

# AI VIET NAM – COURSE 2023

## Module 2 - Exercise 4

Numpy và Vectorization

Ngày 18 tháng 6 năm 2023

**Giới thiệu về bài tập:** Ở phần bài tập này các bạn sẽ được ôn tập cách sử dụng thư viện numpy để thao tác trên 1D và 2D array. Cũng như cách cài đặt giải thuật linear regression cho bài toán advertising theo cách vectorization dùng stochastic gradient descent, m samples (mini-batch gradient descent), and N samples (batch gradient descent).

**Bài tập 1** (Làm quen với numpy):

**Question 1:** Câu nào sau đây là đúng để tạo mảng 1 chiều từ 0 đến 9

- a) `import numpy as np`  
`arr = np.arange(0, 10, 1)`
- b) `import numpy as np`  
`arr = np.arange(1, 10, 1)`
- c) `import numpy as np`  
`arr = np.arange(0, 9, 1)`
- d) `import numpy as np`  
`arr = np.arange(1, 11, 1)`

**Question 2:** Cách tạo một mảng boolean 3x3 với tất cả giá trị là True

- a) `arr = np.ones((3,3)) > 0`
- b) `arr = np.ones((3,3), dtype=bool)`
- c) `arr = np.full((3,3), fill_value=True, dtype=bool)`
- d) a, b and c

**Question 3:** Kết quả của đoạn code sau đây:

```
1 import numpy as np
2 arr = np.arange(0,10)
3 print(arr[arr%2 == 1])
```

- a) `[1 3 5 7 9]`
- b) `[1 2 3 4 5]`
- c) `[2 4 6 8 9]`
- d) `[1 7 5 7 9]`

Question 4: Kết quả của đoạn code sau đây:

```
1 import numpy as np
2 arr = np.arange(0,10)
3 arr[arr%2 ==1] = -1
4 print(arr)
```

- a) [ -1 -1 -1 -1 -1 -1 -1 -1 -1 -1]
- b) [ 0 -1 2 -1 4 -1 6 -1 8 -1]
- c) [ -1 0 -1 -0 -1 0 -1 0 -1 0]
- d) [ -1 1 -1 1 -1 1 -1 1 -1 1]

Question 5: Kết quả của đoạn code sau đây:

```
1 import numpy as np
2 arr = np.arange(10)
3 arr_2d = arr.reshape(2,-1)
4 print(arr_2d)
```

- a) [[0 1 2 3 4 5 6 7 8 9]]
- b) [[0 1 2 3 4]  
[5 6 7 8 9]]
- c) [[0 1 2 3 4]]
- d) [[0 1 2 3]  
[5 6 7 8]]

Question 6: Kết quả của đoạn code sau đây:

```
1 import numpy as np
2 arr1 = np.arange(10).reshape(2,-1)
3 arr2 = np.repeat(1,10).reshape(2,-1)
4 c = np.concatenate([arr1,arr2],axis =0)
5 print("Result: \n", c)
```

- a) [[0 1 2 3 4]  
[5 6 7 8 9]  
[1 1 1 1 1]  
[1 1 1 1 1]]
- b) [[0 1 2 3 4]  
[5 6 7 8 9]]
- c) [[0 1 2 3 4]  
[5 6 7 8 9]  
[1 1 1 1 1]  
[5 6 7 8 9]]
- d) [[0 1 2 3 4]]

Question 7: Kết quả của đoạn code sau đây:

```
1 import numpy as np
2 arr1 = np.arange(10).reshape(2,-1)
3 arr2 = np.repeat(1,10).reshape(2,-1)
```

```

4 c = np.concatenate([arr1,arr2],axis =1)
5 print("C = ", c)

```

- a)  $\begin{bmatrix} 0 & 2 & 2 & 2 & 2 & 1 & 1 & 1 & 1 & 1 \\ 5 & 6 & 7 & 8 & 9 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$
- b)  $\begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 5 & 6 & 7 & 8 & 9 & 5 & 6 & 7 & 8 & 9 \end{bmatrix}$
- c)  $\begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 1 & 1 & 1 & 1 & 1 \\ 5 & 6 & 7 & 8 & 9 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$
- d)  $\begin{bmatrix} 0 & 2 & 2 & 2 & 2 & 1 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$

