

Linear Regression

A Simple Approach

Quang-Vinh Dinh
PhD in Computer Science

Outline

- Simple version of Linear Regression
- Computational Graph
- Mini-batch Training
- Batch Training
- Generalization of Linear Regression
- Loss Functions (extension)

Linear Regression

Introduction

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=?

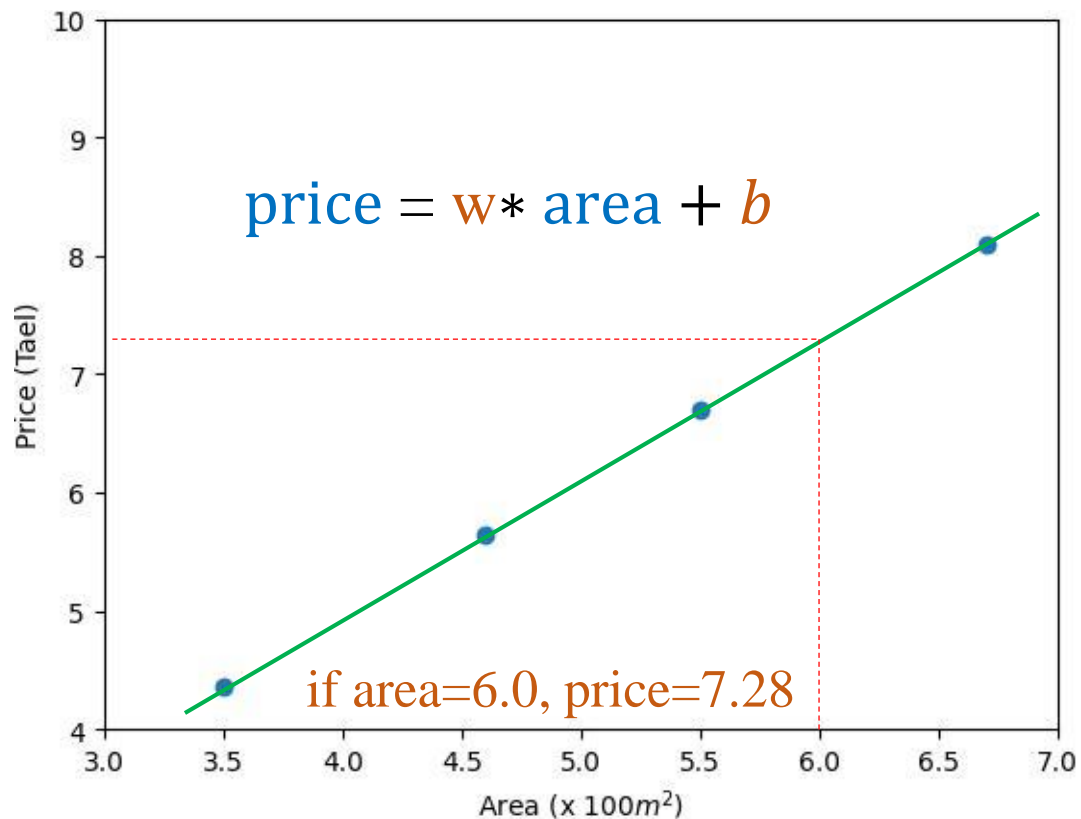
Features													Label
crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	black	lstat	medv
0.00632	18	2.31	0	0.538	6.575	65.2	4.09	1	296	15.3	396.9	4.98	24
0.02731	0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.9	9.14	21.6
0.03237	0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
0.06905	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

