

LINEAR REGRESSION

Linear Regression

- Hồi quy tuyến tính đơn biến (Simple Linear Regression)
- Hồi quy tuyến tính đa biến (Multiple Linear Regression)

Hồi quy tuyến tính đơn biến

- Bài toán đặt ra:
 - Giả sử ta có tập dữ liệu:

Experience	Salary
3	60
4	55
5	66
6	93
7	?

- Hãy xây dựng một chương trình ML để tự động dự đoán tiền lương của nhân viên dựa trên số năm kinh nghiệm của họ.
- Câu hỏi đặt ra là, làm cách nào để xây dựng một chương trình có thể tự động dự đoán tiền lương dựa vào số năm kinh nghiệm?

Hồi quy tuyến tính đơn biến

- Quan sát bảng dữ liệu + Trực quan hóa bảng dữ liệu bằng biểu đồ scatter, ta thấy:
 - Khi kinh nghiệm tăng -> lương tăng
 - Các điểm dữ liệu thường tạo thành dạng gần thẳng

⇒ Mô hình phù hợp nhất: Hồi quy tuyến tính đơn biến (Simple Linear Regression)

⇒ Nghĩa là, phải xây dựng được công thức một đường thẳng có dạng:

$$\text{Salary} = W \cdot \text{Experience} + b$$

Experience	Salary
3	60
4	55
5	66
6	93
7	?

Hồi quy tuyến tính đơn biến

- Với tập dữ liệu đã có, ta có thể xây dựng các đường thẳng sau:
 - $\text{Salary}_1 = W_1 \cdot \text{Experience}_1 + b_1$
 - $\text{Salary}_2 = W_2 \cdot \text{Experience}_2 + b_2$
 - $\text{Salary}_3 = W_3 \cdot \text{Experience}_3 + b_3$
 - $\text{Salary}_4 = W_4 \cdot \text{Experience}_1 + b_4$
- Nghĩa là ta có:
 - $60 = 3w_1 + b_1 \Rightarrow w_1 = 16, b_1 = 12$
 $\Rightarrow \text{Salary} = 16 \cdot \text{Experience} + 12$
 - $55 = 4w_2 + b_2 \Rightarrow w_2 = 12, b_2 = 7$
 $\Rightarrow \text{Salary} = 12 \cdot \text{Experience} + 7$
 - $66 = 5w_3 + b_3 \Rightarrow w_3 = 12, b_3 = 6$
 $\Rightarrow \text{Salary} = 12 \cdot \text{Experience} + 6$
 - $93 = 6w_4 + b_4 \Rightarrow w_4 = 15, b_4 = 3$
 $\Rightarrow \text{Salary} = 15 \cdot \text{Experience} + 3$

Experience	Salary
3	60
4	55
5	66
6	93
7	?

Hồi quy tuyến tính đơn biến

- Nghĩa là ta có:

- $\text{Salary} = 16 \times \text{Experience} + 12$

- $\Rightarrow \text{Experience} = 7 \Rightarrow \text{Salary} = 16 \times 7 + 12 = 124$

- $\text{Salary} = 12 \times \text{Experience} + 7$

- $\Rightarrow \text{Experience} = 7 \Rightarrow \text{Salary} = 12 \times 7 + 7 = 91$

- $\text{Salary} = 12 \times \text{Experience} + 6$

- $\Rightarrow \text{Experience} = 7 \Rightarrow \text{Salary} = 12 \times 7 + 6 = 90$

- $\text{Salary} = 15 \times \text{Experience} + 3$

- $\Rightarrow \text{Experience} = 7 \Rightarrow \text{Salary} = 15 \times 7 + 3 = 108$

- Vậy rút cuộc đường thẳng nào cho câu trả lời phù hợp nhất?

