

AI VIET NAM

@aivietnam.edu.vn

Vectorized Implementation for Linear Regression (Naive and Curious Approach)

Quang-Vinh Dinh
Ph.D. in Computer Science

❖ Linear regression

| Feature | | Label | |
|---------|------|-------|--|
| | | | |
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

House price data

$$\text{Model: } \hat{y} = w_1 x_1 + b$$

$$\text{price} = a * \text{area} + b$$

| Features | | | Label |
|----------|---------|-------------|---------|
| | | | |
| TV | ↕ Radio | ↕ Newspaper | ↕ Sales |
| 230.1 | 37.8 | 69.2 | 22.1 |
| 44.5 | 39.3 | 45.1 | 10.4 |
| 17.2 | 45.9 | 69.3 | 12 |
| 151.5 | 41.3 | 58.5 | 16.5 |
| 180.8 | 10.8 | 58.4 | 17.9 |

Advertising data

$$\text{Model: } \hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$\text{Sale} = w_1 * TV + w_2 * Radio + w_3 * Newspaper + b$$

Linear Regression

1) Pick a sample (x_1, x_2, x_3, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = w_1 * TV + w_2 * R + w_3 * N + b$$

$$\hat{y} = w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w_1} = 2x_1(\hat{y} - y) \quad \frac{\partial L}{\partial w_3} = 2x_3(\hat{y} - y)$$

$$\frac{\partial L}{\partial w_2} = 2x_2(\hat{y} - y) \quad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w_1 = w_1 - \eta \frac{\partial L}{\partial w_1} \quad w_3 = w_3 - \eta \frac{\partial L}{\partial w_3}$$

$$w_2 = w_2 - \eta \frac{\partial L}{\partial w_2} \quad b = b - \eta \frac{\partial L}{\partial b}$$

| Features | | | Label |
|----------|---------|-------------|---------|
| TV | ↕ Radio | ↕ Newspaper | ↕ Sales |
| 230.1 | 37.8 | 69.2 | 22.1 |
| 44.5 | 39.3 | 45.1 | 10.4 |
| 17.2 | 45.9 | 69.3 | 12 |
| 151.5 | 41.3 | 58.5 | 16.5 |
| 180.8 | 10.8 | 58.4 | 17.9 |

Advertising data

Model

$$\text{Sale} = w_1 * TV + w_2 * Radio + w_3 * Newspaper + b$$

$$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

1) Pick a sample (x_1, x_2, x_3, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = w_1 * TV + w_2 * R + w_3 * N + b$$

$$\hat{y} = w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w_1} = 2x_1(\hat{y} - y) \quad \frac{\partial L}{\partial w_3} = 2x_3(\hat{y} - y)$$

$$\frac{\partial L}{\partial w_2} = 2x_2(\hat{y} - y) \quad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w_1 = w_1 - \eta \frac{\partial L}{\partial w_1} \quad w_3 = w_3 - \eta \frac{\partial L}{\partial w_3}$$

$$w_2 = w_2 - \eta \frac{\partial L}{\partial w_2} \quad b = b - \eta \frac{\partial L}{\partial b}$$

```
1  # compute output and loss
2  def predict(x1, x2, x3, w1, w2, w3, b):
3      return w1*x1 + w2*x2 + w3*x3 + b
4  def compute_loss(y_hat, y):
5      return (y_hat - y)**2
6
7  # compute gradient
8  def compute_gradient_wi(xi, y, y_hat):
9      dl_dwi = 2*xi*(y_hat-y)
10     return dl_dwi
11 def compute_gradient_b(y, y_hat):
12     dl_db = 2*(y_hat-y)
13     return dl_db
14
15 # update weights
16 def update_weight_wi(wi, dl_dwi, lr):
17     wi = wi - lr*dl_dwi
18     return wi
19 def update_weight_b(b, dl_db, lr):
20     b = b - lr*dl_db
21     return b
```

Motivation

| Feature | | Label | |
|---------|------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

House price data

| Features | | | Label |
|----------|-------|-----------|-------|
| TV | Radio | Newspaper | Sales |
| 230.1 | 37.8 | 69.2 | 22.1 |
| 44.5 | 39.3 | 45.1 | 10.4 |
| 17.2 | 45.9 | 69.3 | 12 |
| 151.5 | 41.3 | 58.5 | 16.5 |
| 180.8 | 10.8 | 58.4 | 17.9 |

Advertising data

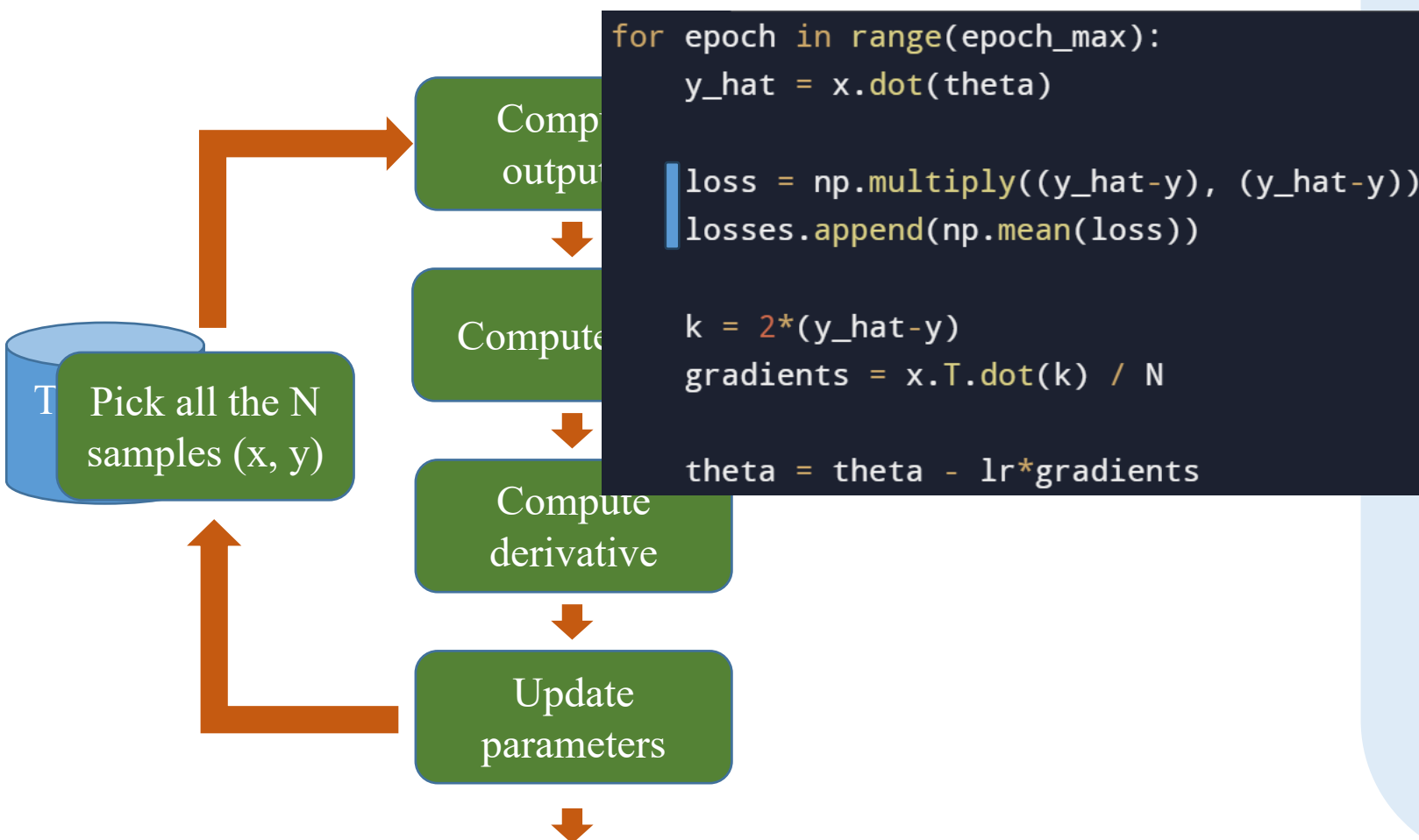
if area=6.0, price=?

if TV=55.0, Radio=34.0,
and Newspaper=62.0,
price=?

| Features | | | | | | | | | | | | | | Label |
|----------|------|-------|------|-------|-------|------|--------|-----|-----|---------|--------|-------|------|-------|
| crim | zn | indus | chas | nox | rm | age | dis | rad | tax | ptratio | black | lstat | medv | |
| 0.00632 | 18 | 2.31 | 0 | 0.538 | 6.575 | 65.2 | 4.09 | 1 | 296 | 15.3 | 396.9 | 4.98 | 24 | |
| 0.02731 | 0 | 7.07 | 0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2 | 242 | 17.8 | 396.9 | 9.14 | 21.6 | |
| 0.03237 | 0 | 2.18 | 0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3 | 222 | 18.7 | 394.63 | 2.94 | 33.4 | |
| 0.06905 | 0 | 2.18 | 0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3 | 222 | 18.7 | 396.9 | 5.33 | 36.2 | |
| 0.08829 | 12.5 | 7.87 | 0 | 0.524 | 6.012 | 66.6 | 5.5605 | 5 | 311 | 15.2 | 395.6 | 12.43 | 22.9 | |