House Price Prediction

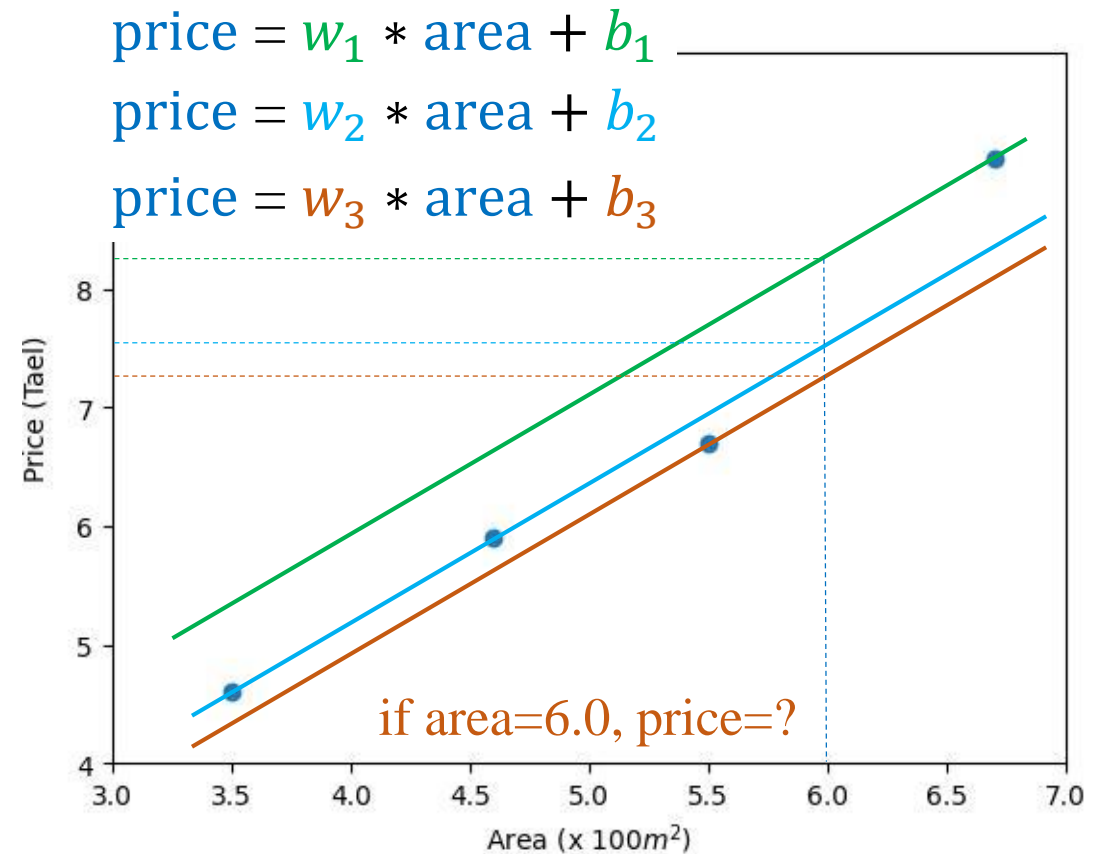
House price data

Feature	Label
area	price
6.7	8.1
4.6	5.6
3.5	4.3
5.5	6.7



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

House price data

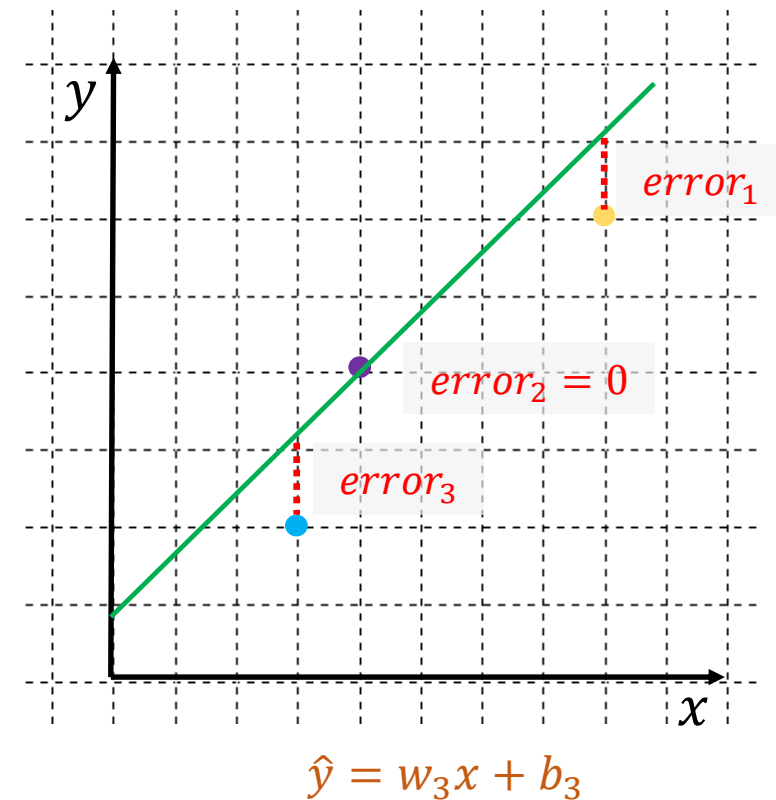
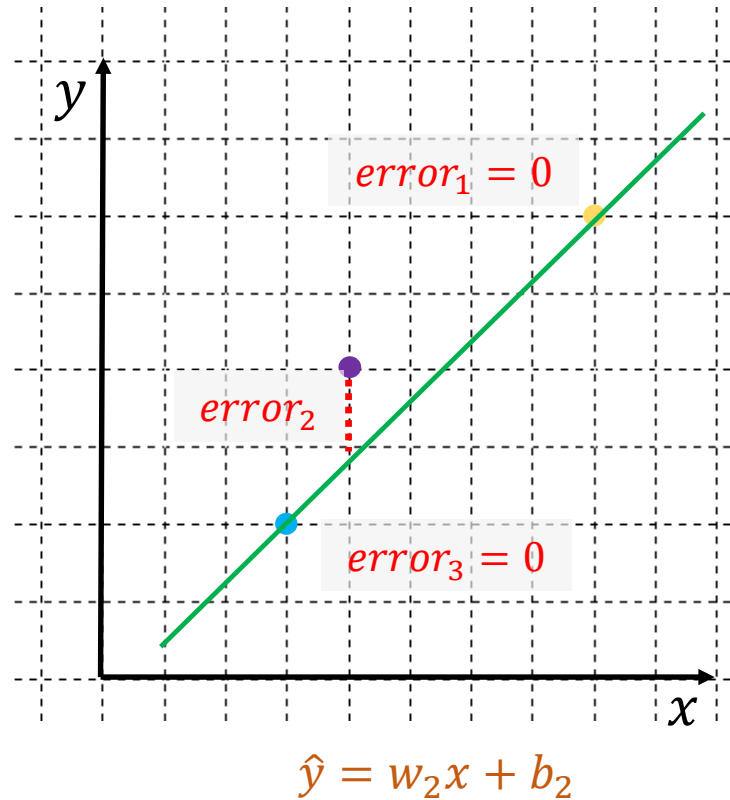
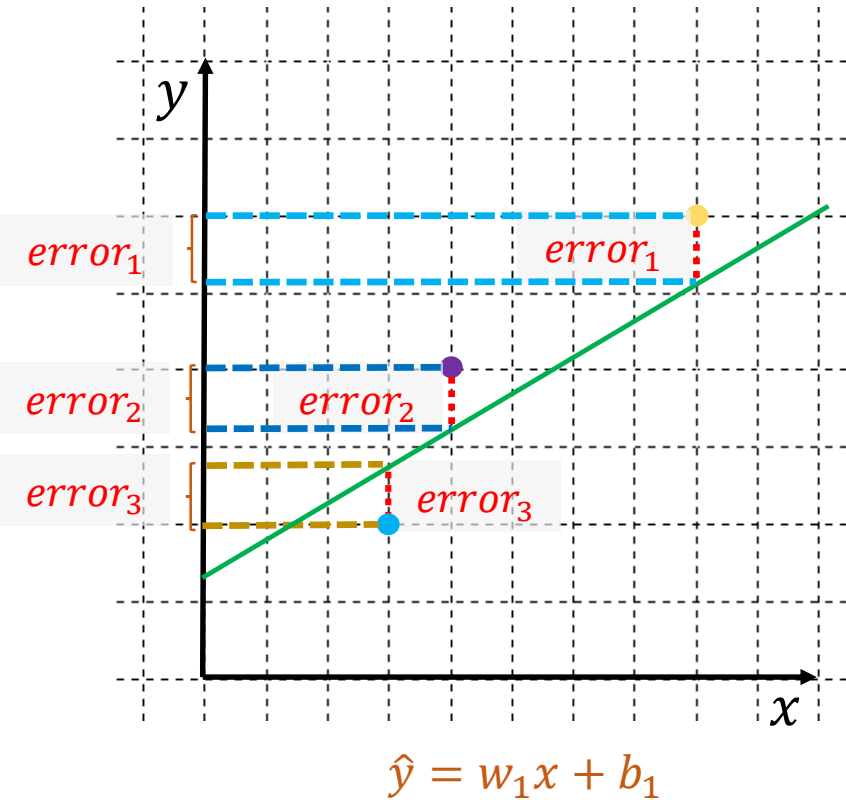


Linear Regression

❖ Area-based house price prediction

● Training data

$$error_i = distance(\hat{y}_i, y_i)$$



Find a and b whose model has the smallest error, where $error = \sum_i error_i$ **How?**

Linear Regression

❖ Area-based house price prediction

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

$$\text{error} = (\text{predicted_price} - \text{real_price})^2$$

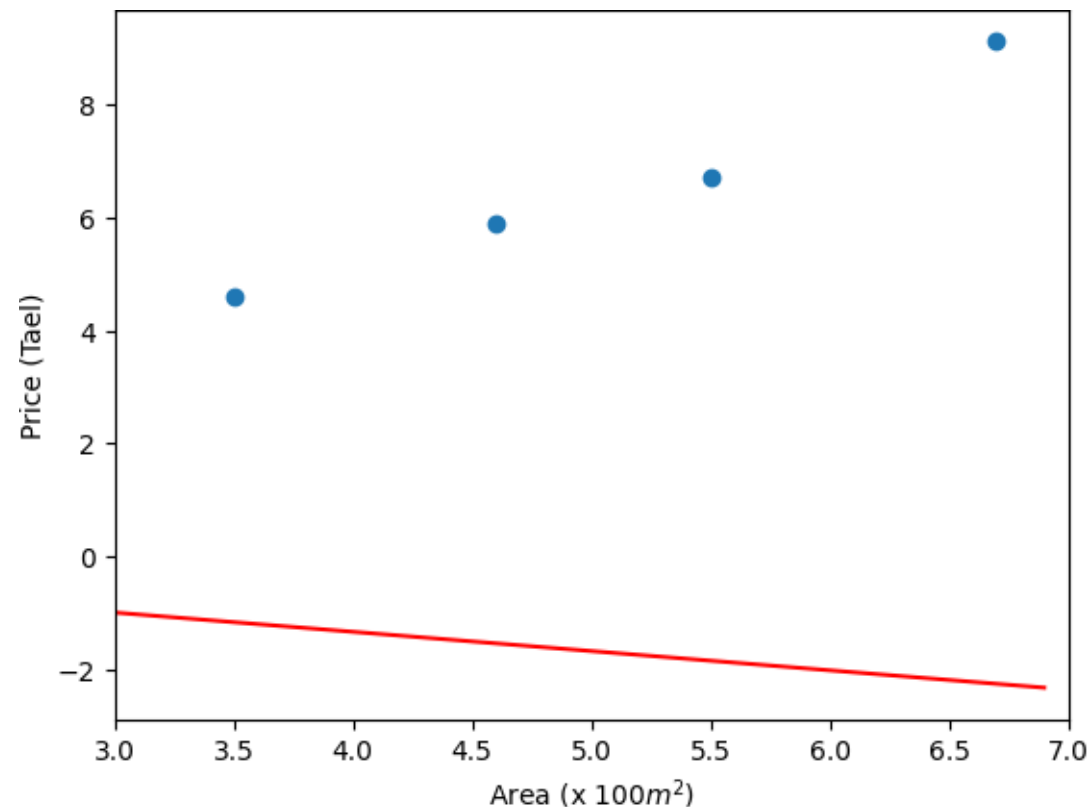
$$\hat{y}_i = wx_i + b$$

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

area	price	predicted	error
6.7	9.1	-2.238	128.55
4.6	5.9	-1.524	55.11
3.5	4.6	-1.15	33.06
5.5	6.7	-1.83	72.76

$$w = -0.34$$

$$b = 0.04$$



Linear Regression

❖ Area-based house price prediction

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

$$\text{error} = (\text{predicted_price} - \text{real_price})^2$$

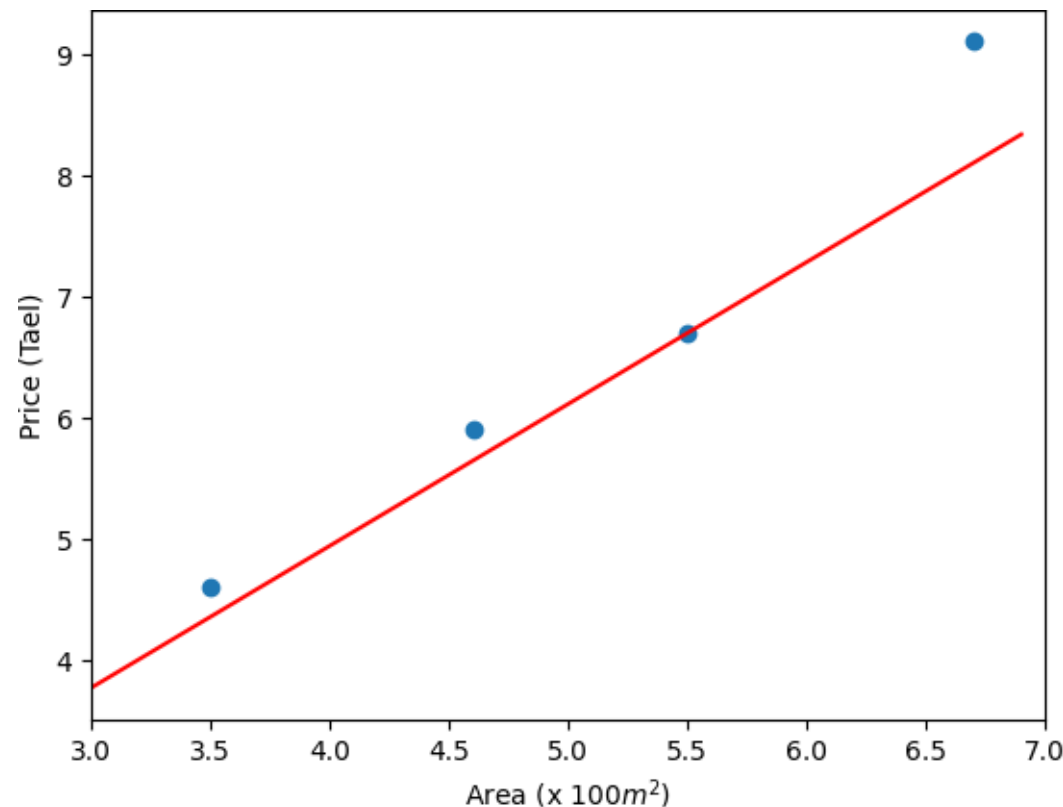
$$\hat{y}_i = wx_i + b$$

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

area	price	predicted	error
6.7	9.1	8.099	1.002
4.6	5.9	5.642	0.066
3.5	4.6	4.355	0.06
5.5	6.7	6.695	0.00002