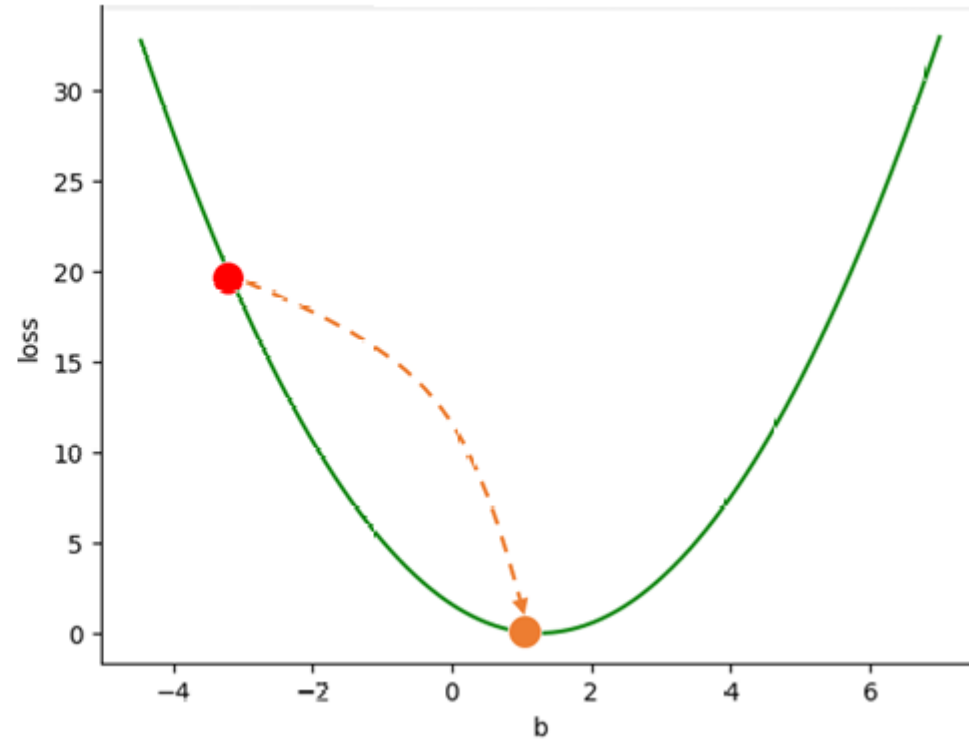
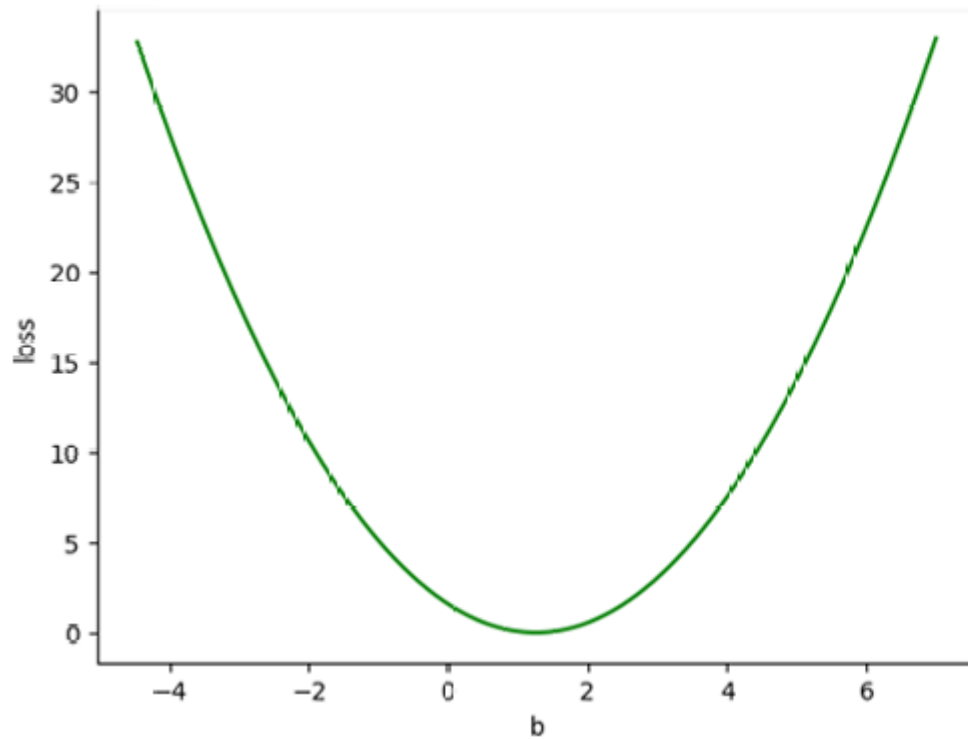
Experience	Salary
3	60
4	55
5	66
6	93
7	?