Boston House Price Data

- ❖ Using vector and matrix for notation
- ❖ Using Numpy for implementation



1) Pick all the N samples from training data

2) Compute output \hat{y}

$$\hat{y} = X\theta$$

3) Compute loss

$$L = (\hat{y} - y)(\hat{y} - y)^T \frac{1}{N}$$

4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = X^T k$$

5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{N}$$

Outline

SECTION 1

1-sample Vectorization

SECTION 2

m-sample Vectorization

SECTION 3

N-sample Vectorization

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |
| x | y |

$$\hat{y} = wx + b \quad \cancel{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

✓

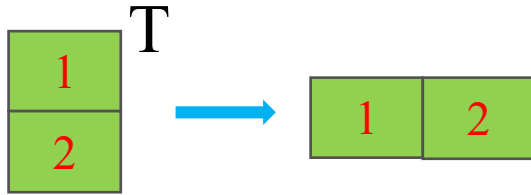
$$X \cdot \theta^T$$

Linear Regression (1-samples)

Transpose

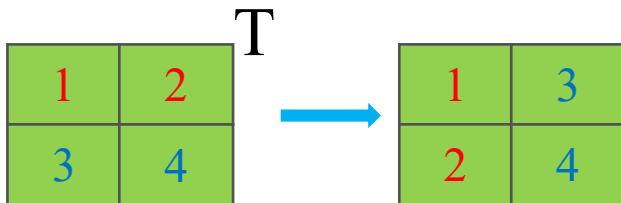
$$\vec{v} = \begin{bmatrix} v_1 \\ \dots \\ v_n \end{bmatrix}$$

$$\vec{v}^T = [v_1 \ \dots \ v_n]$$



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

$$A^T = \begin{bmatrix} a_{11} & \dots & a_{m1} \\ \dots & \dots & \dots \\ a_{1n} & \dots & a_{mn} \end{bmatrix}$$

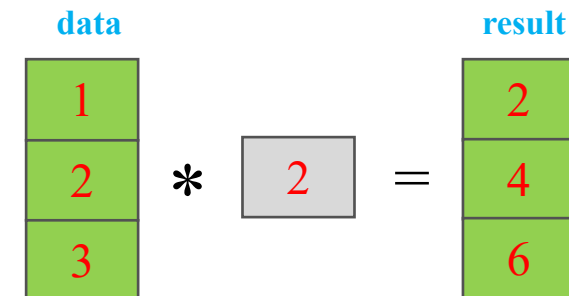


```
2 import numpy as np
3
4 # create data
5 data = np.array([1,2,3])
6 factor = 2
7
8 # broadcasting
9 result_multiplication = data*factor
```

```
[1 2 3]
[2 4 6]
```

Multiply with a number

$$\alpha \vec{u} = \alpha \begin{bmatrix} u_1 \\ \dots \\ u_n \end{bmatrix} = \begin{bmatrix} \alpha u_1 \\ \dots \\ \alpha u_n \end{bmatrix}$$



Linear Regression (1-samples)

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |
| x | y |

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y)$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Traditional

$$\hat{y} = wx + b \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix} \Rightarrow \theta^T = [b \ w]$$

$$\hat{y} = wx + b1 = \begin{bmatrix} b & w \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix} = \theta^T x$$

dot product

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y)$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Traditional

$$\hat{y} = wx + b \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\hat{y} = \theta^T x$$

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

numbers

What will we do?

Linear Regression (1-samples)

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

Traditional

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y) \quad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

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

η is learning rate

$$\hat{y} = wx + b \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y) = 2 \times (\hat{y} - y) \times 1$$

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y) = 2 \times (\hat{y} - y) \times x$$

$$2 \times (\hat{y} - y) \times x$$

$$\begin{bmatrix} 2 \times (\hat{y} - y) \times 1 \\ 2 \times (\hat{y} - y) \times x \end{bmatrix} = 2(\hat{y} - y) \begin{bmatrix} 1 \\ x \end{bmatrix} = 2(\hat{y} - y) \mathbf{x} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix} = L'_{\theta}$$

common factor

$$\rightarrow L'_{\theta} = 2(\hat{y} - y)$$

Linear Regression (1-samples)

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y) \qquad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$\underline{w = w - \eta \frac{\partial L}{\partial w} \qquad b = b - \eta \frac{\partial L}{\partial b}}$$

Traditional

$$\hat{y} = wx + b$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$L'_{\boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix}$$

$$\left\{ \begin{array}{l} b = b - \eta \frac{\partial L}{\partial b} \\ w = w - \eta \frac{\partial L}{\partial w} \\ \boldsymbol{\theta} \quad \boldsymbol{\theta} \quad L'_{\boldsymbol{\theta}} \end{array} \right.$$

$$\rightarrow \boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

Linear Regression (1-samples)

6

1) Pick a sample (x, y) from training data

2) Compute the output \hat{y}

$$\hat{y} = wx + b$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y)$$

$$\frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Traditional

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \theta^T x = x^T \theta$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$L'_\theta = 2x(\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta L'_\theta$$

η is learning rate

Vectorized

$$\hat{y} = wx + b$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |
| x | y |

1

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} = \begin{bmatrix} 1 \\ 6.7 \end{bmatrix}$$

Given $\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$

$$\eta = 0.01$$

1) Pick a sample (\mathbf{x}, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

η is learning rate

Input

\mathbf{x}

Model

$$\boldsymbol{\theta} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix}$$

2

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = -2.238$$

Label

$$y = 9.1$$

Loss

3

$$(\hat{y} - y)^2 = 128.5$$

update

4

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y) = \begin{bmatrix} -22.658 \\ -151.81 \end{bmatrix}$$

5

$$\boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}} = \begin{bmatrix} 0.049 \\ -0.34 \end{bmatrix} - 0.01 \begin{bmatrix} -22.658 \\ -151.81 \end{bmatrix} = \begin{bmatrix} 0.2755 \\ 1.1781 \end{bmatrix}$$

❖ Implementation (vectorization using numpy)

1) Pick a sample (x, y) from training data

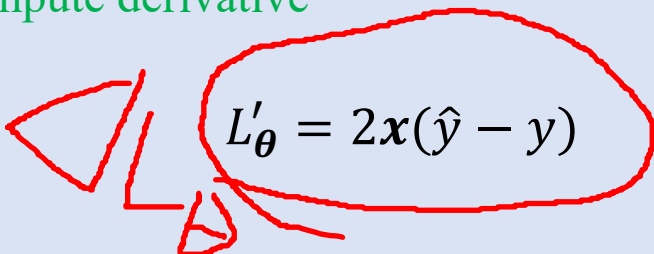
2) Compute output \hat{y}

$$\hat{y} = \theta^T x = x^T \theta$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative


$$L'_\theta = 2x(\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta L'_\theta$$

η is learning rate

```
1 import numpy as np
2
3 # forward
4 def predict(x, theta):
5     return x.dot(theta)
6
7 # compute gradient
8 def gradient(y_hat, y, x):
9     dtheta = 2*x*(y_hat-y)
10
11     return dtheta
12
13 # update weights
14 def update_weight(theta, lr, dtheta):
15     dtheta_new = theta - lr*dtheta
16
17     return dtheta_new
```

$$\hat{y} = wx + b$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |
| x | y |

1

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} = \begin{bmatrix} 1 \\ 3.5 \end{bmatrix}$$