$$w = 1.17$$

$$b = 0.26$$

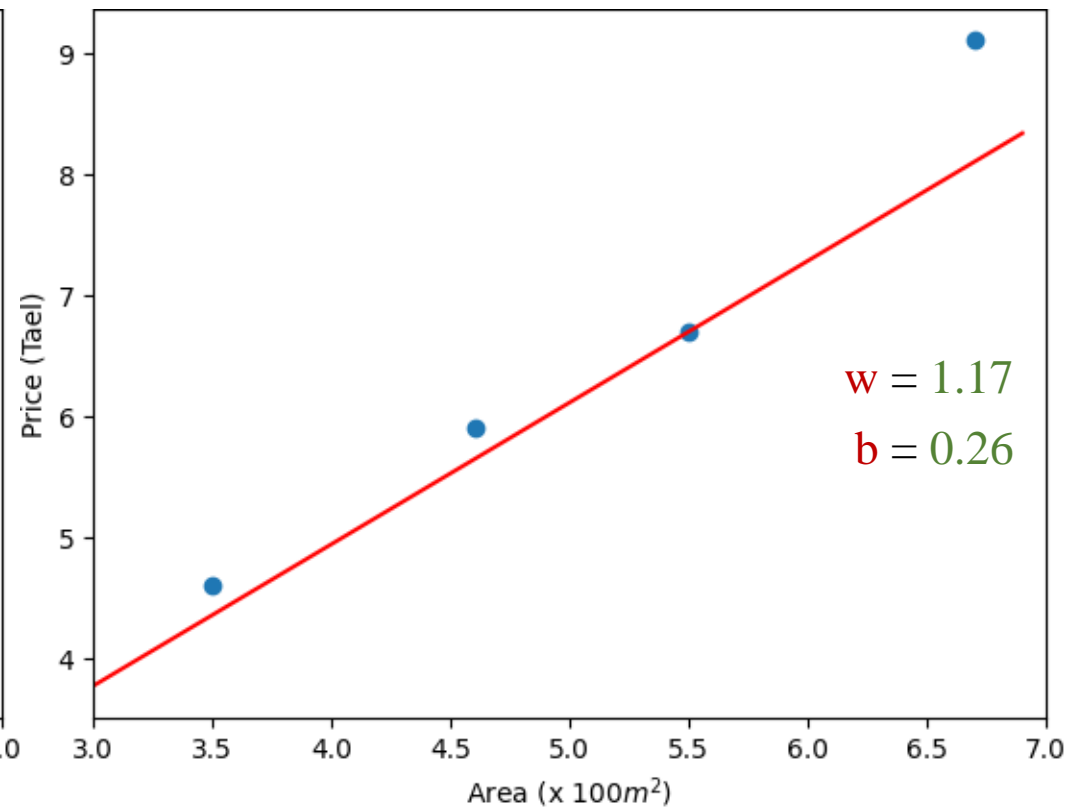
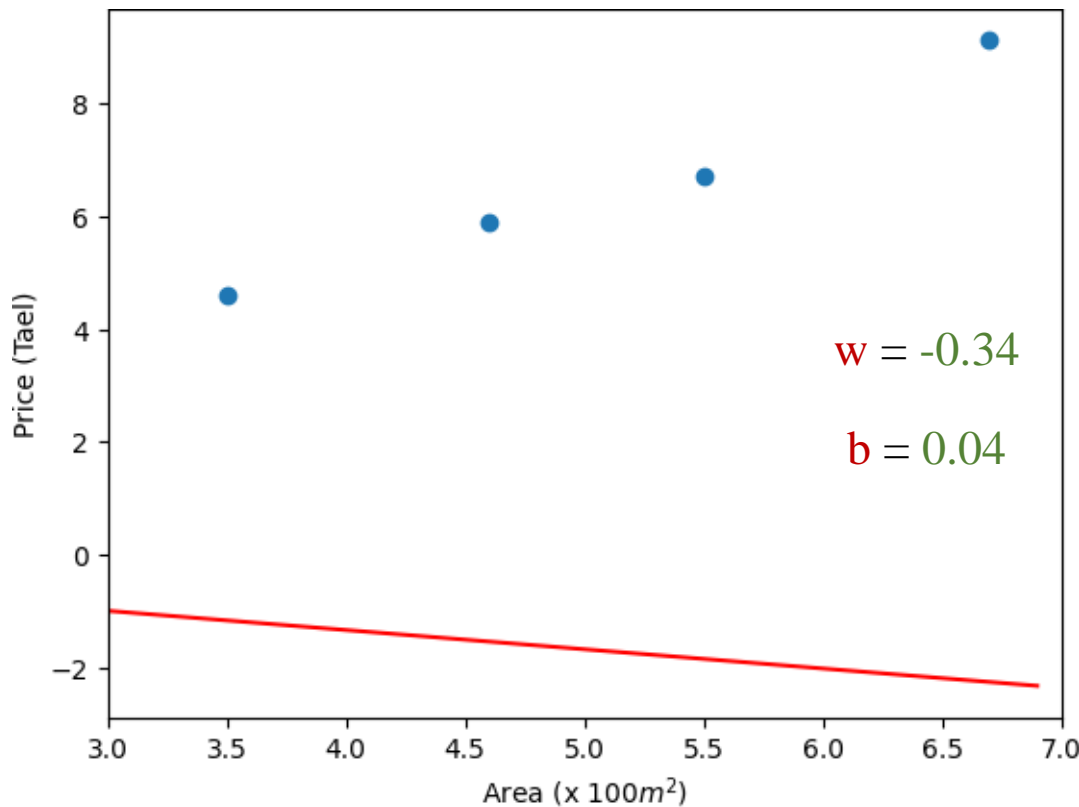