Hồi quy tuyến tính đơn biến

- Mục tiêu của mô hình Simple Linear Regression:
 - Tìm một đường thẳng tốt nhất mô tả mối liên hệ giữa một biến đầu vào (x) và một biến đầu ra (y), để từ đó dự đoán giá trị y cho giá trị x mới.
- Cụ thể:
 - Từ các (x,y) đã có \Rightarrow tìm w, b để xây dựng đường thẳng $\hat{y} = w \times x + b$ sao cho $L = (\hat{y} - y)^2$ nhỏ nhất

Hồi quy tuyến tính đơn biến

- Từ các phân tích, ta có một bài toán mới được phát biểu như sau: Tìm w và b để giá trị loss (L) là nhỏ nhất.



Hồi quy tuyến tính đơn biến: Quy trình huấn luyện

1. Khởi tạo ngẫu nhiên giá trị cho hai tham số là w và b .
2. Với mỗi mẫu dữ liệu thứ i trong bộ dữ liệu, ta áp dụng các bước tính toán sau:
 - (a) Thực hiện dự đoán output \hat{y}_i với input là x_i theo công thức sau:
$$\hat{y}_i = f(x_i) = w \times x_i + b$$
 - (b) Để tính toán sự chênh lệch giữa giá trị dự đoán so với giá trị thực tế, ta đưa y_i và \hat{y}_i vào hàm tính loss như sau:
$$L(\hat{y}_i, y_i) = (\hat{y}_i - y_i)^2$$

Hồi quy tuyến tính đơn biến

- (c) Ta tìm giá trị đạo hàm tại mẫu dữ liệu thứ i cho hai tham số w và b như sau:

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial b} = 2(\hat{y}_i - y_i), \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial w} = 2x_i(\hat{y}_i - y_i)$$

- (d) Với giá trị đạo hàm riêng vừa tìm được cho hai tham số w và b trên mẫu dữ liệu i , ta cập nhật lại giá trị mới cho hai tham số này với công thức sau:

$$w = w - \eta \frac{\partial L}{\partial w}, \quad b = b - \eta \frac{\partial L}{\partial b}$$

η là learning rate (tạm dịch: hệ số học).

- Sau đó, như đã đề cập, ta lặp lại bước 2 cho đến khi xử lý hết tất cả các mẫu dữ liệu trong bộ dữ liệu.

THỰC HÀNH

1. Khởi tạo $w = 10$, $b = 5$ và hệ số học $\eta = 0.01$.

2. Duyệt qua từng mẫu dữ liệu:

1. Mẫu 0:

a. $\hat{y}_0 = f(x_0) = 10 \times 3 + 5 = 35$

b. $L(\hat{y}_0, y_0) = (35 - 60)^2 = 625$

c. $\frac{\partial L}{\partial w} = 2x_0(\hat{y}_0 - y_0) = 2 \times 3(35 - 60) = -150$

$$\frac{\partial L}{\partial b} = 2(\hat{y}_0 - y_0) = 2(35 - 60) = -50$$

d. $w = w - \eta \frac{\partial L}{\partial w} = 10 - 0.01 \times (-150) = 11.5$

$$b = b - \eta \frac{\partial L}{\partial b} = 5 - 0.01 \times (-50) = 5.5$$

Sau mẫu 0: $w = 11.5$, $b = 5.5$ và hệ số học $\eta = 0.01$

(Dùng bộ này để tính mẫu 1)

Experience	Salary
3	60
4	55
5	66
6	93
7	?