Given $\boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$

$$\eta = 0.01$$

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

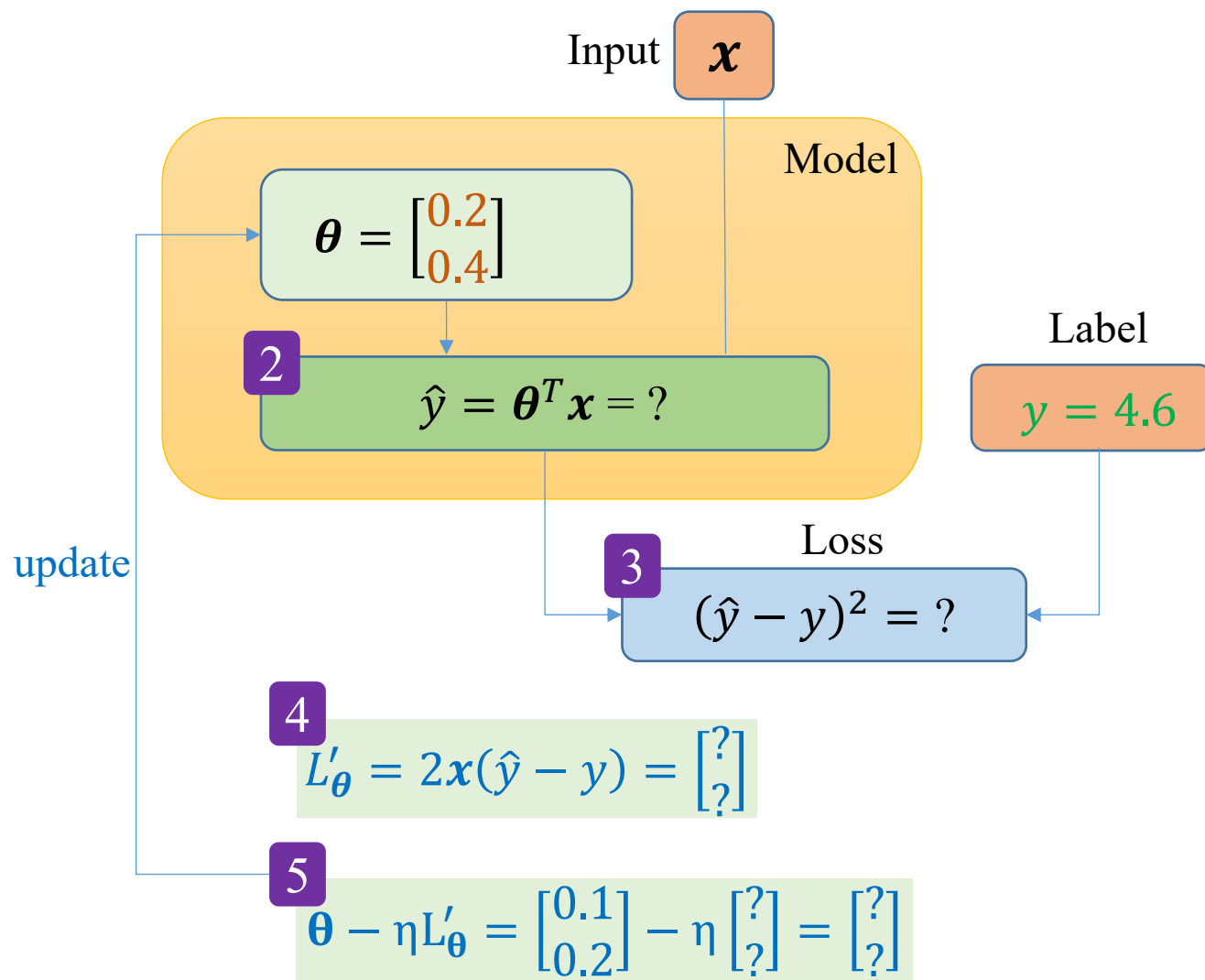
4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta L'_{\boldsymbol{\theta}}$$

η is learning rate




```
3 data = np.genfromtxt('data.csv', delimiter=',')
4 N = 4
5
6 areas = data[:, 0].reshape(N, 1)
7 prices = data[:, 1].reshape(N,)
8
9 # vector [1, area]
10 features = np.hstack([np.ones((N,1)), areas])
```

```
1 # forward
2 def predict(x, theta):
3     return x.T.dot(theta)
4
5 # compute gradient
6 def gradient(y_hat, y, x):
7     dtheta = 2*x*(y_hat-y)
8     return dtheta
9
10 # update weights
11 def update_weight(theta, lr, dtheta):
12     dtheta_new = theta - lr*dtheta
13     return dtheta_new
```

```
1 lr = 0.01
2 epoch_max = 10
3
4 # [b, w]
5 theta = np.array([0.049, -0.34])
6
7 for epoch in range(epoch_max):
8     for i in range(N):
9         # get a sample
10         x = features[i,:]
11         y = prices[i]
12
13         # predict y_hat
14         y_hat = predict(x, theta)
15
16         # compute loss
17         loss = (y_hat-y)*(y_hat-y)
18
19         # compute gradient
20         dtheta = gradient(y_hat, y, x)
21
22         # update weights
23         theta = update_weight(theta, lr, dtheta)
```

Advertising Problem

| Features | | | Label |
|----------|---------|-------------|---------|
| TV | ↕ Radio | ↕ Newspaper | ↕ Sales |
| 230.1 | 37.8 | 69.2 | 22.1 |
| 44.5 | 39.3 | 45.1 | 10.4 |
| 17.2 | 45.9 | 69.3 | 12 |
| 151.5 | 41.3 | 58.5 | 16.5 |
| 180.8 | 10.8 | 58.4 | 17.9 |

Advertising data

if TV=55.0, Radio=34.0,
and Newspaper=62.0,
price=?

$$\hat{y} = w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + b$$

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \theta^T x = x^T \theta$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

4) Compute derivative

$$L'_\theta = 2x(\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta L'_\theta$$

η is learning rate

Outline

SECTION 1

1-sample Vectorization

Parameter Initialization

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

SECTION 2

m-sample Vectorization

Way 1 for constructing matrix \mathbf{x}

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

SECTION 3

N-sample Vectorization

Way 2 for constructing matrix \mathbf{x}

$$\mathbf{X} = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix}$$

❖ Vectorization

| | Feature | Label |
|------------------|---------|-------|
| | area | price |
| House price data | 6.7 | 9.1 |
| | 4.6 | 5.9 |
| | 3.5 | 4.6 |
| | 5.5 | 6.7 |

Model

price = w * area + b

$\hat{y} = wx + b$

$$y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

$X^T = \begin{bmatrix} 1 & 1 \end{bmatrix}$

Parameter Initialization

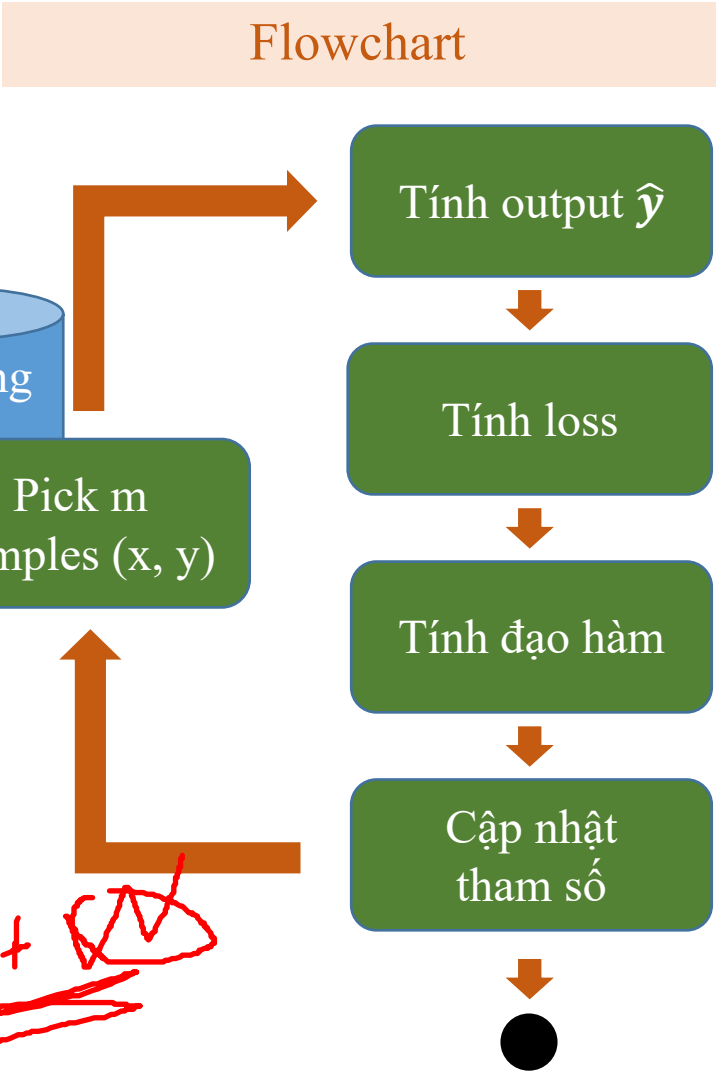
$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$