Linear Regression

❖ Area-based house price prediction

$$\hat{y}_i = wx_i + b$$
$$L(\hat{y}_i, y_i) = (\hat{y}_i - y_i)^2$$

How to change w and b so that $L(\hat{y}_i, y_i)$ reduces



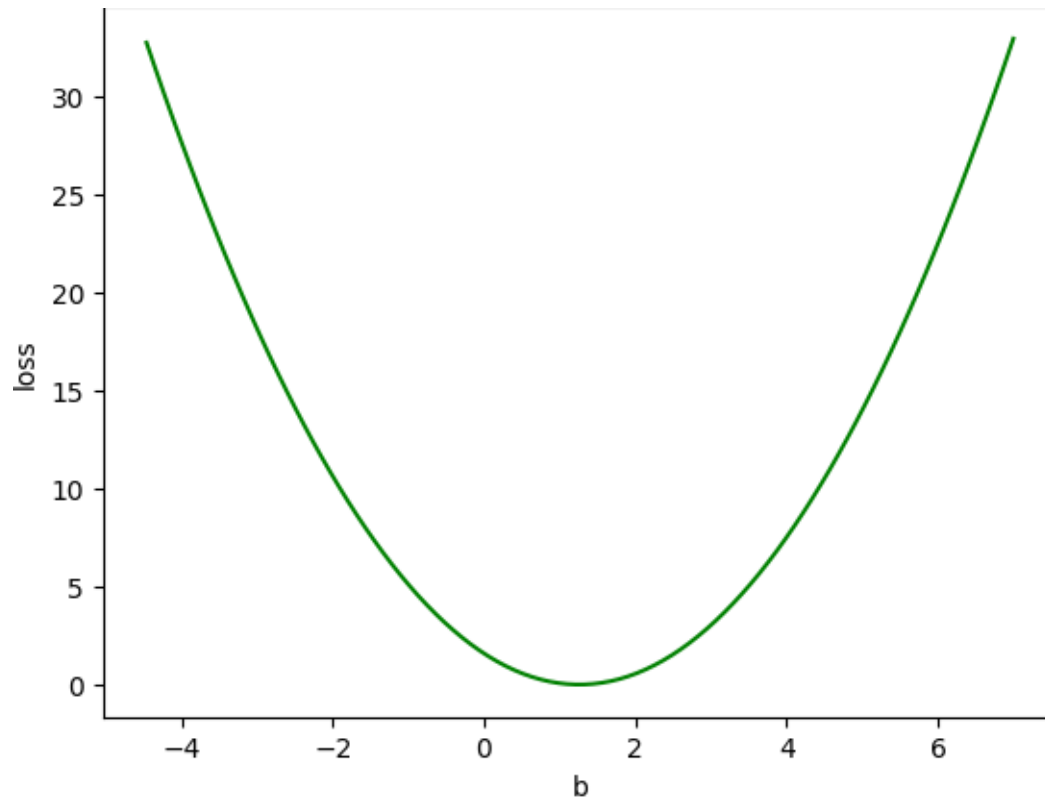
Linear Regression

$$\hat{y}_i = wx_i + b$$

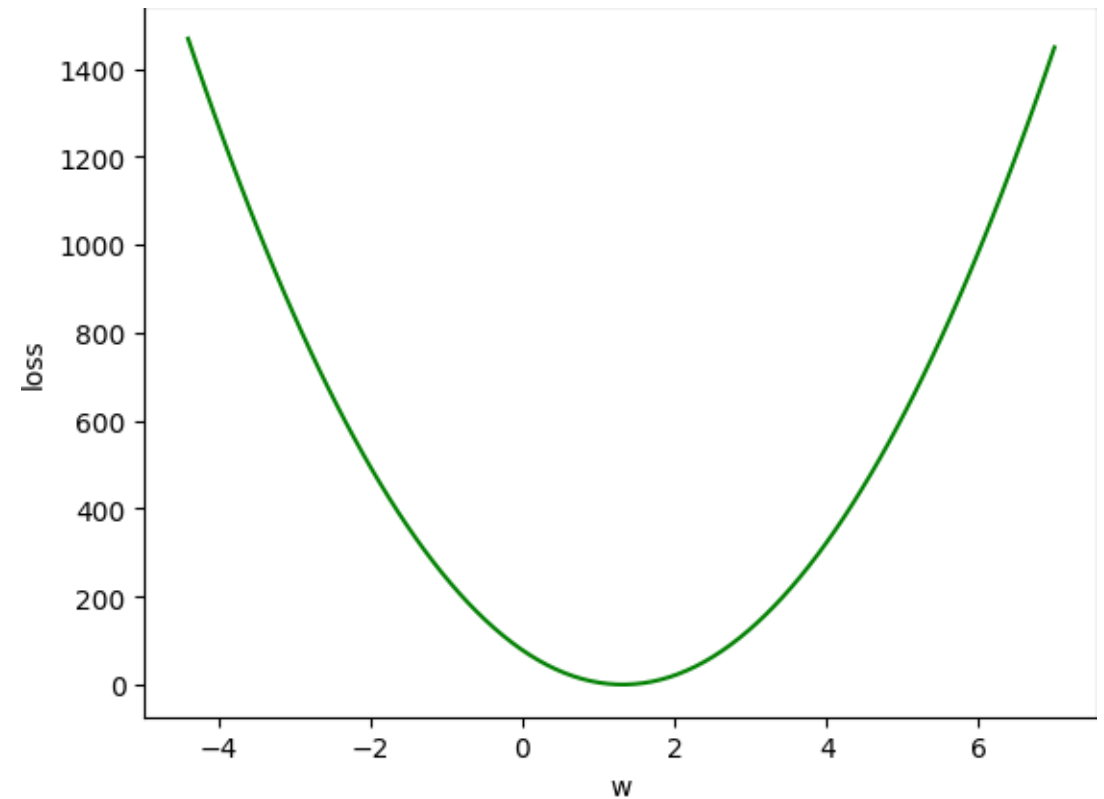
$$L(\hat{y}_i, y_i) = (\hat{y}_i - y_i)^2$$

❖ Understanding the loss function

How to change w and b so that $L(\hat{y}_i, y_i)$ reduces



Different b values with a fixed w value



Different w values with a fixed b value

Linear Regression

Linear equation

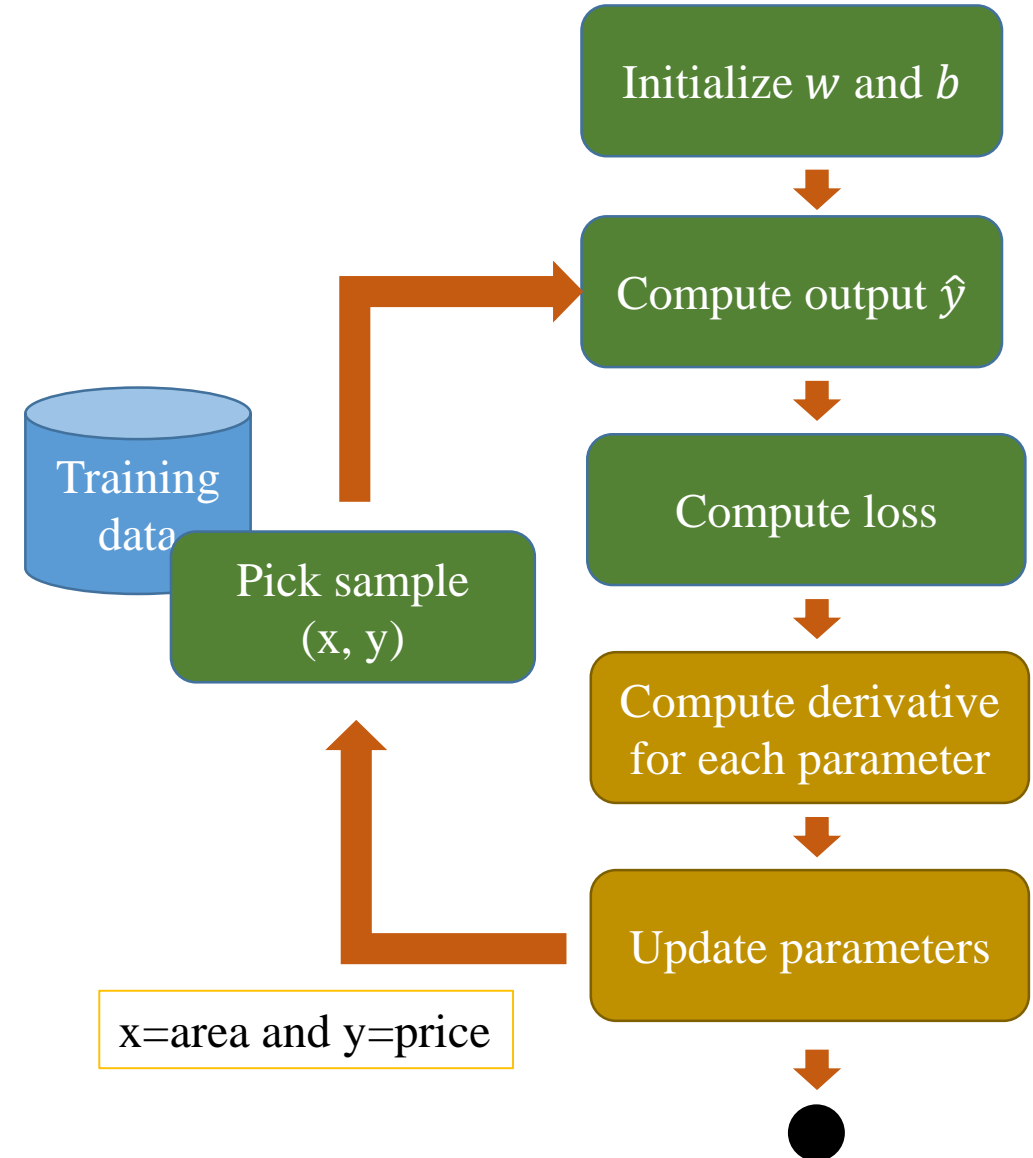
$$\hat{y} = wx + b$$

where \hat{y} is a predicted value,
 w and b are parameters
and x is input feature

Error (loss) computation

Idea: compare predicted values \hat{y} and label values y
Squared loss

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



Linear Regression

Linear equation

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

where \hat{y} is a predicted value,

w and b are parameters

and x is input feature

Error (loss) computation

Idea: compare predicted values \hat{y} and label values y

Squared loss

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

Find better w and b

Use gradient descent to minimize the loss function

Compute derivate for each parameter

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

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

Update parameters

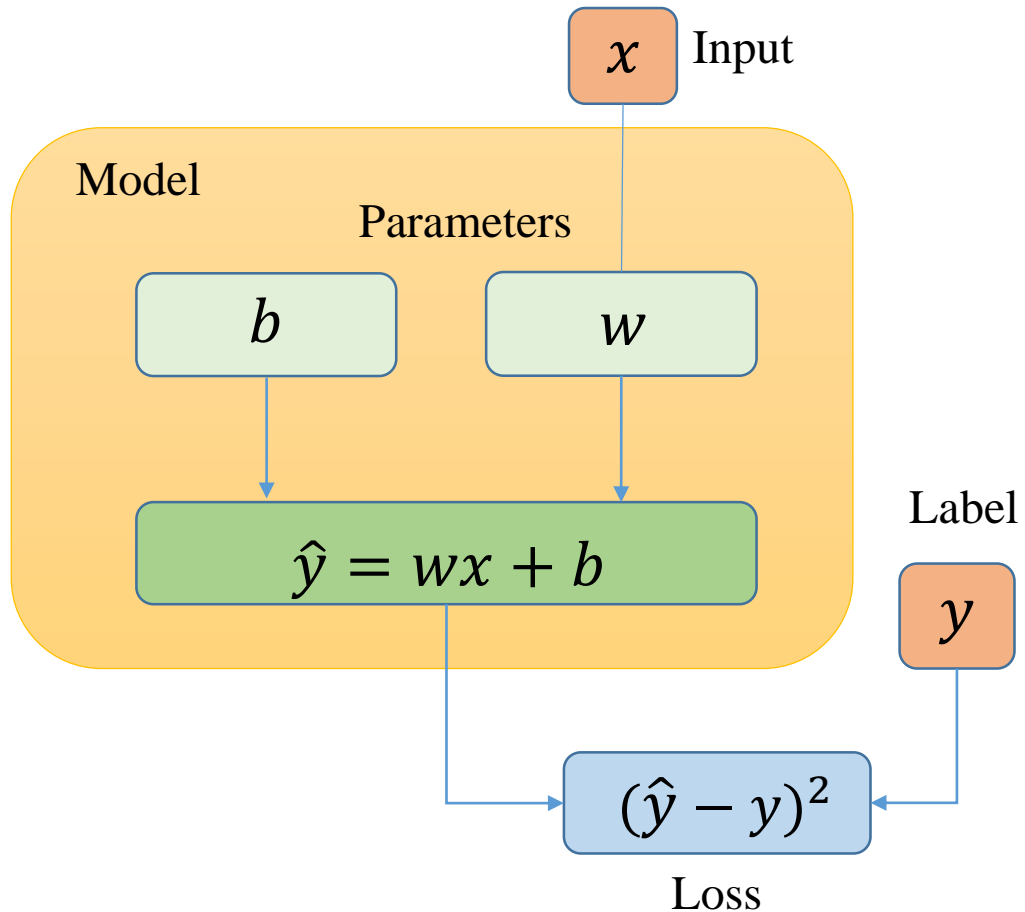
$$w = w - \eta \frac{\partial L}{\partial w} \quad b = b - \eta \frac{\partial L}{\partial b}$$

