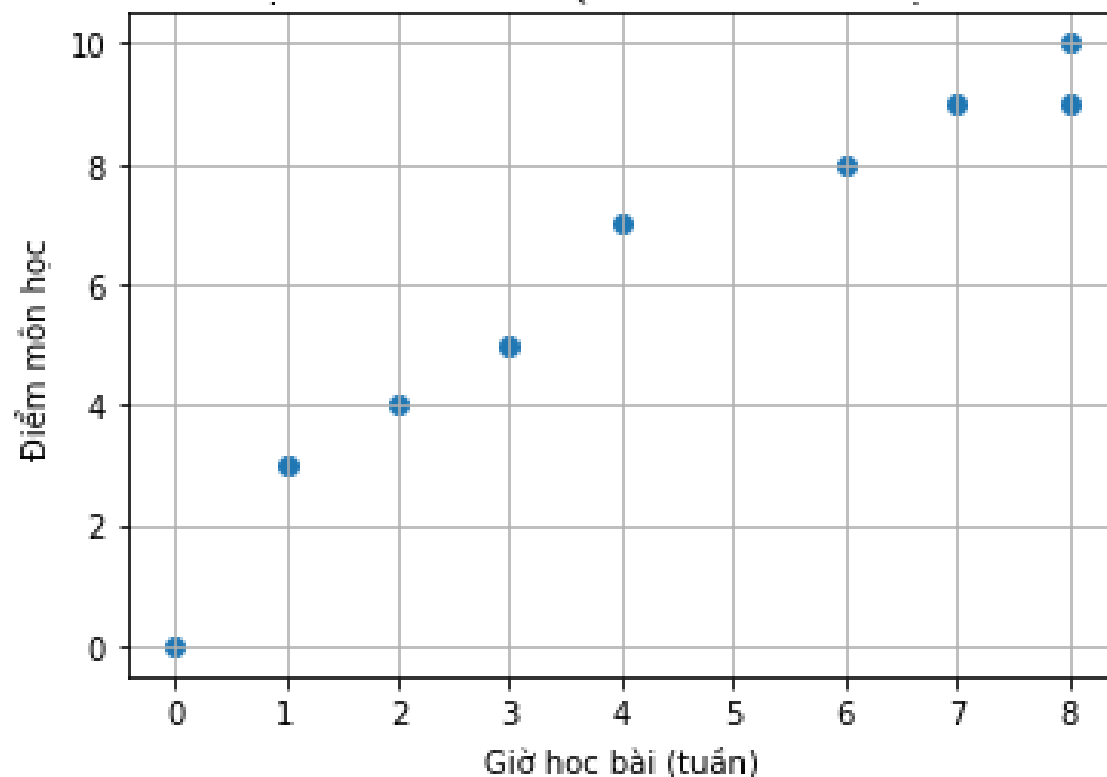
Hệ số tương quan (2)

```
#corrcoef: Hệ số tương quan
#Thời gian dành cho học bài
a_giohoc = np.array([4,7,1,2,8,0,3,8,6])
#Điểm thi nhận được:
b_diem = np.array([7,9,3,4,9,0,5,10,8])
co = np.corrcoef(a_giohoc,b_diem)
print(type(co))
print('Hệ số tương quan: \n', co)
```

```
<class 'numpy.ndarray'>
Hệ số tương quan:
[[1.          0.96995403]
 [0.96995403  1.          ]]
```

Ví dụ về mối tương quan giữa **thời gian dành cho việc học bài** với **điểm thi nhận được**!