VIẾT CHƯƠNG TRÌNH

Xây dựng mô hình

1. Khởi tạo w và b

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

```
1 INPUT_DIM = X_train.shape[1] # X is 1-dimensional
2 OUTPUT_DIM = y_train.shape[1] # y is 1-dimensional
```

```
1 # Initialize random weights
2 W = 0.01 * np.random.randn(INPUT_DIM, OUTPUT_DIM)
3 b = np.zeros((1, 1))
4 print (f"W: {W.shape}")
5 print (f"b: {b.shape}")
```

VIẾT CHƯƠNG TRÌNH

Xây dựng mô hình

2. Modeling

2.1. Tính \hat{y}_0

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

```
1 # Forward pass  $[NX1] \cdot [1X1] = [NX1]$ 
2 y_pred = np.dot(X_train, W) + b
```

2.2. Tính L

```
1 # Loss
2 N = len(y_train)
3 loss = (1/N) * np.sum((y_train - y_pred)**2)
4 print(f"loss: {loss:.2f}")
```

VIẾT CHƯƠNG TRÌNH

Xây dựng mô hình

2. Modeling

2.3. Tính gradient

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

```
1  # Backpropagation
2  dW = -(2/N) * np.sum((y_train - y_pred) * X_train)
3  db = -(2/N) * np.sum((y_train - y_pred) * 1)
```

2.4. Cập nhật w và b

```
1  LEARNING_RATE = 1e-1
```

```
1  # Update weights
2  W += -LEARNING_RATE * dW
3  b += -LEARNING_RATE * db
```

VIẾT CHƯƠNG TRÌNH

Huấn luyện mô hình theo cơ chế full-sample

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

```
1 | NUM_EPOCHS = 100
```

```
1 | # Initialize random weights
2 | W = 0.01 * np.random.randn(INPUT_DIM, OUTPUT_DIM)
3 | b = np.zeros((1, ))
4 |
5 | # Training loop
6 | for epoch_num in range(NUM_EPOCHS):
7 |
8 |     # Forward pass  $[NX1] \cdot [1X1] = [NX1]$ 
9 |     y_pred = np.dot(X_train, W) + b
10 |
11 |     # Loss
12 |     loss = (1/len(y_train)) * np.sum((y_train - y_pred)**2)
13 |
14 |     # Show progress
15 |     if epoch_num%10 == 0:
16 |         print(f"Epoch: {epoch_num}, loss: {loss:.3f}")
17 |
18 |     # Backpropagation
19 |     dW = -(2/N) * np.sum((y_train - y_pred) * X_train)
20 |     db = -(2/N) * np.sum((y_train - y_pred) * 1)
21 |
22 |     # Update weights
23 |     W += -LEARNING_RATE * dW
24 |     b += -LEARNING_RATE * db
```

VIẾT CHƯƠNG TRÌNH

Xây dựng mô hình

4. Testing (Evaluation)

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

```
1  # Predictions
2  pred_train = W*X_train + b
3  pred_test = W*X_test + b
```

```
1  # Train and test MSE
2  train_mse = np.mean((y_train - pred_train) ** 2)
3  test_mse = np.mean((y_test - pred_test) ** 2)
4  print (f"train_MSE: {train_mse:.2f}, test_MSE: {test_mse:.2f}")
```