η is learning rate

Linear Regression

❖ Toy example

Diagram



Cheat sheet

Compute the output \hat{y}

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

Compute the loss

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

Compute derivative

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

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

Update parameters

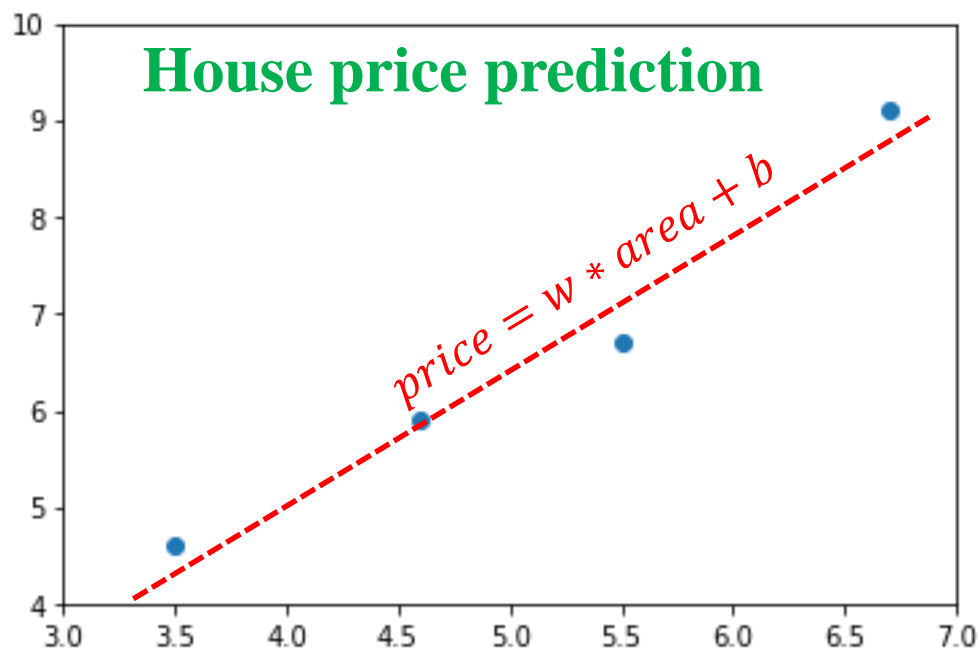
$$w = w - \eta \frac{\partial L}{\partial w}$$

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

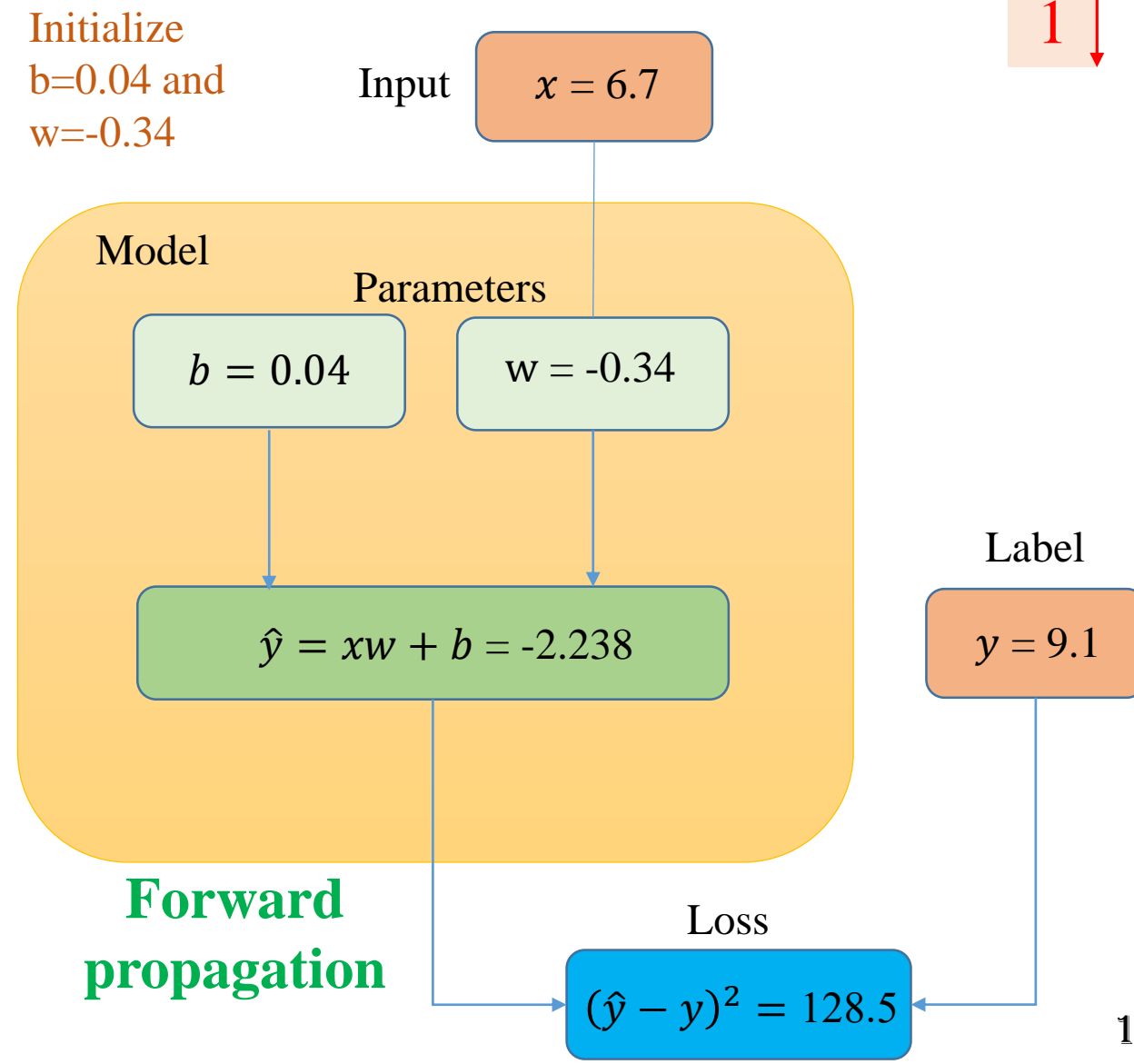
Linear Regression

Given
sample
data

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



Initialize
 $b=0.04$ and
 $w=-0.34$



Linear Regression

2

Input

$x = 6.7$

Backpropagation

Model

Parameters

$b = 0.26676$

$w = 1.17929$

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

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

$$\hat{y} = xw + b = -2.238$$

$$\begin{aligned} \frac{\partial L}{\partial w} &= 2x(\hat{y} - y) \\ &= -151.9292 \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial b} &= 2(\hat{y} - y) \\ &= -22.676 \end{aligned}$$