**Question 8: Kết quả của đoạn code sau đây:**

```

1 import numpy as np
2 arr = np.array([1,2,3])
3 print(np.repeat(arr,3))
4 print(np.tile(arr,3))

```

- a)  $\begin{bmatrix} 1 & 1 & 1 & 2 & 2 & 2 & 3 & 3 & 3 \\ 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \end{bmatrix}$
- b)  $\begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \end{bmatrix}$
- c)  $\begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 \end{bmatrix}$
- d)  $\begin{bmatrix} 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 \end{bmatrix}$

**Question 9: Kết quả của đoạn code sau đây:**

```

1 import numpy as np
2 a = np.array([2,6,1,9,10,3,27])
3 index = np.where((a>=5)&(a<=10))
4 print("result", a[index])

```

- a)  $[3 \quad 9 \quad 10]$
- b)  $[1 \quad 9 \quad 10]$
- c)  $[6 \quad 9 \quad 10]$
- d)  $[2 \quad 1 \quad 10]$

**Question 10: Kết quả của đoạn code sau đây:**

```

1 import numpy as np
2
3 def maxx(x,y):
4     if x >= y :
5         return x
6     else:
7         return y
8
9 a = np.array([5,7,9,8,6,4,5])

```

```

10 b = np.array([6,3,4,8,9,7,1])
11
12 pair_max = np.vectorize(maxx,otypes=[float])
13 print(pair_max(a,b))

```

- a) [5, 7, 9, 8, 6, 4, 5]
- b) [6, 3, 4, 8, 9, 7, 1]
- c) [6. 8. 9. 8. 7. 7. 5.]
- d) [6. 7. 9. 8. 9. 7. 5.]

**Question 11: Kết quả của đoạn code sau đây:**

```

1 import numpy as np
2
3 a = np.array([5,7,9,8,6,4,5])
4 b = np.array([6,3,4,8,9,7,1])
5
6 print("Result",np.where(a<b, b, a))

```

- a) [6. 7. 9. 8. 9. 7. 5.]
- b) [5, 7, 9, 8, 6, 4, 5]
- c) [6, 3, 4, 8, 9, 7, 1]
- d) [6. 8. 9. 8. 7. 7. 5.]

**Bài tập 2** (Làm quen với Ma trận):

**Matrix Properties:** Các bạn chứng minh các tính chất của ma trận sau:

1.  $AB \neq BA$ , ( $A \in \mathbb{R}^{m \times m}$  and  $B \in \mathbb{R}^{m \times m}$ ) ( $A$  và  $B$  là ma trận vuông)
2.  $(A + B)^T = A^T + B^T$ , ( $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{m \times n}$ )
3.  $(AB)^T = B^T A^T$ , ( $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times k}$ )

**Example:**

- $A + B = B + A$ , ( $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{m \times n}$ )

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & \dots & b_{1n} \\ \dots & \dots & \dots \\ b_{m1} & \dots & b_{mn} \end{bmatrix}$$

$$A + B = \begin{bmatrix} (a_{11} + b_{11}) & \dots & (a_{1n} + b_{1n}) \\ \dots & \dots & \dots \\ (a_{m1} + b_{m1}) & \dots & (a_{mn} + b_{mn}) \end{bmatrix} = \begin{bmatrix} (b_{11} + a_{11}) & \dots & (b_{1n} + a_{1n}) \\ \dots & \dots & \dots \\ (b_{m1} + a_{m1}) & \dots & (b_{mn} + a_{mn}) \end{bmatrix} = B + A$$

**Bài tập 3** (Hiện thực giải thuật linear regression theo phương pháp vectorization):

**Linear Regression:** Các bạn thực hiện train linear regression model trên tập data advertising.csv theo các yêu cầu sau. Các bạn sẽ dựa trên 3 thông tin đầu vào là TV, Radio, Newspaper để dự đoán Sale.

**Giới thiệu về tập data:** Data có 200 samples (rows), gồm 4 cột thông tin Tv, Radio, Newspaper, và Sales. Đề bài yêu cầu dùng thông tin ở 3 cột đầu tiên (Tv, Radio, Newspaper) để dự đoán được cột cuối cùng (Sale) dùng linear regression model.