BIỂU ĐỒ THỂ HIỆN MỐI TƯƠNG QUAN GIỮA GIỜ HỌC BÀI VÀ ĐIỂM THI



THỰC HÀNH 3

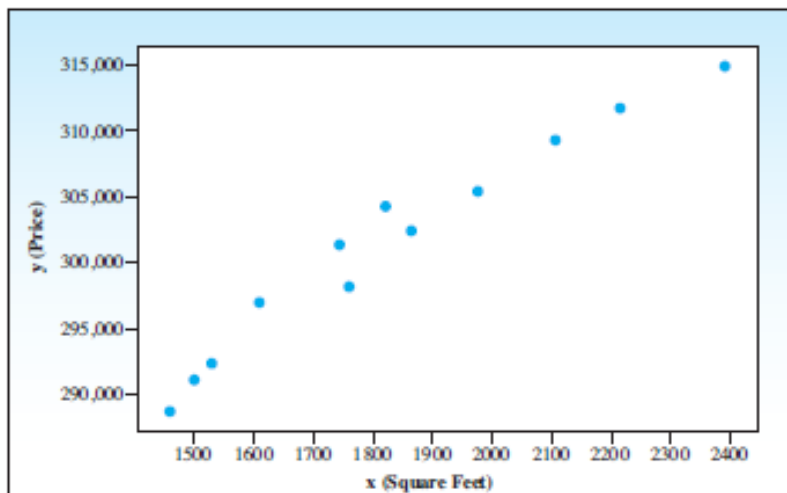
Thực hành 3

?

Yêu cầu: Xác định hệ số tương quan giữa diện tích (1) | Khoảng cách từ trung tâm thành phố (2) và giá bán nhà theo bảng số liệu dưới đây:

Square Feet, x	Price, y	Square Feet, x	Price, y
1460	\$288,700	1977	\$305,400
2108	309,300	1610	297,000
1743	301,400	1530	292,400
1499	291,100	1759	298,200
1864	302,400	1821	304,300
2391	314,900	2216	311,700

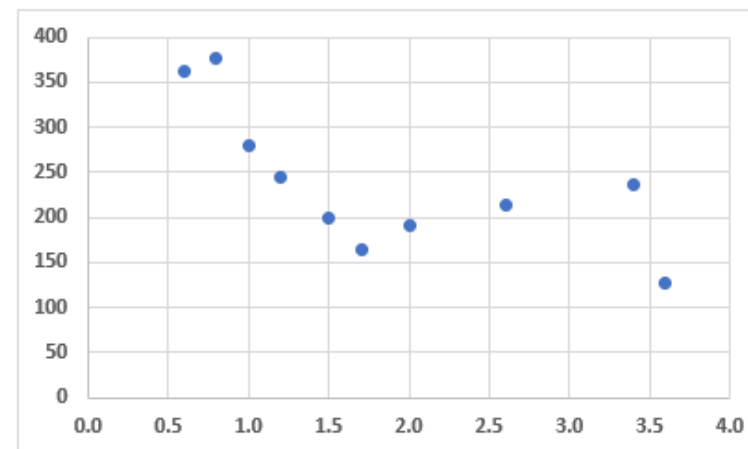
Plot of data for Exercise 12.42



01

Distance (in Miles from the city centre)	Prices (in \$ x 1000)
2.6	214
0.8	376
1.0	280
0.6	362
1.5	200
2.0	190
3.4	236
1.2	244
3.6	128
1.7	165

02



THỰC HÀNH 4

Thực hành 4

Mô tả file dữ liệu: Temp.txt

Temp - Notepad

File Edit Format View Help

25.65	24.79	24.01	25.06	25.48	24.97
25.31	24.21	24.02	24.93	25.16	24.83
25.05	23.73	23.89	24.79	24.80	24.55
24.79	23.36	23.83	24.84	24.74	24.48
24.59	23.05	23.69	24.82	24.80	24.38
24.40	22.80	23.52	24.79	24.87	24.40
24.38	22.79	23.68	25.10	24.71	24.41
26.72	25.61	24.92	26.56	25.03	24.91
28.84	26.93	26.51	26.53	25.75	25.85
30.29	28.72	27.48	26.95	26.64	26.79
31.35	29.97	26.96	27.23	27.68	27.53
32.05	28.93	26.86	27.38	28.43	28.98
31.31	28.94	26.65	27.47	28.29	29.24
30.95	30.25	27.83	27.44	28.00	30.66
30.56	30.62	26.49	27.16	27.67	30.97
31.13	30.58	26.29	26.68	27.29	30.59

Thực hành 4

Mô tả file dữ liệu: Temp.txt

- File dữ liệu lưu trữ nhiệt độ (°C) của 6 thành phố lớn dọc theo nước Việt Nam là: Hà Nội, Vinh, Đà Nẵng, Nha trang, Hồ Chính Minh và Cà Mau
- Thời gian từ 0h ngày 15/09/2024 tới 23h ngày 22/09/2024