Label

$y = 9.1$

Loss

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

$\eta = 0.01$

3

Input

$x = 6.7$

Forward propagation

Model

Parameters

$b = 0.26676$

$w = 1.17929$

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

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

$$\hat{y} = xw + b = -2.238$$

Label

$y = 9.1$

Loss

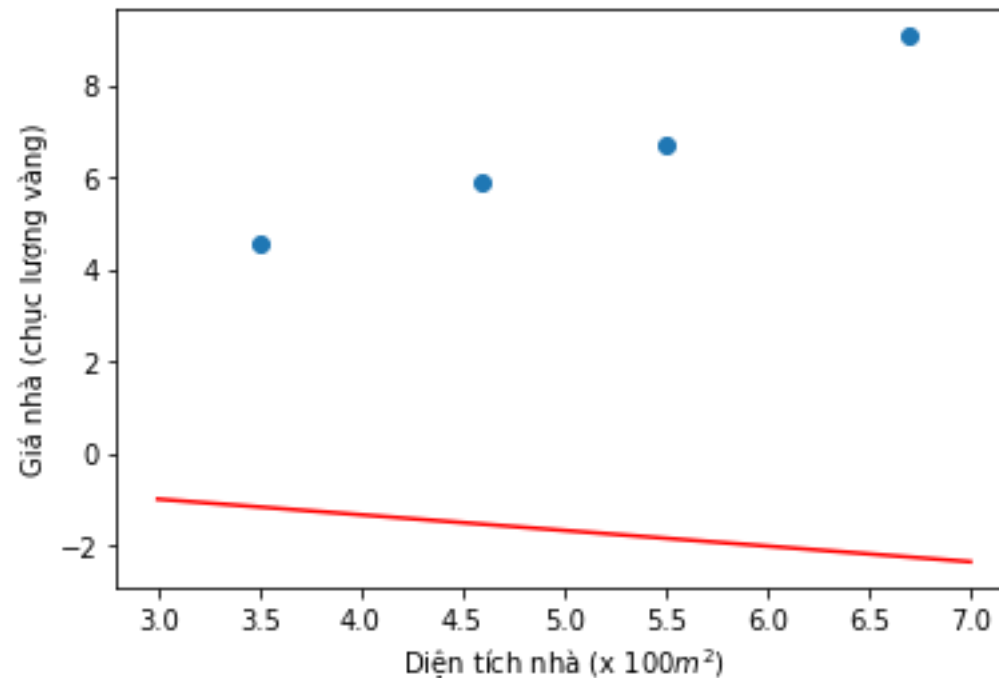
$$(\hat{y} - y)^2 = 0.868$$

New w and b help
the loss reduce

Linear Regression

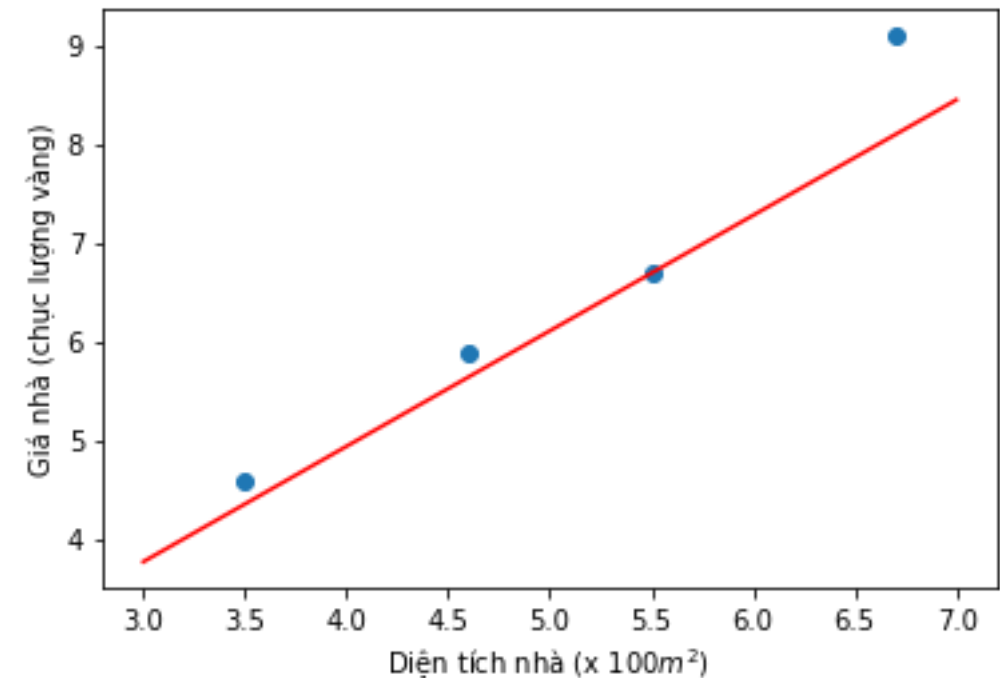
❖ Toy example

Model prediction before and after the first update



$$w = -0.34 \quad b = 0.04 \quad L = 128.55$$

Before updating

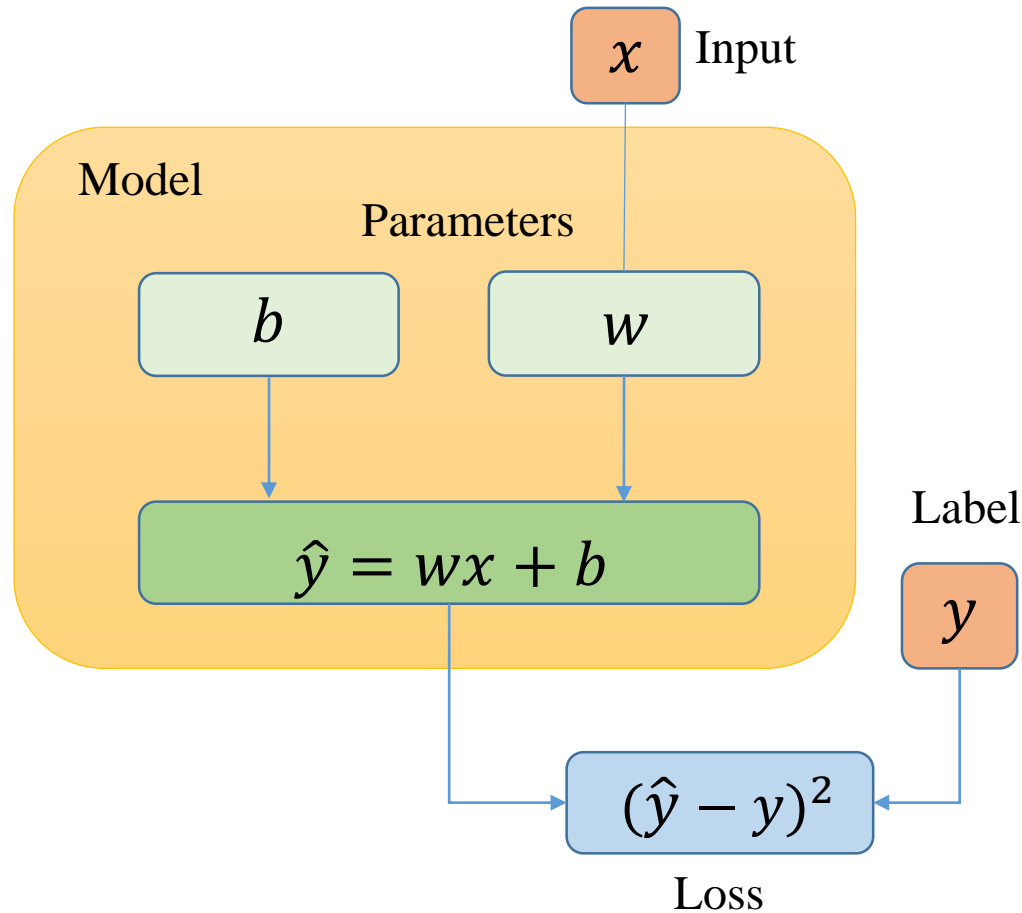


$$w = 1.179292 \quad b = 0.26676 \quad L = 0.868$$

After updating

Linear Regression

❖ Summary (simple version)



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

Linear Regression

❖ Implementation

Cheat sheet

Compute the output \hat{y}

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

Compute the loss

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

Compute derivative

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

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

Update parameters

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

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

```
1 # forward
2 def predict(x, w, b):
3     return x*w + b
4
5 # compute gradient
6 def gradient(y_hat, y, x):
7     dw = 2*x*(y_hat-y)
8     db = 2*(y_hat-y)
9
10    return (dw, db)
11
12 # update weights
13 def update_weight(w, b, lr, dw, db):
14     w_new = w - lr*dw
15     b_new = b - lr*db
16
17    return (w_new, b_new)
```

Linear Regression

❖ Code for one update

```
1 # forward
2 def predict(x, w, b):
3     return x*w + b
4
5 # compute gradient
6 def gradient(y_hat, y, x):
7     dw = 2*x*(y_hat-y)
8     db = 2*(y_hat-y)
9
10    return (dw, db)
11
12 # update weights
13 def update_weight(w, b, lr, dw, db):
14     w_new = w - lr*dw
15     b_new = b - lr*db
16
17    return (w_new, b_new)
```

```
1 # test sample
2 x = 6.7
3 y = 9.1
4
5 # init weights
6 b = 0.04
7 w = -0.34
8 lr = 0.01
9
10 # predict y_hat
11 y_hat = predict(x, w, b)
12 print('y_hat: ', y_hat)
13
14 # compute loss
15 loss = (y_hat-y)*(y_hat-y)
16 print('Loss: ', loss)
17
18 # compute gradient
19 (dw, db) = gradient(y_hat, y, x)
20 print('dw: ', dw)
21 print('db: ', db)
22
23 # update weights
24 (w, b) = update_weight(w, b, lr, dw, db)
25 print('w_new: ', w)
26 print('b_new: ', b)
```