Để chuẩn hoá data đầu vào, AIVN cung cấp trước cho các bạn function đọc dữ liệu và chuẩn hoá `mean_normalization(X)` như bên dưới:

```

1 # dataset
2 data = genfromtxt('advertising.csv', delimiter=',', skip_header=1)
3 N = data.shape[0]
4 X = data[:, :3]
5 y = data[:, 3:]
6
7 # Normalize input data by using mean normalizaton
8 def mean_normalization(X):
9     N = len(X)
10    maxi = np.max(X)
11    mini = np.min(X)
12    avg = np.mean(X)
13    X = (X-avg) / (maxi-mini)
14    X_b = np.c_[np.ones((N, 1)), X]
15    return X_b, maxi, mini, avg
16
17 X_b, maxi, mini, avg = mean_normalization(X)

```

Yêu cầu của bài tập này là các bạn lần lượt hiện thực lại giải thuật linear regression để dự đoán Sales dựa vào các yêu cầu sau:

1. Hoàn thành function `stochastic_gradient_descent()` để huấn luyện data sử dụng Stochastic Gradient Descent. Lưu ý các bạn cần tận dụng tối đa vectorization để hoàn thiện bài tập này.

- **input:** (4 inputs) `X_b`, `y`, `n_epochs`, `learning_rate`
- **output:** `thetas_path`, `losses`

```

1     def stochastic_gradient_descent(X_b, y, n_epochs=50, learning_rate=0.00001):
2
3     # thetas = np.random.randn(4, 1) # uncomment this line for real application
4     thetas = np.asarray([[1.16270837], [-0.81960489], [1.39501033],
5                          [0.29763545]])
6
7     thetas_path = [thetas]
8     losses = []
9
10    for epoch in range(n_epochs):
11        for i in range(N):
12            # select random number in N
13            # random_index = np.random.randint(N) #In real application, you
14            # should use this code
15            random_index = i # This code is used for this assignment only
16
17            xi = X_b[random_index:random_index+1]
18            yi = y[random_index:random_index+1]
19
20            # Compute output
21            *****Your code here *****
22
23            # Compute loss li
24            *****Your code here *****
25
26            # Compute gradient for loss

```

```

25         *****Your code here *****
26
27         # Compute gradient
28         *****Your code here *****
29
30         # update theta
31         *****Your code here *****
32
33         # logging
34         *****Your code here *****
35
36     return thetas_path, losses
37

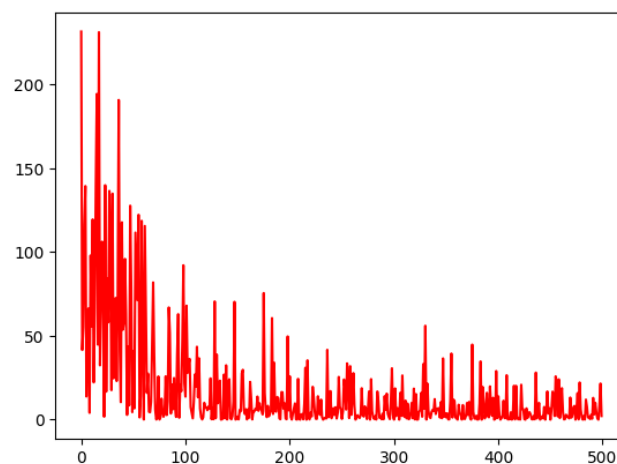
```

Hình 1 là kết quả sau khi thực thi đoạn code sau:

```

1  sgd_theta, losses = stochastic_gradient_descent(X_b, y, n_epochs=50,
2      learning_rate=0.01 )
3  x_axis = list(range(500))
4  plt.plot(x_axis, losses[:500], color="r")
5  plt.show()

```



Hình 1: Kết quả loss values sử dụng Stochastic Gradient Descent

**Question 12: Kết quả của đoạn code sau đây:**

```

1  sgd_theta, losses = stochastic_gradient_descent(X_b, y, n_epochs=1, learning_rate
2      =0.01 )
3  print(np.sum(losses))

```

- a) 7754.64
- b) 6754.64
- c) 8754.64
- d) 9754.64

2. Hoàn thành fuction **mini\_batch\_gradient\_descent()** để huấn luyện data sử dụng Mini-batch Gradient Descent. Lưu ý các bạn cần tận dụng tối đa vectorization để hoàn thiện bài tập này.