Way 1 for constructing matrix x

$X = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$

Way 2 for constructing matrix x

$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix}$

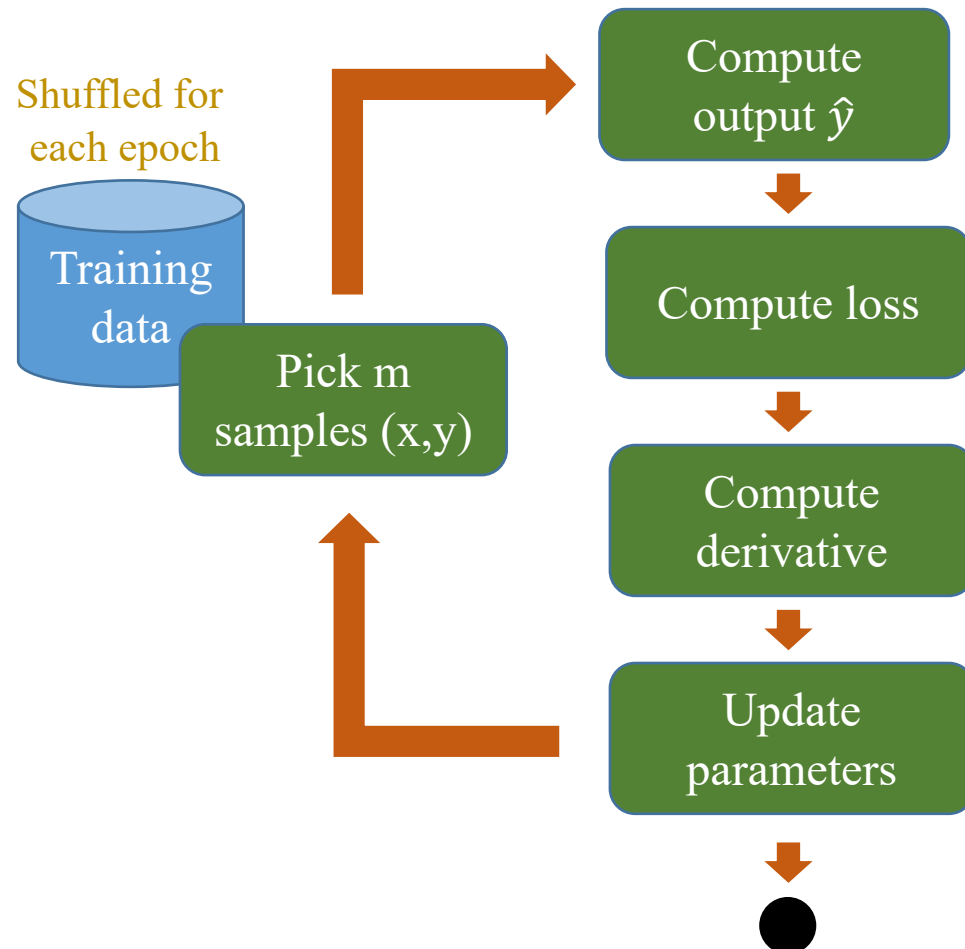


Linear Regression (m-samples)

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |

❖ House price prediction

❖ m-sample training ($1 < m < N$)



1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = wx^{(i)} + b$$

for $0 \leq i < m$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$\begin{cases} L'_w{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \\ L'_b{}^{(i)} = 2(\hat{y}^{(i)} - y^{(i)}) \end{cases}$$

for $0 \leq i < m$

5) Update parameters

$$w = w - \eta \frac{\sum_i L'_w{}^{(i)}}{m}$$

$$b = b - \eta \frac{\sum_i L'_b{}^{(i)}}{m}$$

Learning rate η

Linear Regression (m-samples)

| Feature | | Label |
|---------|-------|-------|
| area | price | |
| 6.7 | 9.1 | |
| 4.6 | 5.9 | |

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = wx^{(i)} + b \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_w{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)})$$
$$L'_b{}^{(i)} = 2(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$w = w - \eta \frac{\sum_i L'_w{}^{(i)}}{m}$$
$$b = b - \eta \frac{\sum_i L'_b{}^{(i)}}{m} \quad \text{Learning rate } \eta$$

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \theta^T x^{(i)} = (x^{(i)})^T \theta \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_\theta{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\theta = \theta - \eta \frac{\sum_i L'_\theta{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

Implementation - Differences

| | Feature | Label | |
|--|---------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |

1) Pick a sample (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = \theta^T x = x^T \theta$$

3) Compute loss

$$L = (\hat{y} - y)^2$$

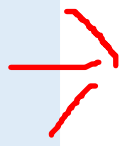
4) Compute derivative

$$L'_{\theta} = 2x(\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta L'_{\theta}$$

η is learning rate



1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \theta^T x^{(i)} = (x^{(i)})^T \theta \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'^{(i)}_{\theta} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\theta = \theta - \eta \frac{\sum_i L'^{(i)}_{\theta}}{m} \quad \eta \text{ is learning rate}$$

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \theta^T x^{(i)} = (x^{(i)})^T \theta \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_\theta = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\theta = \theta - \eta \frac{\sum_i L'_\theta}{m} \quad \eta \text{ is learning rate}$$

```
1 # vector [x, b]
2 data = np.c_[areas, np.ones((data_size, 1))]
3
4 # init weight
5 lr = 0.01
6 theta = np.array([-0.34, 0.04]) #[w, b]
7
8 # number of epochs
9 epoch_max = 10
10
11 # mini-batch size
12 m = 2
```


Linear Regression

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \theta^T x^{(i)} = (x^{(i)})^T \theta \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\theta}{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\theta = \theta - \eta \frac{\sum_i L'_{\theta}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

```
14 for epoch in range(epoch_max):
15     for j in range(0, data_size, m):
16
17         # some variables
18         sum_of_losses = 0
19         gradients = np.zeros((2,))
20         for index in range(j, j+m):
21             # get mini-batch
22             x_i = data[index]
23             y_i = prices[index]
24
25             # predict y_hat_i
26             y_hat_i = x_i.dot(theta)
27
28             # compute loss
29             l_i = (y_hat_i - y_i)*(y_hat_i - y_i)
30
31             # compute gradient
32             gradient_i = x_i*2*(y_hat_i - y_i)
33
34             # accumulate gradients
35             gradients = gradients + gradient_i
36             sum_of_losses = sum_of_losses + l_i
37
38         # normalize
39         sum_of_losses = sum_of_losses/2
40         gradients = gradients/2
```

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \theta^T x^{(i)} = (x^{(i)})^T \theta \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\theta}{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\theta = \theta - \eta \frac{\sum_i L'_{\theta}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

More vectorization

| | Feature | Label | |
|--|---------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

way 1

$$X = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$X = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T X$$

$$= [-0.34 \quad 0.049] \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} = [-2.238 \quad -1.524]$$

| Feature | | Label | |
|---------|--|-------|--|
| area | | price | |
| 6.7 | | 9.1 | |
| 4.6 | | 5.9 | |
| 3.5 | | 4.6 | |
| 5.5 | | 6.7 | |
| | | | |

Linear Regression (m-samples)

Feature Label

| area | price |
|------|-------|
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

Way 1.1

1) Pick m samples $(\mathbf{x}^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$X = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T X = [-2.238 \quad -1.524]$$

$$L = \frac{1}{m} [(\hat{y}^{(0)} - y^{(0)})^2 + (\hat{y}^{(1)} - y^{(1)})^2]$$

$$= \frac{1}{m} [(\hat{y}^{(0)} - y^{(0)}) \quad (\hat{y}^{(1)} - y^{(1)})] \begin{bmatrix} (\hat{y}^{(0)} - y^{(0)}) \\ (\hat{y}^{(1)} - y^{(1)}) \end{bmatrix}$$

$$= \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$= \frac{1}{2} (128.5 + 55.11) = 91.8$$

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

$$= \begin{bmatrix} -22.676 & -14.848 \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

| Feature | | Label | |
|---------|------|-------|-------|
| | area | | price |
| | 6.7 | | 9.1 |
| | 4.6 | | 5.9 |
| | 3.5 | | 4.6 |
| | 5.5 | | 6.7 |
| | | | |

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

$$= \begin{bmatrix} -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} = \begin{bmatrix} -22.676 & -14.848 \\ -22.676 & -14.848 \end{bmatrix}$$

| Feature | | Label | |
|---------|--|-------|--|
| area | | price | |
| 6.7 | | 9.1 | |
| 4.6 | | 5.9 | |
| 3.5 | | 4.6 | |
| 5.5 | | 6.7 | |
| | | | |

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

Take a look
and improve

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T) = 2(\hat{\mathbf{y}} - \mathbf{y}^T) \\ = \begin{bmatrix} -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} = \begin{bmatrix} -22.676 & -14.848 \\ -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix}$$

Gradient for w from $x^{(i)}$

Gradient for b from $x^{(i)}$

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

Take a look
and improve

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\begin{aligned} \mathbf{k} &= 2(\hat{\mathbf{y}} - \mathbf{y}^T) = 2(\hat{\mathbf{y}} - \mathbf{y}^T) \\ &= \begin{bmatrix} -22.676 & -14.848 \end{bmatrix} \end{aligned}$$

$$\mathbf{x} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} = \begin{bmatrix} -22.676 & -14.848 \\ -22.676 & -14.848 \end{bmatrix}$$

Take a look
and improve

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix}$$

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix}$$

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix}$$

| Feature | | Label | |
|---------|------|-------|-------|
| | area | | price |
| | 6.7 | | 9.1 |
| | 4.6 | | 5.9 |
| | 3.5 | | 4.6 |
| | 5.5 | | 6.7 |
| | | | |

Linear Regression (m-samples)

Way 1.1

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}{}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix}$$

$$\eta = 0.01$$

$$m = 2$$

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m}$$

$$= \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix} - 0.005 \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix} = \begin{bmatrix} 0.761 \\ 0.227 \end{bmatrix}$$

| Feature | | Label | |
|---------|------|-------|-------|
| | area | | price |
| | 6.7 | | 9.1 |
| | 4.6 | | 5.9 |
| | 3.5 | | 4.6 |
| | 5.5 | | 6.7 |
| | | | |

Linear Regression (m-samples)

2

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < m$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}^{(i)}}{m} \quad \eta \text{ is learning rate}$$

Generalized formula

1) Pick m samples (\mathbf{x}, \mathbf{y}) from training data

2) Compute output $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X}$$

3) Compute loss

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

4) Tính đạo hàm

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

\odot is element-wise
multiplication

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

5) Cập nhật tham số

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m} \quad \eta \text{ is learning rate}$$