VIẾT CHƯƠNG TRÌNH

Huấn luyện mô hình theo cơ chế one-sample

```
w, b = initialize_params()
N = len(y_train)
LEARNING_RATE = 1e-2
epoch_max = 100

for epoch in range(epoch_max):
    for i in range(N):
        # get a sample
        x = X_train[i]          # vector shape (d,)
        y = y_train[i]

        # compute output
        y_pred = np.dot(w, x) + b

        # compute gradients
        dw = 2*(y_pred - y) * x
        db = 2*(y_pred - y)

        # update parameters
        w -= LEARNING_RATE * dw
        b -= LEARNING_RATE * db

    # print progress
    y_pred_all = X_train.dot(w) + b
    loss = np.mean((y_pred_all - y_train)**2)

    if epoch % 10 == 0:
        print(f"Epoch {epoch}, loss = {loss:.4f}")
```

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

Hồi quy tuyến tính đa biến

- Bài toán đặt ra:
 - Giả sử ta có tập dữ liệu:

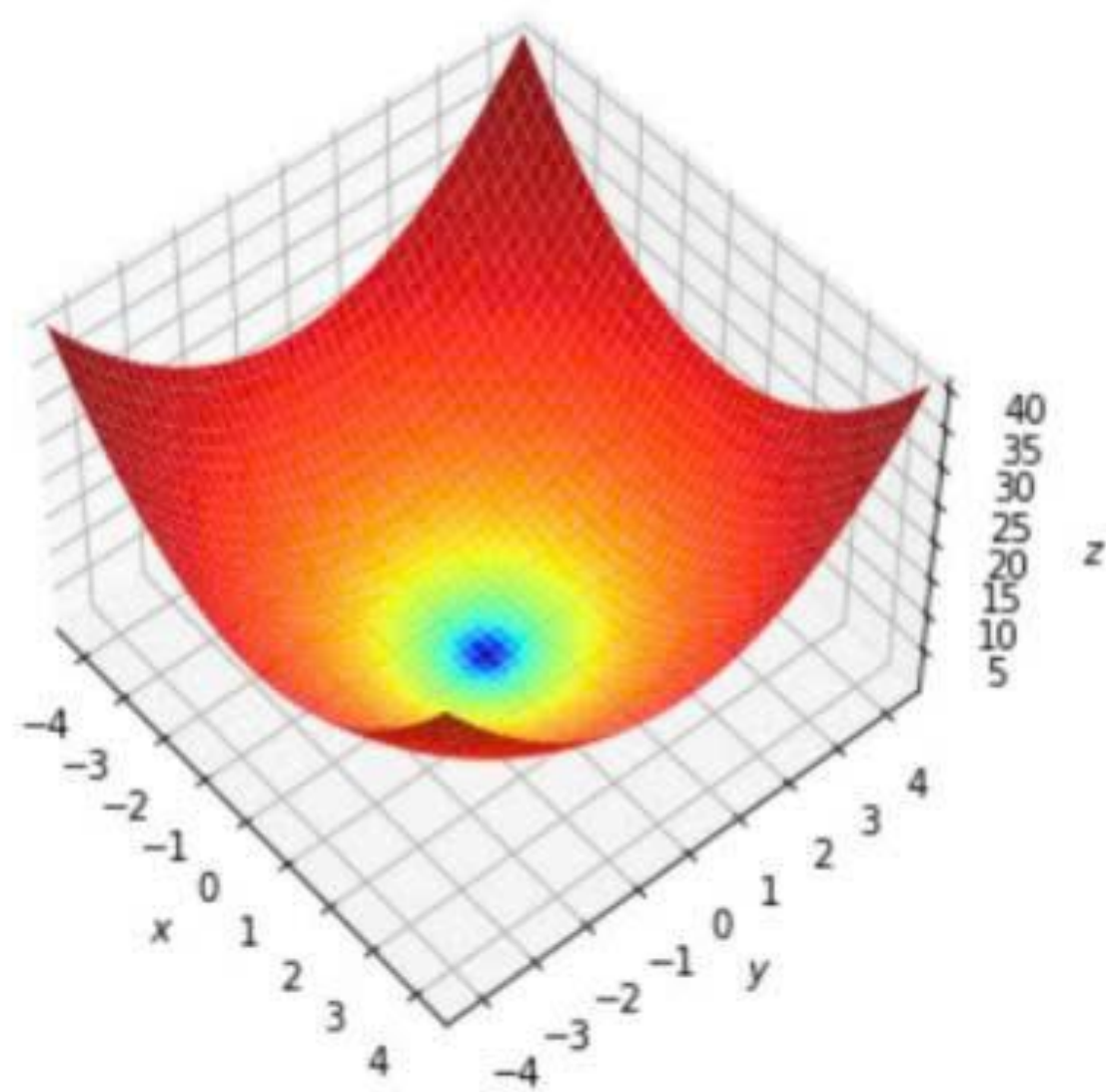
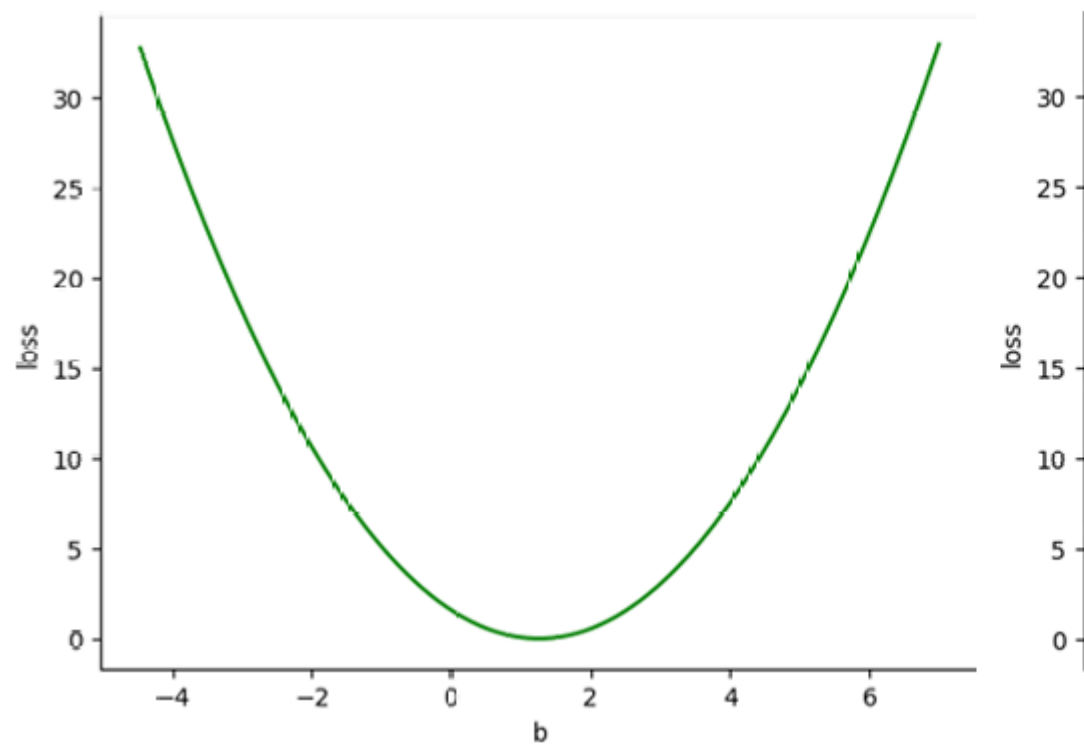
Experience	Degree	Salary
3	1	60
4	1	55
5	2	66
6	2	93
7	?	?

- Hãy xây dựng một chương trình ML để tự động dự đoán tiền lương của nhân viên dựa trên số năm kinh nghiệm và bậc lương của họ.
- Câu hỏi đặt ra là, làm cách nào để xây dựng một chương trình có thể tự động dự đoán tiền lương dựa vào số năm kinh nghiệm và bậc lương?