- **input:** (5 inputs) X\_b, y, n\_epochs, minibatch\_size, learning\_rate
- **output:** thetas\_path, losses

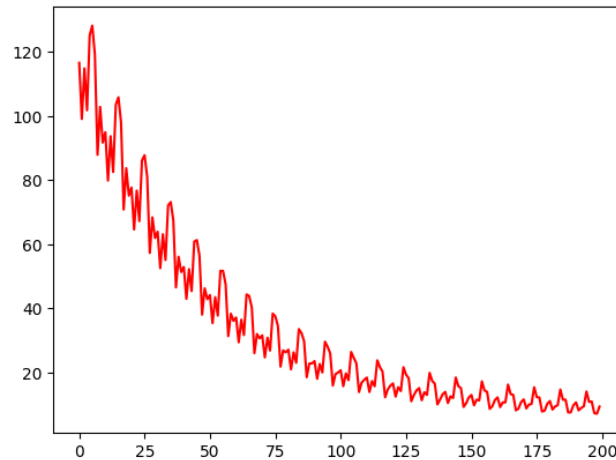
```

1     def mini_batch_gradient_descent(X_b, y, n_epochs=50, minibatch_size =
2     20, learning_rate=0.01):
3
4     # thetas = np.random.randn(4, 1)
5     thetas = np.asarray([[1.16270837], [-0.81960489], [1.39501033],
6     [0.29763545]])
7
8     thetas_path = [thetas]
9     losses = []
10
11     for epoch in range(n_epochs):
12         # shuffled_indices = np.random.permutation(N) # uncomment this code for
13         # real application
14
15         shuffled_indices = np.asarray([21, 144, 17, 107, 37, 115, 167, 31, 3,
16         132, 179, 155, 36, 191, 182, 170, 27, 35, 162, 25, 28, 73, 172, 152, 102, 16,
17         185, 11, 1, 34, 177, 29, 96, 22, 76, 196, 6, 128, 114, 117, 111, 43, 57, 126,
18         165, 78, 151, 104, 110, 53, 181, 113, 173, 75, 23, 161, 85, 94, 18, 148, 190,
19         169, 149, 79, 138, 20, 108, 137, 93, 192, 198, 153, 4, 45, 164, 26, 8, 131,
20         77, 80, 130, 127, 125, 61, 10, 175, 143, 87, 33, 50, 54, 97, 9, 84, 188, 139,
21         195, 72, 64, 194, 44, 109, 112, 60, 86, 90, 140, 171, 59, 199, 105, 41, 147,
22         92, 52, 124, 71, 197, 163, 98, 189, 103, 51, 39, 180, 74, 145, 118, 38, 47,
23         174, 100, 184, 183, 160, 69, 91, 82, 42, 89, 81, 186, 136, 63, 157, 46, 67,
24         129, 120, 116, 32, 19, 187, 70, 141, 146, 15, 58, 119, 12, 95, 0, 40, 83, 24,
25         168, 150, 178, 49, 159, 7, 193, 48, 30, 14, 121, 5, 142, 65, 176, 101, 55,
26         133, 13, 106, 66, 99, 68, 135, 158, 88, 62, 166, 156, 2, 134, 56, 123, 122,
27         154])
28
29         X_b_shuffled = X_b[shuffled_indices]
30         y_shuffled = y[shuffled_indices]
31
32         for i in range(0, N, minibatch_size):
33             xi = X_b_shuffled[i:i+minibatch_size]
34             yi = y_shuffled[i:i+minibatch_size]
35
36             # compute output
37             *****Your code here *****
38
39             # compute loss
40             *****Your code here *****
41
42             # compute derivative of loss
43             *****Your code here *****
44
45             # compute derivative of parameters
46             *****Your code here *****
47
48             # update parameters
49             thetas = thetas - learning_rate*gradients
50             thetas_path.append(thetas)
51
52             loss_mean = np.sum(loss)/minibatch_size
53             losses.append(loss_mean)
54
55     return thetas_path, losses

```

Hình 2 là kết quả sau khi thực thi đoạn code sau:

```
1 mbgd_thetas, losses = mini_batch_gradient_descent(X_b, y, n_epochs=50,
2 minibatch_size = 20, learning_rate=0.01)
3 x_axis = list(range(200))
4 plt.plot(x_axis, losses[:200], color="r")
5 plt.show()
```



Hình 2: Kết quả loss values sử dụng Mini batch Gradient Descent

**Question 13: Kết quả của đoạn code sau đây:**

```
1 mbgd_thetas, losses = mini_batch_gradient_descent(X_b, y, n_epochs=50,
2 minibatch_size = 20, learning_rate=0.01)
3 print(round(sum(losses),2))
```

- a) 7865.65
- b) 6865.65
- c) 5865.65
- d) 8865.65

3. Hoàn thành function **batch\_gradient\_descent()** để huấn luyện data sử dụng batch Gradient Descent. Lưu ý các bạn cần tận dụng tối đa vectorization để hoàn thiện bài tập này.