More generalized formula

Linear Regression

1) Pick m samples (\mathbf{x}, \mathbf{y}) from training data

2) Compute output $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X}$$

3) Compute loss

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

4) Compute derivative

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

\odot is element-wise
multiplication

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m}$$

η is learning rate

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 data = np.genfromtxt('data.csv', delimiter=',')
5 areas = data[:,0]
6 prices = data[:,1:]
7 N = areas.size
8
9 # vector [x, b]^T
10 data = np.vstack([areas, np.ones((N,))])
11
12 # [w, b]
13 theta = np.array([[ -0.34],
14                  [ 0.04]])
15
16 # params
17 lr = 0.01
18 epoch_max = 1
19 m = 2
20
21 # logging
22 losses = []
```

1) Pick m samples (\mathbf{x}, \mathbf{y}) from training data

2) Compute output $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X}$$

3) Compute loss

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

4) Compute derivative

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

\odot is element-wise
multiplication

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m}$$

η is learning rate

More generalized formula

```

1 for epoch in range(epoch_max):
2     for i in range(0, N, m):
3         # get m samples
4         x = data[:, i:i+m]
5         y = prices[i:i+m, :]
6
7         # predict y_hat
8         y_hat = theta.T.dot(x)
9
10        # compute loss
11        loss = np.multiply((y_hat-y.T), (y_hat-y.T))
12        losses.append(np.mean(loss))
13
14        # compute gradient
15        k = 2*(y_hat-y.T)
16        gradients = np.multiply(np.vstack((k, k)), x)
17        gradients = gradients.dot(np.ones((m, 1))) / m
18
19        # update weights
20        theta = theta - lr*gradients

```

Linear Regression (m-samples)

Way 1.2

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix} \quad L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T) \\ = \begin{bmatrix} -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} k \\ k \end{bmatrix} = \begin{bmatrix} -22.676 & -14.848 \\ -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} k \\ k \end{bmatrix} \odot \mathbf{X} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix}$$

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} k \\ k \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix}$$

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

Take a look
and improve

Linear Regression (m-samples)

Way 1.2

1) Pick m samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} = (\mathbf{x}^{(i)})^T \boldsymbol{\theta} \quad \text{for } 0 \leq i < m$$

3) Compute loss

$$L = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})^2$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < m$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}}{m} \quad \eta \text{ is learning rate}$$

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix} \quad L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X} = [-2.238 \quad -1.524]$$

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T) \\ = \begin{bmatrix} -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} = \begin{bmatrix} -22.676 & -14.848 \\ -22.676 & -14.848 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix}$$

$$L'_{\boldsymbol{\theta}} = \left(\begin{bmatrix} \mathbf{k} \\ \mathbf{k} \end{bmatrix} \odot \mathbf{X} \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -151.92 & -68.301 \\ -22.676 & -14.848 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix}$$

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

$$\mathbf{X} = \begin{bmatrix} 6.7 & 4.6 \\ 1 & 1 \end{bmatrix}$$

$$L'_{\boldsymbol{\theta}} = \mathbf{X} \mathbf{k}^T$$

Linear Regression

Way 1.2

1) Pick m samples (\mathbf{x}, \mathbf{y}) from training data

2) Compute output $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X}$$

3) Compute loss

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

4) Compute derivative

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

$$L'_{\boldsymbol{\theta}} = \mathbf{X} \mathbf{k}^T$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m}$$

η is learning rate

More generalized formula

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 data = np.genfromtxt('data.csv', delimiter=',')
5 areas = data[:,0]
6 prices = data[:,1:]
7 N = areas.size
8
9 # vector [x, b]^T
10 data = np.vstack([areas, np.ones((N,))])
11
12 # [w, b]
13 theta = np.array([[ -0.34],
14                  [ 0.04]])
15
16 # params
17 lr = 0.01
18 epoch_max = 1
19 m = 2
20
21 # logging
22 losses = []
```


Way 1.2

1) Pick m samples (\mathbf{x}, \mathbf{y}) from training data

2) Compute output $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \boldsymbol{\theta}^T \mathbf{X}$$

3) Compute loss

$$L = \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y}^T)(\hat{\mathbf{y}} - \mathbf{y}^T)^T$$

4) Compute derivative

$$\mathbf{k} = 2(\hat{\mathbf{y}} - \mathbf{y}^T)$$

$$L'_{\boldsymbol{\theta}} = \mathbf{X}\mathbf{k}^T$$

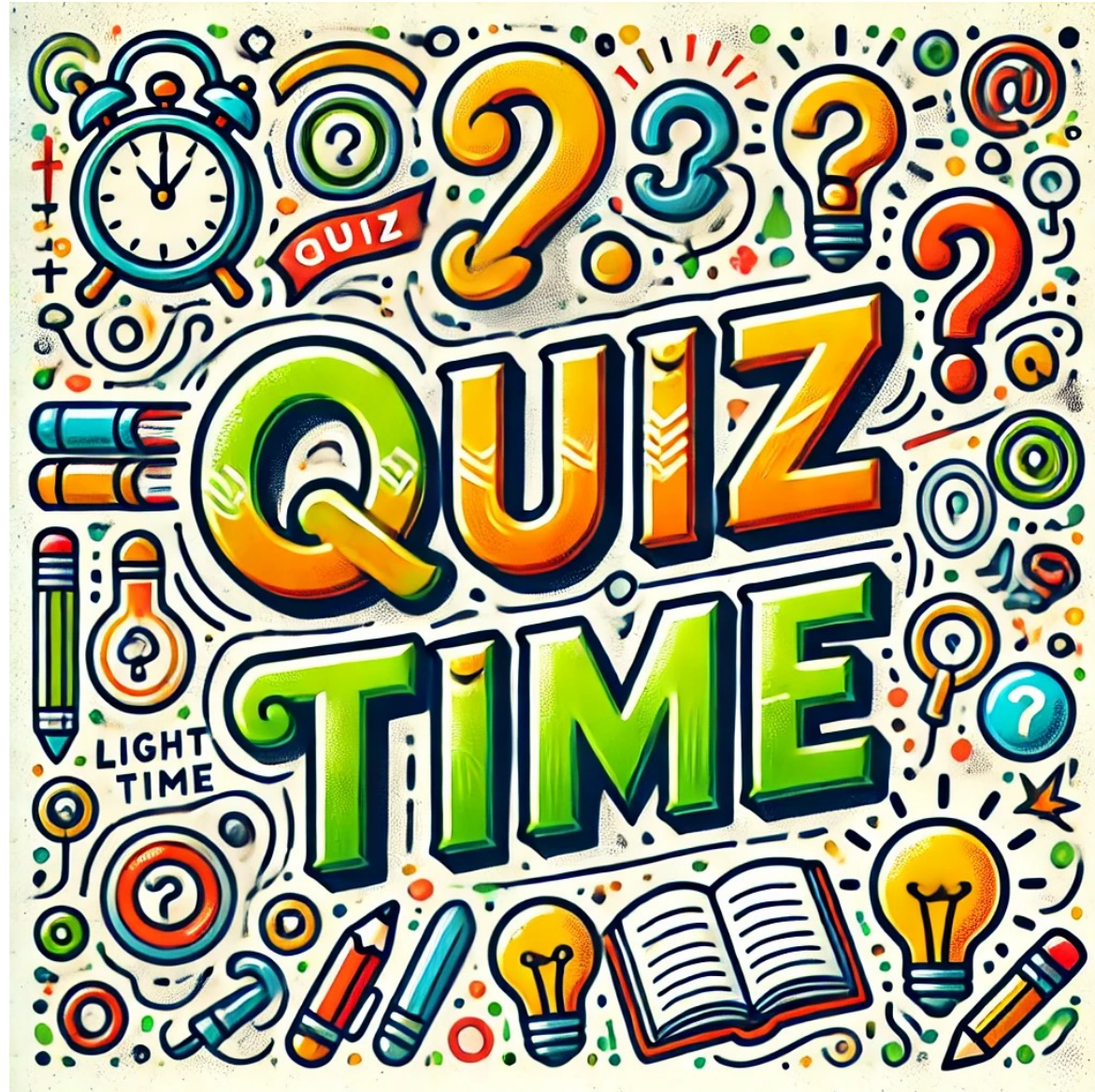
5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{L'_{\boldsymbol{\theta}}}{m}$$

η is learning rate

More generalized formula

```
1 for epoch in range(epoch_max):
2     for i in range(0, N, m):
3         # get m samples
4         x = data[:, i:i+m]
5         y = prices[i:i+m, :]
6
7         # predict y_hat
8         y_hat = theta.T.dot(x)
9
10        # compute loss
11        loss = np.multiply((y_hat-y.T), (y_hat-y.T))
12        losses.append(np.mean(loss))
13
14        # compute gradient
15        k = 2*(y_hat-y.T)
16        gradients = x.dot(k.T) / m
17
18        # update weights
19        theta = theta - lr*gradients
```



❖ What about this arrangement?