Thực hành 4

Mô tả file dữ liệu: Temp.txt

Hà Nội	Vinh	Đà Nẵng	Nha Trang	HCM	Cà Mau	
25.65	24.79	24.01	25.06	25.48	24.97	Time: 0h 15/09
25.31	24.21	24.02	24.93	25.16	24.83	1h 15/09
25.05	23.73	23.89	24.79	24.80	24.55	
24.79	23.36	23.83	24.84	24.74	24.48	
24.59	23.05	23.69	24.82	24.80	24.38	
24.40	22.80	23.52	24.79	24.87	24.40	
24.38	22.79	23.68	25.10	24.71	24.41	
26.72	25.61	24.92	26.56	25.03	24.91	
28.84	26.93	26.51	26.53	25.75	25.85	
30.29	28.72	27.48	26.95	26.64	26.79	

Thực hành 4

Yêu cầu 1) Đọc dữ liệu lưu trữ trong file Temp.txt vào biến data_numpy, cho biết kích thước, số chiều, kiểu dữ liệu và số phần tử của biến data_numpy.

```
1 print(data_numpy)
2 print('-----')
3 print('Kích thước biến:',data_numpy.shape)
4 print('Số chiều của biến:',data_numpy.ndim)
5 print('Kiểu dữ liệu của các phần tử:',data_numpy.dtype)
6 print('Số phần tử:',data_numpy.size)
```

```
[[25.65 24.79 24.01 25.06 25.48 24.97]
 [25.31 24.21 24.02 24.93 25.16 24.83]
 [25.05 23.73 23.89 24.79 24.8 24.55]
 ...
 [24.81 24.47 23.4 25.86 25.05 25.29]
 [23.97 24.22 22.95 25.74 24.92 24.87]
 [22.84 23.99 22.59 25.5 24.77 24.57]]
```

```
-----
Kích thước biến: (192, 6)
Số chiều của biến: 2
Kiểu dữ liệu của các phần tử: float64
Số phần tử: 1152
```

Thực hành 4

Yêu cầu 2) Tìm nhiệt độ cao nhất (Max) – Thấp nhất (Min) – Nhiệt độ trung bình của cả 6 thành phố.

Yêu cầu 3) Tìm nhiệt độ cao nhất (Max) – Thấp nhất (Min) – Nhiệt độ trung bình của từng thành phố và hiển thị kết quả.

---THÔNG KÊ CHO CẢ 6 THÀNH PHỐ---

Nhiệt độ cao nhất: 33.45
Nhiệt độ thấp nhất: 20.93
Nhiệt độ trung bình: 26.50222222222222

1) Hà Nội
Nhiệt độ cao nhất: 33.45
Nhiệt độ thấp nhất: 21.68
Nhiệt độ trung bình: 27.71229166666667
2) Vinh (Nghệ An)
Nhiệt độ cao nhất: 32.57
Nhiệt độ thấp nhất: 22.6
Nhiệt độ trung bình: 26.719895833333336
3) Đà Nẵng
Nhiệt độ cao nhất: 29.88
Nhiệt độ thấp nhất: 20.93
Nhiệt độ trung bình: 25.522499999999997

4) Nha Trang
Nhiệt độ cao nhất: 28.68
Nhiệt độ thấp nhất: 24.5
Nhiệt độ trung bình: 26.166875000000005
5) TP Hồ Chí Minh
Nhiệt độ cao nhất: 31.06
Nhiệt độ thấp nhất: 23.22
Nhiệt độ trung bình: 26.159218749999997
6) Cà Mau
Nhiệt độ cao nhất: 31.37
Nhiệt độ thấp nhất: 23.99
Nhiệt độ trung bình: 26.732552083333333

Thực hành 4

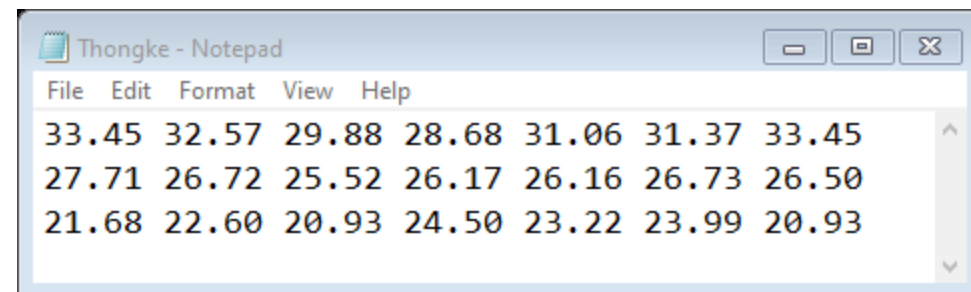
Yêu cầu 4) Tạo một ma trận **data_thongke** gồm 3 hàng x 7 cột; các hàng lần lượt lưu trữ dữ liệu như sau:

- hàng 0: Nhiệt độ lớn nhất (Max)
- hàng 1: Nhiệt độ trung bình (Mean), làm tròn đến 2 số sau dấu phẩy
- hàng 2: Nhiệt độ nhỏ nhất (Min)

Các cột lần lượt theo thứ tự của 6 thành phố và cột cuối cùng là cột thống kê chung cho cả 6 thành phố. Lưu ra file **thongke.txt**

```
1 print(data_thongke)
2 print(type(data_thongke))
3 print('Kích thước:', data_thongke.shape)

[[33.45 32.57 29.88 28.68 31.06 31.37 33.45]
 [27.71 26.72 25.52 26.17 26.16 26.73 26.5 ]
 [21.68 22.6  20.93 24.5  23.22 23.99 20.93]]
<class 'numpy.ndarray'>
Kích thước: (3, 7)
```



Thongke - Notepad

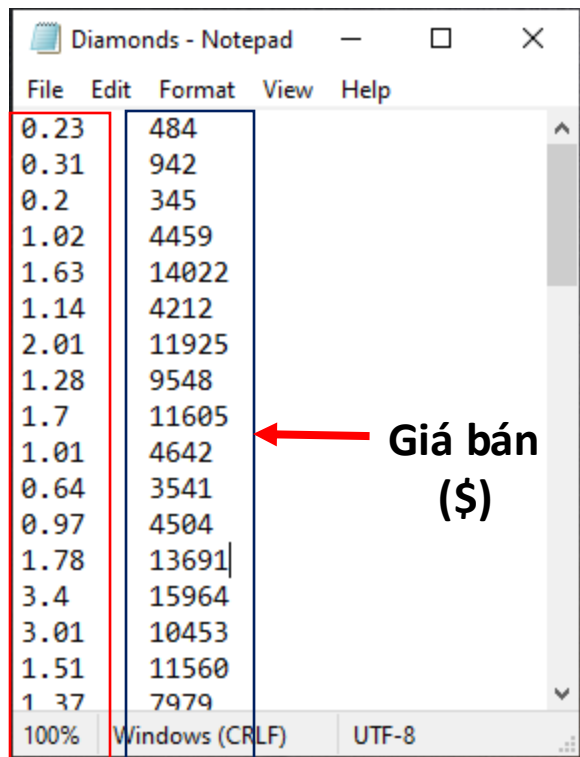
File Edit Format View Help

33.45 32.57 29.88 28.68 31.06 31.37 33.45
27.71 26.72 25.52 26.17 26.16 26.73 26.50
21.68 22.60 20.93 24.50 23.22 23.99 20.93

THỰC HÀNH 5

Thực hành 5

File dữ liệu: Diamonds.txt



0.23	484
0.31	942
0.2	345
1.02	4459
1.63	14022
1.14	4212
2.01	11925
1.28	9548
1.7	11605
1.01	4642
0.64	3541
0.97	4504
1.78	13691
3.4	15964
3.01	10453
1.51	11560
1.37	7979

Trọng lượng
(carat)

Giá bán
(\$)

Mô tả file dữ liệu:

- File dữ liệu lưu trữ thông số 50 viên kim cương bao gồm: Trọng lượng (carat) và Giá bán (\$) tương ứng



Thực hành 5

Học viên thực hiện các yêu cầu sau:

Yêu cầu 1) Đọc dữ liệu lưu trữ trong file Diamonds.txt vào biến kiểu mảng **data_diamond**, cho biết kích thước, số chiều, kiểu dữ liệu và số phần tử của biến data_diamond

```
[3.4000e-01 7.6500e+02]  
[4.1000e-01 8.2700e+02]  
[7.5000e-01 3.1200e+03]  
[1.0700e+00 5.2200e+03]  
[1.3400e+00 7.4270e+03]  
[1.7500e+00 9.8900e+03]]
```

```
-----  
Kích thước biến data_diamond: (50, 2)  
Số chiều của biến data_diamond: 2  
Kiểu dữ liệu của các phần tử: float64  
Số phần tử: 100
```


Thực hành 5

Yêu cầu 2) Tách mảng data_diamond thành 2 vector: diamond_size và diamond_price lưu trữ trọng lượng và giá bán.