Hồi quy tuyến tính đa biến

- Mục tiêu của mô hình Multiple Linear Regression:
 - Tìm một mặt phẳng tốt nhất mô tả mối liên hệ giữa các biến đầu vào (x_1, x_2) và một biến đầu ra (y), để từ đó dự đoán giá trị y cho các giá trị x mới.
- Cụ thể:
 - Từ các (x, y) đã có \Rightarrow tìm w_1, w_2, b để xây dựng mặt phẳng $\hat{y} = w_1x_1 + w_2x_2 + b$ sao cho $L = (\hat{y} - y)^2$ nhỏ nhất
 - Hoặc viết gọn hơn: $\hat{y} = w^T x + b$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix}$$



Hồi quy tuyến tính đa biến: Quy trình huấn luyện

1. Khởi tạo ngẫu nhiên giá trị cho các tham số là w_1, w_2, b .

2. Lặp qua Dữ liệu: Với mỗi mẫu dữ liệu (x_i, y_i)

- Dự đoán output \hat{y}_i theo công thức sau: $\hat{y}_i = w^T x_i + b$ $\hat{y}_i = \sum_j w_j x_{ij} + b$

- Tính loss như sau: $L(\hat{y}_i, y_i) = (\hat{y}_i - y_i)^2$

- Tính đạo hàm:

$$\frac{\partial L_i}{\partial w_j} = 2x_{ij}(\hat{y}_i - y_i)$$

$$\frac{\partial L_i}{\partial b} = 2(\hat{y}_i - y_i)$$

- Cập nhật: w_1, w_2, b $w_j = w_j - \eta \frac{\partial L_i}{\partial w_j}, \quad b = b - \eta \frac{\partial L_i}{\partial b}$

3. Lặp lại bước 2 qua nhiều lần (epochs) cho đến khi loss đạt mức tối ưu (giá trị nhỏ nhất).

THỰC HÀNH

1. Khởi tạo giá trị cho $w_1 = 10, w_2 = 2, b = 5$ và hệ số học $\eta = 0.01$.
2. Duyệt qua từng mẫu dữ liệu:

1. Mẫu 0: $x_{01}=3, x_{02}=1, y_0=60$

$$a. \hat{y}_0 = f(x_{0j}) = w_1 x_{01} + w_2 x_{02} + b \\ = 10 \times 3 + 2 \times 1 + 5 = 37$$

$$b. L(\hat{y}_0, y_0) = (60-37)^2 = 529$$

Experience	Degree	Salary
3	1	60
4	1	55
5	2	66
6	2	93
7	?	?

THỰC HÀNH

1. Khởi tạo giá trị cho $w_1 = 10, w_2 = 2, b = 5$ và hệ số học $\eta = 0.01$.
2. Duyệt qua từng mẫu dữ liệu:

1. Mẫu 0: $x_{01}=3, x_{02}=1, y_0=60$

c. Đạo hàm (theo từng tham số)

$$\frac{\partial L_0}{\partial w_1} = 2x_{01}(\widehat{y_0} - y_0) = 2 \times 3(60 - 37) = -138$$

$$\frac{\partial L_0}{\partial w_2} = 2x_{02}(\widehat{y_0} - y_0) = 2 \times 1(60 - 37) = -46$$

$$\frac{\partial L_0}{\partial b} = 2(\widehat{y_0} - y_0) = 2(60 - 37) = -46$$

Experience	Degree	Salary
3	1	60
4	1	55
5	2	66
6	2	93
7	?	?