| Feature | | Label | |
|---------|------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix}$$

$$y = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$\hat{y} = X\theta = \begin{bmatrix} -2.238 \\ -1.524 \end{bmatrix}$$

$$L = \frac{1}{m}(\hat{y} - y)^T(\hat{y} - y) = \frac{1}{2}(128.5 + 55.11) = 91.8$$

$$k = 2(\hat{y} - y) = \begin{bmatrix} -22.676 \\ -14.848 \end{bmatrix}$$

$$L'_\theta = k^T X = \begin{bmatrix} -22.676 & -14.848 \end{bmatrix} \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix} = \begin{bmatrix} -220.23 & -37.524 \end{bmatrix}$$

$$\theta = \theta - \eta \left(\frac{L'_\theta}{m} \right)^T$$

$$= \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix} - 0.005 \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix} = \begin{bmatrix} 0.761 \\ 0.227 \end{bmatrix}$$

❖ What about this arrangement?

| Feature | | Label | |
|---------|------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix}$$

$$y = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

1) Pick m samples (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = X\theta$$

3) Compute loss

$$L = \frac{1}{m} (\hat{y} - y)^T (\hat{y} - y)$$

4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = k^T X$$

5) Update parameters

$$\theta = \theta - \eta \left(\frac{L'_{\theta}}{m} \right)^T \quad \eta \text{ is learning rate}$$

Linear Regression (m-samples)

34

❖ What about this arrangement?

| Feature | | Label | |
|---------|------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix}$$

$$y = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

$$\hat{y} = X\theta = \begin{bmatrix} -2.238 \\ -1.524 \end{bmatrix}$$

$$L = \frac{1}{m} (\hat{y} - y)^T (\hat{y} - y)$$

$$k = 2(\hat{y} - y) = \begin{bmatrix} -22.676 \\ -14.848 \end{bmatrix}$$

$$L'_\theta = X^T k = \begin{bmatrix} 6.7 \\ 1 \end{bmatrix} \begin{bmatrix} 4.6 \\ 1 \end{bmatrix} \begin{bmatrix} -22.676 \\ -14.848 \end{bmatrix} = \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix}$$

$$\theta = \theta - \eta \frac{L'_\theta}{m}$$

$$= \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix} - 0.005 \begin{bmatrix} -220.23 \\ -37.524 \end{bmatrix} = \begin{bmatrix} 0.761 \\ 0.227 \end{bmatrix}$$

❖ What about this arrangement?

| Feature | | Label | |
|---------|------|-------|--|
| | area | price | |
| | 6.7 | 9.1 | |
| | 4.6 | 5.9 | |
| | 3.5 | 4.6 | |
| | 5.5 | 6.7 | |
| | | | |

$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \end{bmatrix}$$

$$y = \begin{bmatrix} 9.1 \\ 5.9 \end{bmatrix}$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

1) Pick m samples (x, y) from training data

2) Compute output \hat{y}

$$\hat{y} = X\theta$$

3) Compute loss

$$L = \frac{1}{m}(\hat{y} - y)^T(\hat{y} - y)$$

4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = X^T k$$

5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{m}$$

η is learning rate

Outline

SECTION 1

1-sample Vectorization

Parameter Initialization

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

SECTION 2

m-sample Vectorization

Way 1 for constructing matrix \mathbf{x}

$$\mathbf{x} = \begin{bmatrix} 6.7 & 4.6 & 3.5 & 5.5 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

SECTION 3

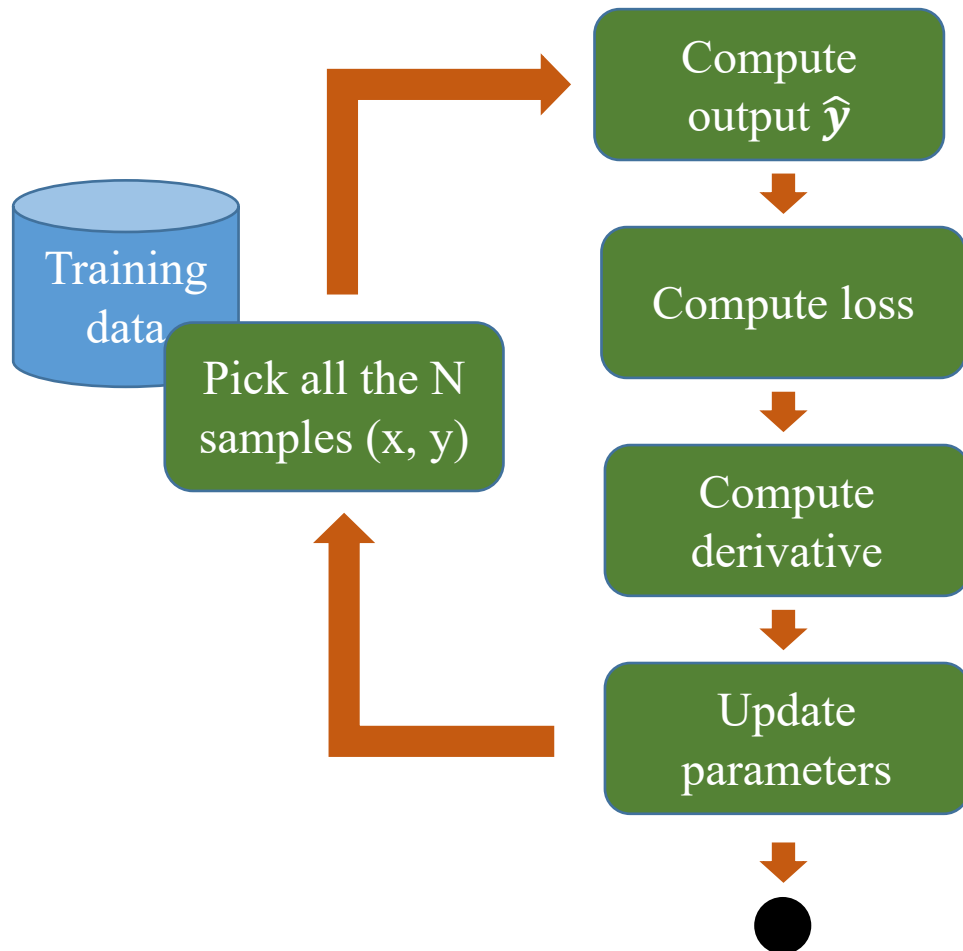
N-sample Vectorization

Way 2 for constructing matrix \mathbf{x}

$$\mathbf{x} = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \\ 3.5 & 1 \\ 5.5 & 1 \end{bmatrix}$$

❖ House price prediction

❖ N-sample training



1) Pick all the N samples $(x^{(i)}, y^{(i)})$ from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = wx^{(i)} + b \quad \text{for } 0 \leq i < N$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < N$$

4) Compute derivative

$$L'_w{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)})$$

$$L'_b{}^{(i)} = 2(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < N$$

5) Update parameters

$$w = w - \eta \frac{\sum_i L'_w{}^{(i)}}{N}$$

$$b = b - \eta \frac{\sum_i L'_b{}^{(i)}}{N}$$

Learning rate η

1) Pick all the N samples from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = wx^{(i)} + b \quad \text{for } 0 \leq i < N$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < N$$

4) Compute derivative

$$\begin{aligned} L'_w{}^{(i)} &= 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \\ L'_b{}^{(i)} &= 2(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < N \end{aligned}$$

5) Update parameters

$$\begin{aligned} w &= w - \eta \frac{\sum_i L'_w{}^{(i)}}{N} \\ b &= b - \eta \frac{\sum_i L'_b{}^{(i)}}{N} \end{aligned} \quad \eta \text{ is learning rate}$$

Friendly version

1) Pick all the N samples from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} \quad \text{for } 0 \leq i < N$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < N$$

4) Compute derivative

$$L'_\theta{}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < N$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_\theta{}^{(i)}}{N} \quad \eta \text{ is learning rate}$$

Generalized formula

1) Pick all the N samples from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} \quad \text{for } 0 \leq i < N$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < N$$

4) Compute derivative

$$L'_{\boldsymbol{\theta}}^{(i)} = 2\mathbf{x}^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < N$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \frac{\sum_i L'_{\boldsymbol{\theta}}^{(i)}}{N} \quad \eta \text{ is learning rate}$$

Generalized formula

```
1 import numpy as np
2 from numpy import genfromtxt
3
4 data = genfromtxt('data.csv', delimiter=',')
5 areas = data[:,0]
6 prices = data[:,1]
7 data_size = areas.size
8
9 # vector [x, b]
10 data = np.c_[areas, np.ones((data_size, 1))]
11
12 n_epochs = 10
13 lr = 0.01
14
15 theta = np.array([[-0.34],[0.04]])
```

Linear Regression

1) Pick all the N samples from training data

2) Compute output $\hat{y}^{(i)}$

$$\hat{y}^{(i)} = \theta^T x^{(i)} \quad \text{for } 0 \leq i < N$$

3) Compute loss

$$L^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2 \quad \text{for } 0 \leq i < N$$

4) Compute derivative

$$L'_{\theta}{}^{(i)} = 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \quad \text{for } 0 \leq i < N$$