Outline

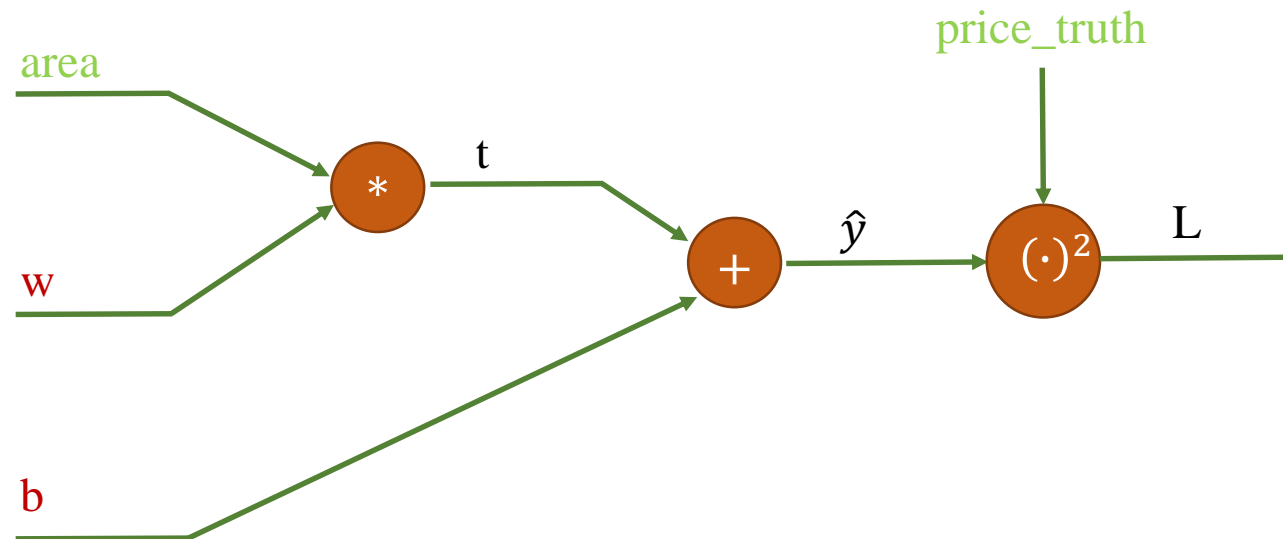
- **Simple version of Linear Regression**
- **Computational Graph**
- **Mini-batch Training**
- **Batch Training**
- **Generalization of Linear Regression**
- **Loss Functions**