THỰC HÀNH

1. Khởi tạo giá trị cho $w_1 = 10, w_2 = 2, b = 5$ và hệ số học $\eta = 0.01$.
2. Duyệt qua từng mẫu dữ liệu:

1. Mẫu 0: $x_{01}=3, x_{02}=1, y_0=60$

d. Cập nhật

$$w_1 = w_1 - \eta \frac{\partial L_0}{\partial w_1} = 10 - 0.01 \times (-138) = 11.38$$

$$w_2 = w_2 - \eta \frac{\partial L_0}{\partial w_2} = 2 - 0.01 \times (-46) = 2.46$$

$$b = b - \eta \frac{\partial L_0}{\partial b} = 5 - 0.01 \times (-46) = 5.46$$

Experience	Degree	Salary
3	1	60
4	1	55
5	2	66
6	2	93
7	?	?

Sau mẫu 0: $w_1 = 11.38, w_2 = 2.46, b = 5.46$ và hệ số học $\eta = 0.01$
(Dùng bộ này để tính mẫu 1)

VIẾT CHƯƠNG TRÌNH

Huấn luyện mô hình theo cơ chế one-sample

```
# Initialize
w1, w2, w3, b = initialize_params()
LEARNING_RATE = 1e-2
epoch_max = 100
N = len(y_train)
for epoch in range(epoch_max):
    # Shuffle data
    indices = np.random.permutation(N)

    for i in indices:
        # Get one sample
        x1 = X_train[i][0]
        x2 = X_train[i][1]
        x3 = X_train[i][2]
        y = y_train[i]

        # compute output
        y_pred = w1*x1 + w2*x2 + w3*x3 + b

        # Compute gradients for 1 sample
        dw1 = 2*(y_pred - y)*x1
        dw2 = 2*(y_pred - y)*x2
        dw3 = 2*(y_pred - y)*x3
        db = 2*(y_pred - y)

        # update parameters
        w1 -= LEARNING_RATE * dw1
        w2 -= LEARNING_RATE * dw2
        w3 -= LEARNING_RATE * dw3
        b -= LEARNING_RATE * db

    # Compute full MSE loss after epoch
    y_pred_all = w1*X_train[:,0] + w2*X_train[:,1] + w3*X_train[:,2] + b
    loss = np.mean((y_pred_all - y_train)**2)

    # Print progress
    if epoch % 10 == 0:
        print(f"Epoch {epoch}, loss = {loss:.4f}")
```

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

VIẾT CHƯƠNG TRÌNH

Huấn luyện mô hình theo cơ chế full-sample

```
w1, w2, w3, b = initialize_params()
N = len(y_train)
LEARNING_RATE = 1e-2
epoch_max = 100

for epoch in range(epoch_max):
    # compute output
    y_pred = w1*x1 + w2*x2 + w3*x3 + b

    # compute loss (MSE)
    loss = (1/N) * np.sum((y_train - y_pred)**2)

    # Show progress
    if epoch % 10 == 0:
        print(f"Epoch: {epoch}, loss: {loss:.3f}")

    # compute gradients
    dw1 = (2/N) * np.sum((y_pred - y_train) * x1)
    dw2 = (2/N) * np.sum((y_pred - y_train) * x2)
    dw3 = (2/N) * np.sum((y_pred - y_train) * x3)
    db = (2/N) * np.sum((y_pred - y_train))

    # update parameters
    w1 -= LEARNING_RATE * dw1
    w2 -= LEARNING_RATE * dw2
    w3 -= LEARNING_RATE * dw3
    b -= LEARNING_RATE * db
```

Đọc và hoàn thiện các đoạn code sau để có một chương trình huấn luyện mô hình hoàn chỉnh

Ridge regression và Lasso regression

1. Ridge regression:

$$\mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha ||\mathbf{w}||_2^2$$

2. Lasso regression:

$$\mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha ||\mathbf{w}||_1$$

Experience	Degree	Salary
3	1	60
4	1	55
5	2	66
6	2	93
7	3	?

Experience	Degree	Salary
3	1	60
4	1	55
5	2	66
6	2	93
7	1	48
8	1	51
7	?	?