5) Update parameters

$$\theta = \theta - \eta \frac{\sum_i L'_{\theta}{}^{(i)}}{N} \quad \eta \text{ is learning rate}$$

Generalized formula

```
1 losses = [] # for debug
2 for epoch in range(n_epochs):
3     sum_of_losses = 0
4     gradients = np.zeros((2,1))
5
6     for index in range(data_size):
7         # get data
8         x_i = data[index:index+1]
9         y_i = prices[index:index+1]
10
11         # compute output y_hat_i
12         y_hat_i = x_i.dot(theta)
13
14         # compute loss
15         l_i = (y_hat_i - y_i)*(y_hat_i - y_i)
16
17         # compute gradient
18         g_l_i = 2*(y_hat_i - y_i)
19         gradient = x_i.T.dot(g_l_i)
20
21         # accumulate gradient
22         gradients = gradients + gradient
23         sum_of_losses = sum_of_losses + l_i
24
25     # normalize
26     sum_of_losses = sum_of_losses/data_size
27     gradients = gradients/data_size
28
29     # for debug
30     losses.append(sum_of_losses[0][0])
31
32     # update
33     theta = theta - lr*gradients
```

❖ Vectorization

House price data

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

Model

$$\text{price} = w * \text{area} + b$$

$$\hat{y} = wx + b$$

$$y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

Parameter Initialization

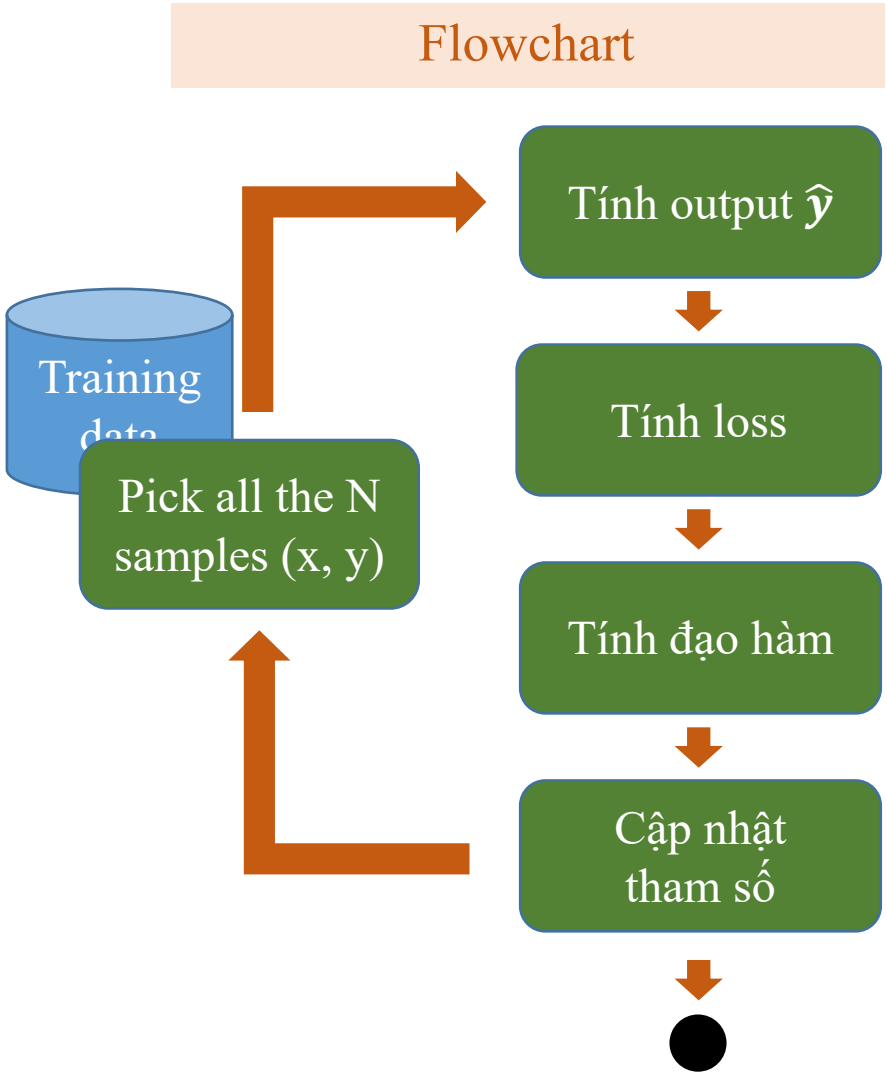
$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$

Way 1 for constructing matrix x

$$X = \begin{bmatrix} 6.7 & 4.6 & 3.5 & 5.5 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

Way 2 for constructing matrix x

$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \\ 3.5 & 1 \\ 5.5 & 1 \end{bmatrix}$$



Linear Regression [1]

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

House Price Data

Way 1

$$X = \begin{bmatrix} 6.7 & 4.6 & 3.5 & 5.5 \\ 1 & 1 & 1 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

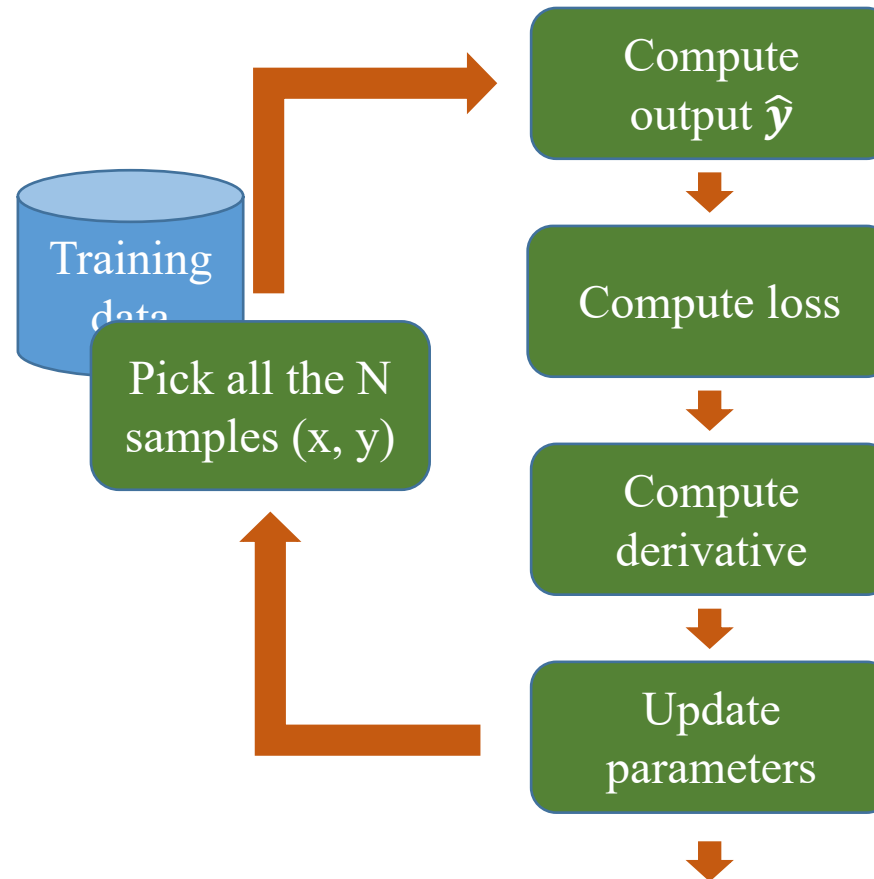
Model

$$\text{price} = w * \text{area} + b$$

$$\hat{y} = wx + b$$

Parameter Initialization

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$



1) Pick all the N samples from training data

2) Compute output \hat{y}

$$\hat{y} = \theta^T X$$

3) Compute loss

$$L = \frac{1}{N} (\hat{y} - y)(\hat{y} - y)^T$$

4) Compute derivative

$$k = 2(\hat{y} - y^T)$$

$$L'_{\theta} = Xk^T$$

5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{N}$$

η is learning rate

Vectorization Approach

Linear Regression [2]