- **input:** (4 inputs) X\_b, y, n\_epochs, learning\_rate
- **output:** thetas\_path, losses

```
1 def batch_gradient_descent(X_b, y, n_epochs=100, learning_rate=0.01):
2
3     # thetas = np.random.randn(4, 1) # uncomment this line for real application
4     thetas = np.asarray([[1.16270837], [-0.81960489], [1.39501033],
5                          [0.29763545]])
6
7     thetas_path = [thetas]
8     losses = []
9
10    for i in range(n_epochs):
```



```

10     # compute output
11     *****Your code here *****
12
13
14     # Compute loss
15     *****Your code here *****
16
17
18     # Compute losses's derivative
19     *****Your code here *****
20
21
22     # computer parameters' derivative
23     *****Your code here *****
24
25     # Update parameters
26     thetas = thetas - learning_rate*gradients
27     thetas_path.append(thetas)
28
29     mean_loss = np.sum(loss)/N
30     losses.append(mean_loss)
31
32     return thetas_path, losses

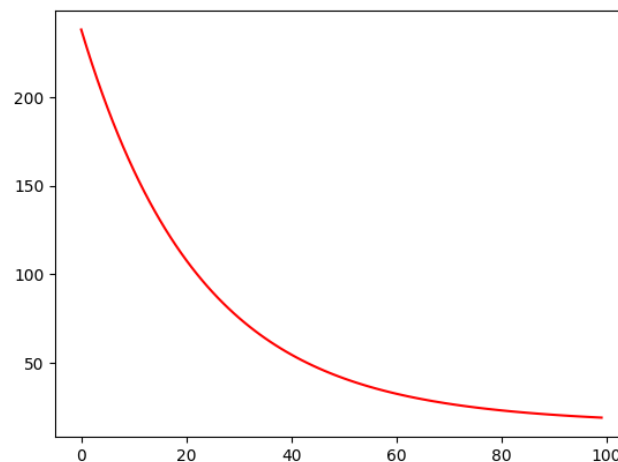
```

Hình 3 là kết quả sau khi thực thi đoạn code sau:

```

1 bgd_thetas, losses = batch_gradient_descent(X_b, y, n_epochs=100, learning_rate
    =0.01)
2
3 x_axis = list(range(100))
4 plt.plot(x_axis, losses[:100], color="r")
5 plt.show()

```



Hình 3: Kết quả loss values sử dụng batch Gradient Descent

**Question 14:** Kết quả của đoạn code sau đây:

```

1 bgd_thetas, losses = batch_gradient_descent(X_b, y, n_epochs=100, learning_rate
    =0.01)
2 print(round(sum(losses),2))

```

a) 7716.46

- b) 8716.46
- c) 6716.46
- d) 5716.46

**NOTE: YÊU CẦU CỦA ĐỀ BÀI LÀ PHẢI THỰC HIỆN THEO VECTORIZATION**

- **X\_b**: là thông tin Tv, Radio, Newspaper (thông tin model nhận vào và sử dụng để predict Sale) (đã được normalize)
- **y**: là thông tin Sale (thông tin mong muốn model dự đoán đúng)
- **n\_epochs**: Số lần train toàn bộ sample trong data
- **minibatch\_size**: Số lượng sample sẽ được train trong 1 step (Chỉ sử dụng ở câu b)
- **learning\_rate**: Tốc độ học
- **thetas\_path**: List weights của model từ lúc khởi tạo cho đến sau mỗi lần cập nhật weights
- **losses**: List loss của mỗi step sau khi cập nhật

**Bài tập 4** (Ưu nhược điểm của giải thuật linear regression): (Thảo luận và trao đổi trực tiếp trên lớp)  
Linear gression có thể ứng dụng cho bài toán phân loại không? nếu có gì như thế nào. Tìm hiểu về "Normal Equation in Linear Regression"