Vector diamond_size:

```
[0.23 0.31 0.2 1.02 1.63 1.14 2.01 1.28 1.7 1.01 0.64 0.97 1.78 3.4  
3.01 1.51 1.37 1.5 0.54 0.72 1.13 2.24 3.01 4.5 0.92 1.05 0.55 0.74  
0.91 1.23 1.52 0.91 0.43 1.24 1.77 1.79 2.05 2.03 2.01 1. 0.9 1.01  
1.14 1.53 0.34 0.41 0.75 1.07 1.34 1.75]
```

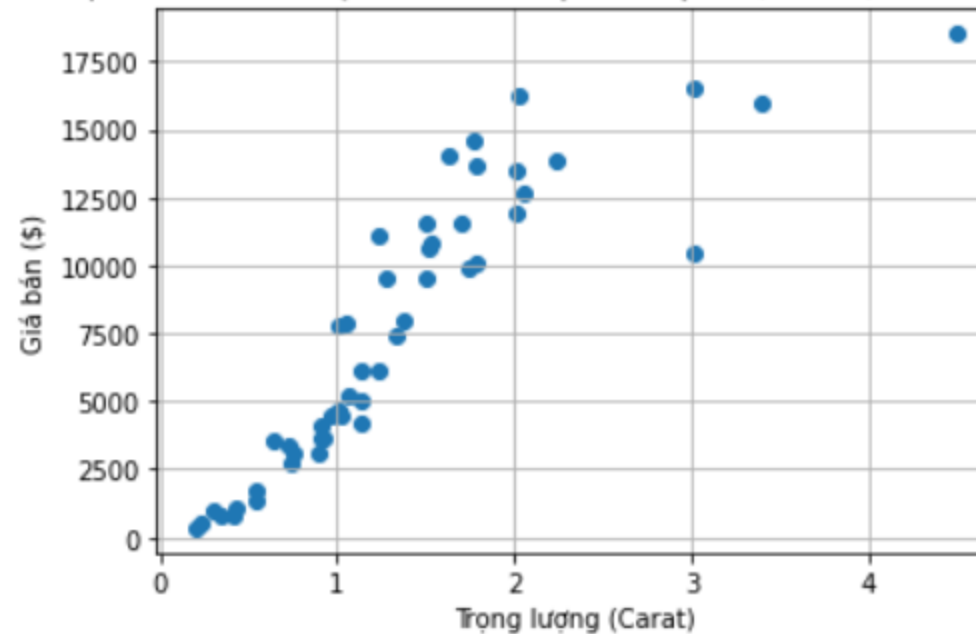
Vector diamond_price:

```
[ 484.  942.  345. 4459. 14022. 4212. 11925. 9548. 11605. 4642.  
3541. 4504. 13691. 15964. 10453. 11560. 7979. 9533. 1723. 3344.  
6133. 13827. 16538. 18531. 3625. 7879. 1319. 2761. 3620. 6165.  
10640. 4138. 1094. 11130. 14561. 10108. 12654. 16280. 13498. 4586.  
3105. 7745. 5047. 10830. 765. 827. 3120. 5220. 7427. 9890.]
```

Thực hành 5

Yêu cầu 3) Vẽ đồ thị thể hiện mối quan hệ giữa kích thước và giá bán kim cương. Xác định hệ số tương quan tương ứng giữa 2 thông số này.

BIỂU ĐỒ THỂ HIỆN MỐI TƯƠNG QUAN GIỮA TRỌNG LƯỢNG (CARAT) VÀ GIÁ BÁN KIM CƯƠNG (\$)




Hệ số tương quan giữa trọng lượng và giá bán kim cương: 0.8814849023922127

Thực hành 5

Yêu cầu 4) Cho biết kích thước và giá trung bình của 50 viên kim cương. Hiển thị giá bán của viên kim cương có trọng lượng 3.01 carat.

```
Trọng lượng trung bình: 1.3448  
Giá trung bình: 7550.78
```

```
Viên kim cương trọng lượng 3.01 carat có giá bán:  
Giá bán 1 : 10453.0  
Giá bán 2 : 16538.0
```



Q & A
Thank you!