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

House Price Data

Way 2

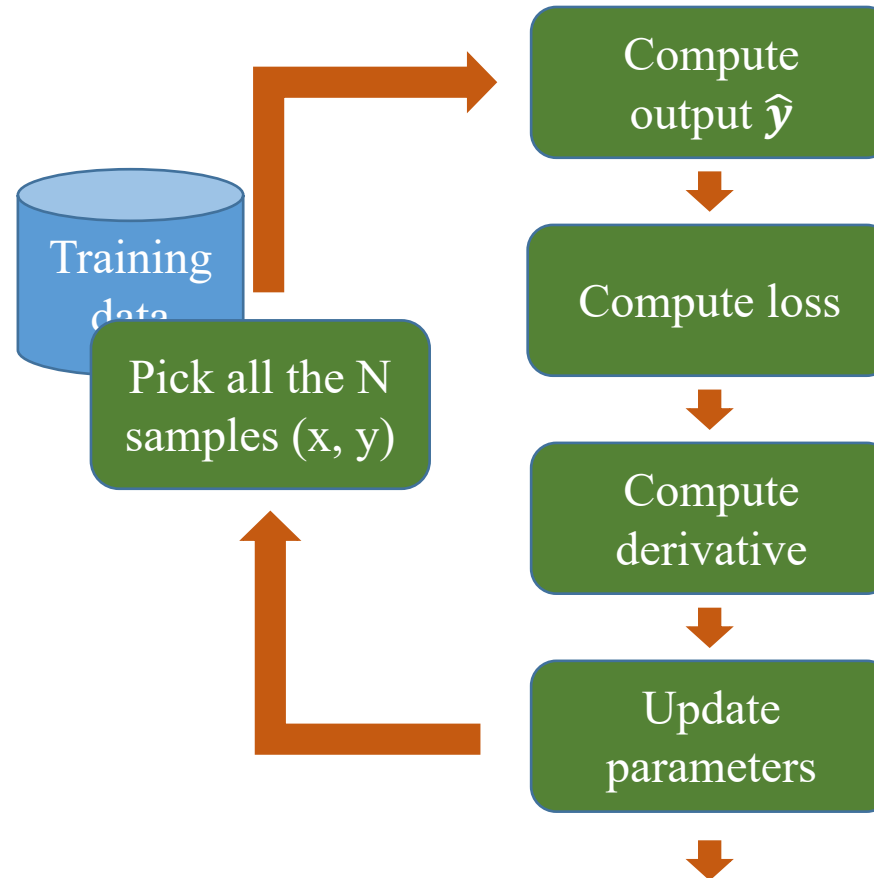
$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \\ 3.5 & 1 \\ 5.5 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

Model

$$\text{price} = w * \text{area} + b$$
$$\hat{y} = wx + b$$

Parameter Initialization

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$



1) Pick all the N samples from training data

2) Compute output \hat{y}

$$\hat{y} = X\theta$$

3) Compute loss

$$L = \frac{1}{N} (\hat{y} - y)^T (\hat{y} - y)$$

4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = k^T X$$

5) Update parameters

$$\theta = \theta - \eta \left(\frac{L'_{\theta}}{N} \right)^T$$

η is learning rate

Vectorization Approach

Linear Regression [3]

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

House Price Data

Way 2

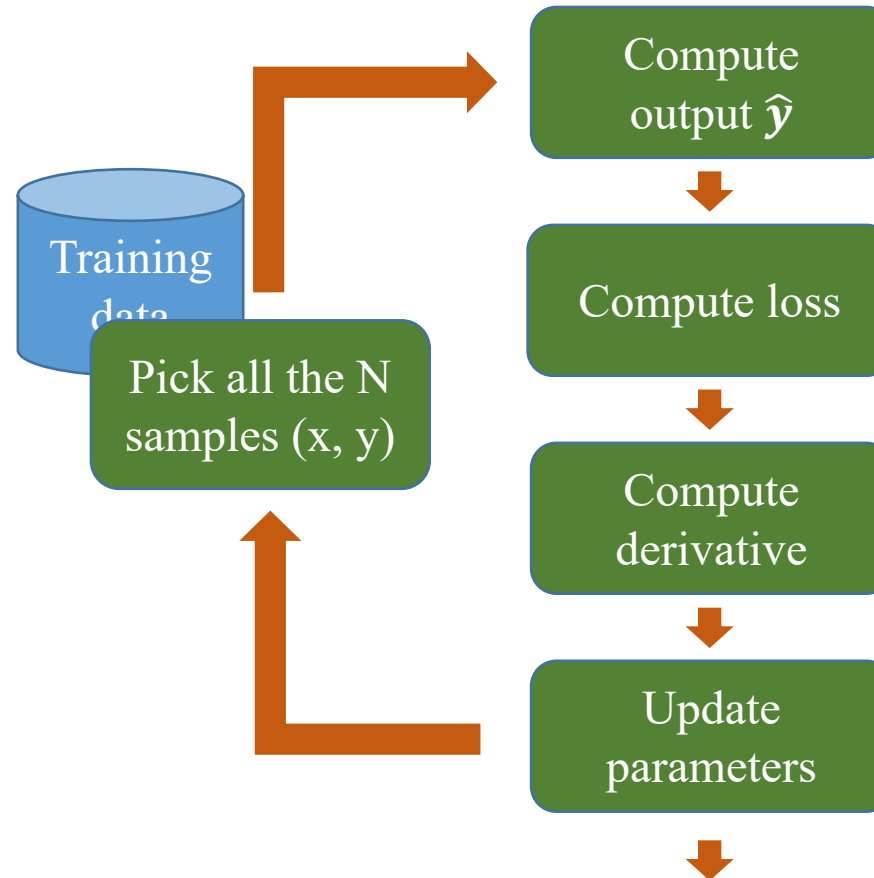
$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \\ 3.5 & 1 \\ 5.5 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

Model

$$\text{price} = w * \text{area} + b$$
$$\hat{y} = wx + b$$

Parameter Initialization

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$



1) Pick all the N samples from training data

2) Compute output \hat{y}

$$\hat{y} = X\theta$$

3) Compute loss

$$L = \frac{1}{N} (\hat{y} - y)^T (\hat{y} - y)$$

4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = X^T k$$

5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{N}$$

η is learning rate

Vectorization Approach

Summary

| Feature | Label |
|---------|-------|
| area | price |
| 6.7 | 9.1 |
| 4.6 | 5.9 |
| 3.5 | 4.6 |
| 5.5 | 6.7 |

House Price Data

Way 3

$$X = \begin{bmatrix} 6.7 & 1 \\ 4.6 & 1 \\ 3.5 & 1 \\ 5.5 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 9.1 \\ 5.9 \\ 4.6 \\ 6.7 \end{bmatrix}$$

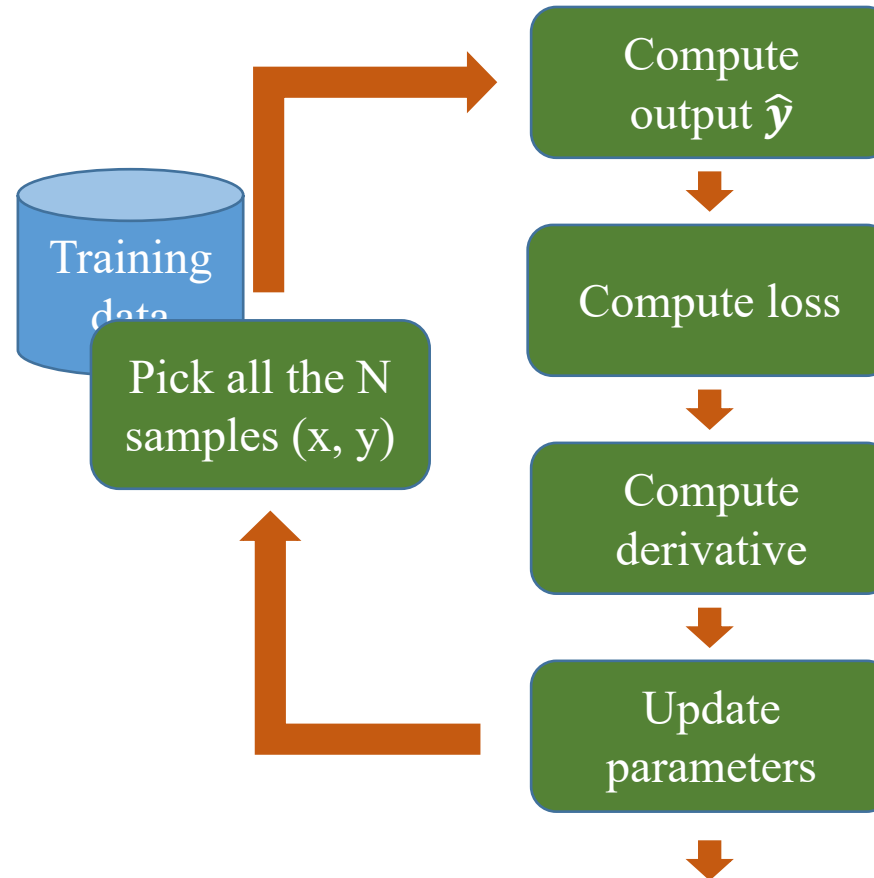
Model

$$\text{price} = w * \text{area} + b$$

$$\hat{y} = wx + b$$

Parameter Initialization

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} -0.34 \\ 0.049 \end{bmatrix}$$



- 1) Pick all the N samples from training data
- 2) Compute output \hat{y}

$$\hat{y} = X\theta$$

- 3) Compute loss

$$L = \frac{1}{N} (\hat{y} - y)^T (\hat{y} - y)$$

- 4) Compute derivative

$$k = 2(\hat{y} - y)$$

$$L'_{\theta} = X^T k$$

- 5) Update parameters

$$\theta = \theta - \eta \frac{L'_{\theta}}{N} \quad \eta \text{ is learning rate}$$

```
for epoch in range(epoch_max):
    y_hat = x.dot(theta)

    loss = np.multiply((y_hat-y), (y_hat-y))
    losses.append(np.mean(loss))

    k = 2*(y_hat-y)
    gradients = x.T.dot(k) / N

    theta = theta - lr*gradients
```