Computational graph

- ❖ House price prediction
 - ❖ One-sample training

$$price = w * area + b$$

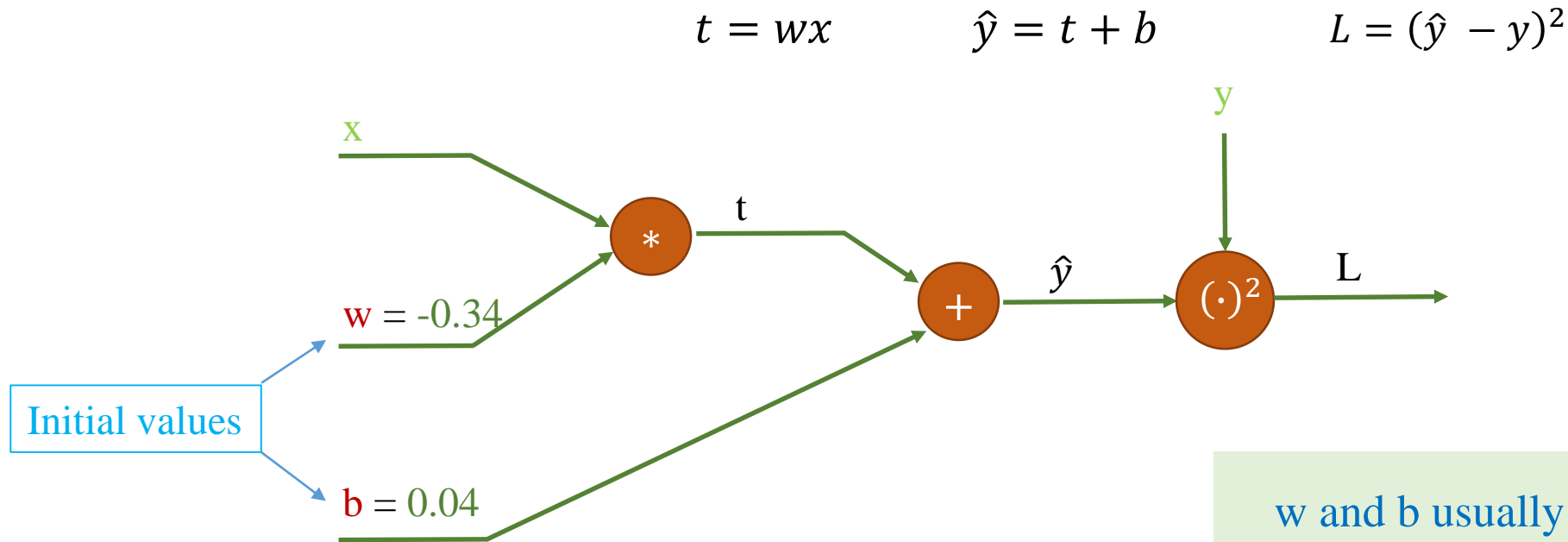
$$t = w * area$$



Computational graph

❖ House price prediction

❖ One-sample training



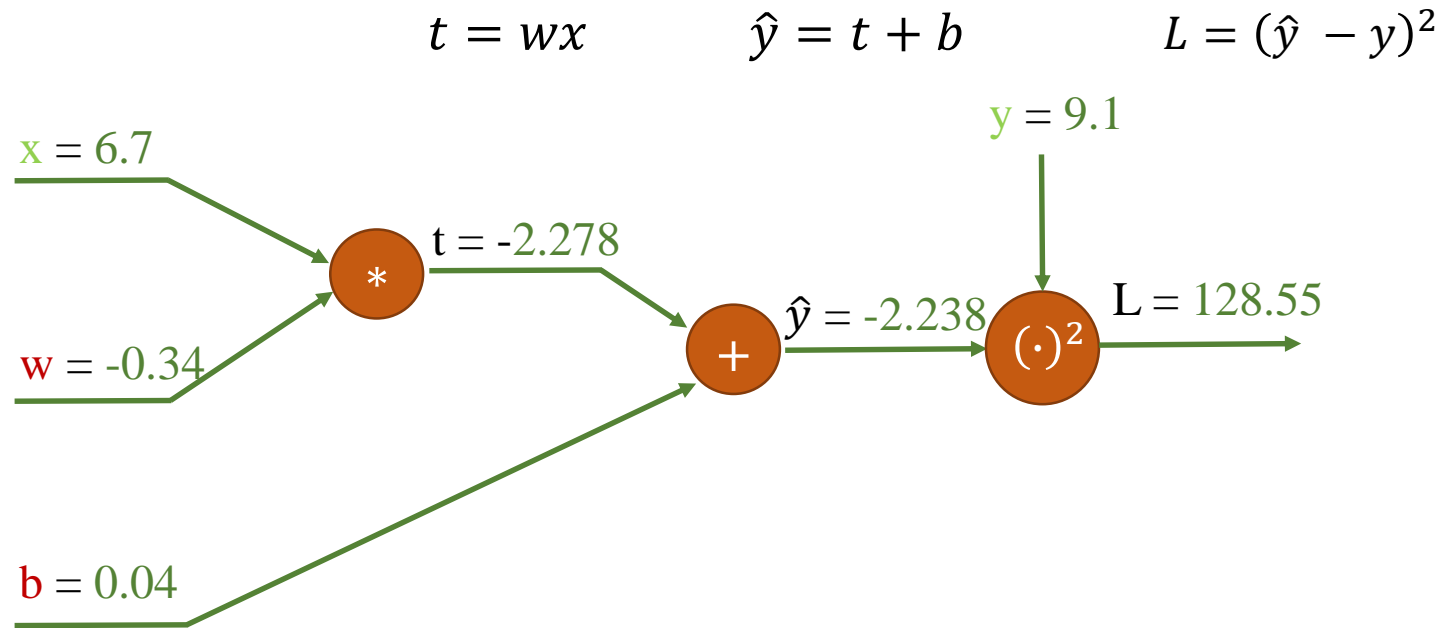
w and b usually be set by random values

For example, $N(0, \sigma)$ where σ is small

Computational graph

❖ House price prediction

❖ One-sample training

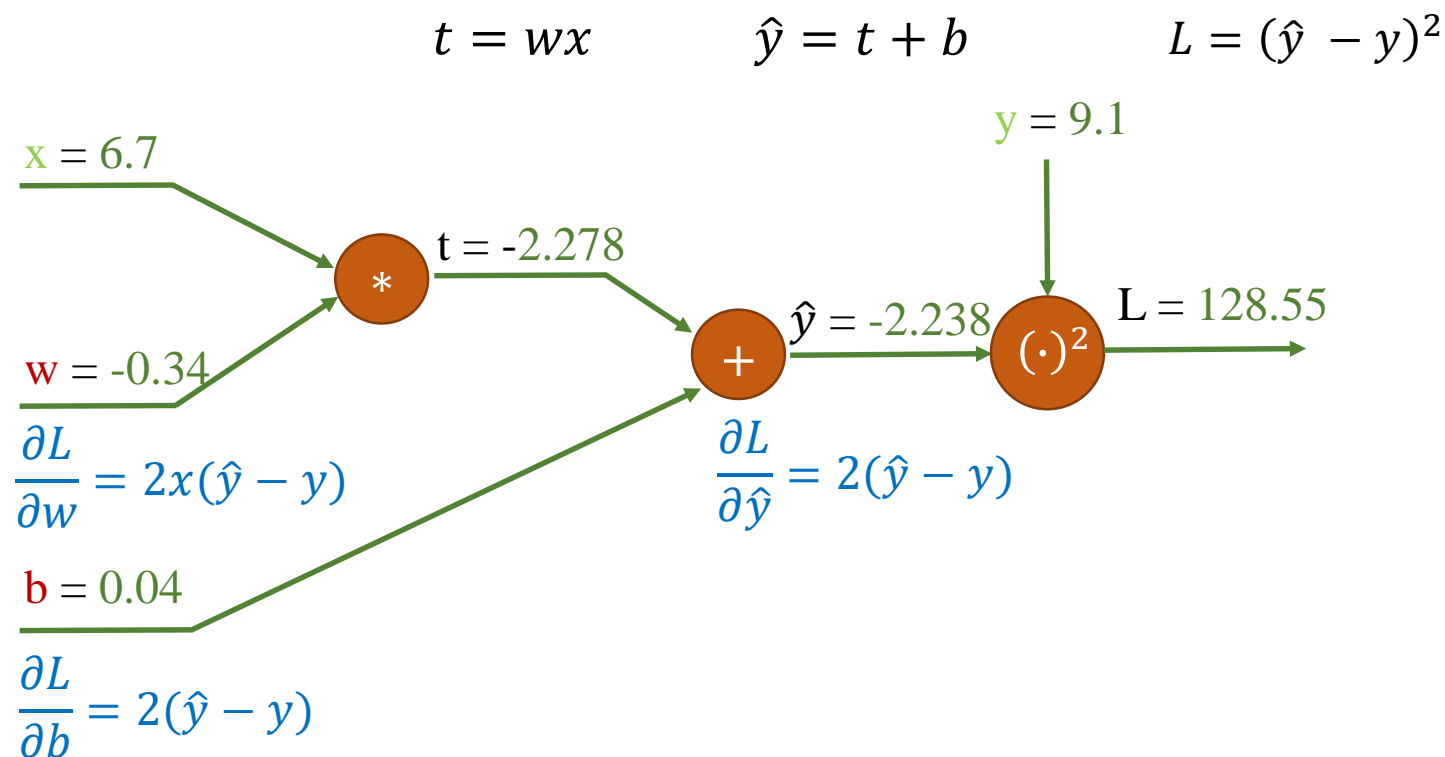


Feature		Label	
area		price	
6.7	9.1		
4.6	5.9		
3.5	4.6		
5.5	6.7		

Computational graph

❖ House price prediction

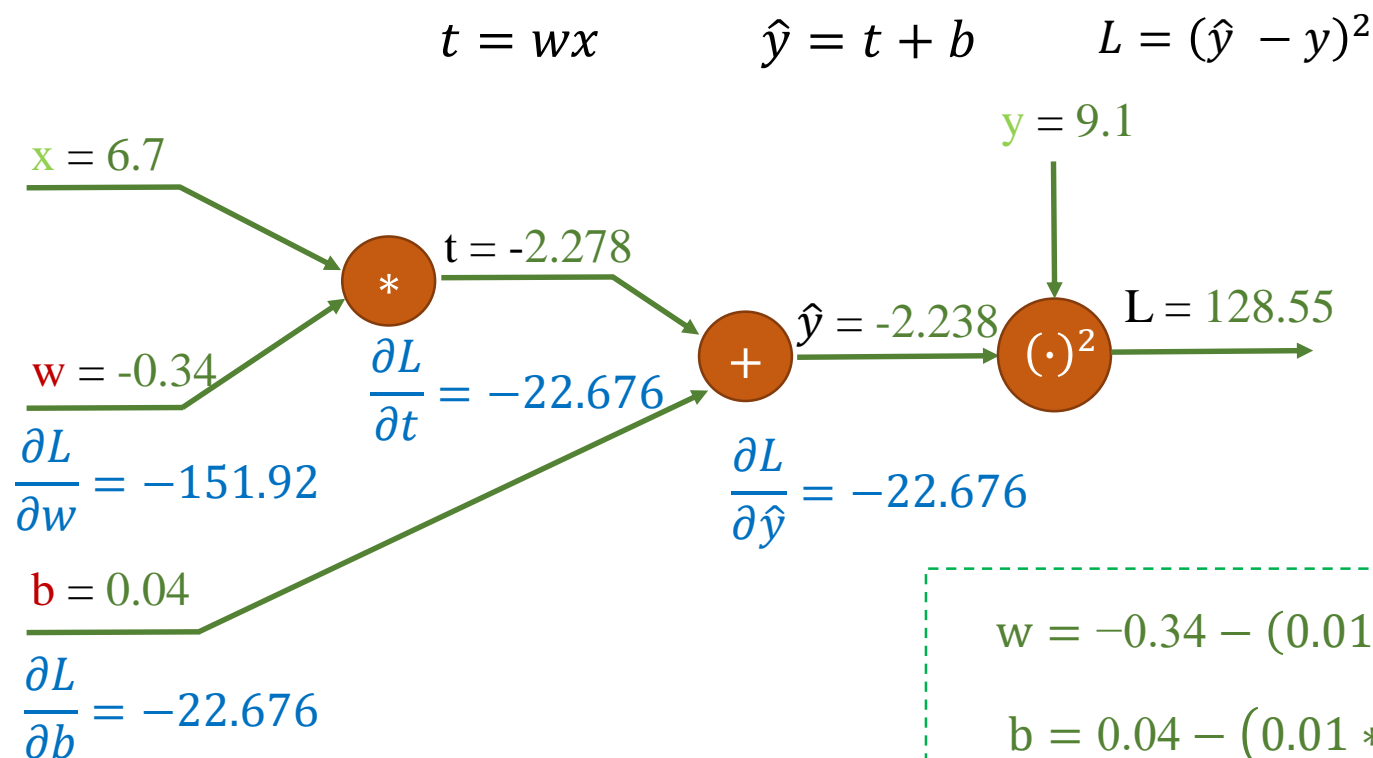
❖ One-sample training



Computational graph

❖ House price prediction

❖ One-sample training



Cách cập nhật a và b

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

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

Learning rate $\eta = 0.01$

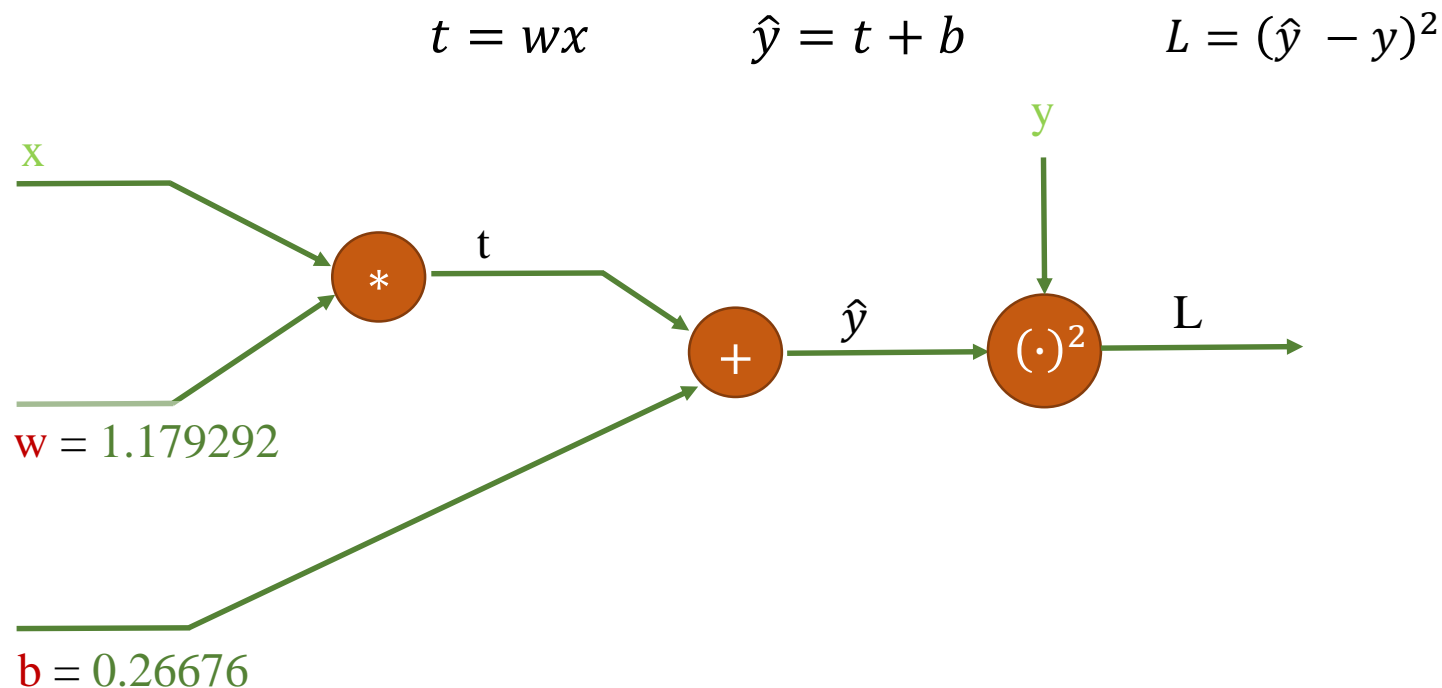
$$w = -0.34 - (0.01 * (-151.9)) = 1.179$$

$$b = 0.04 - (0.01 * (-22.67)) = 0.266$$

Computational graph

❖ House price prediction

❖ One-sample training

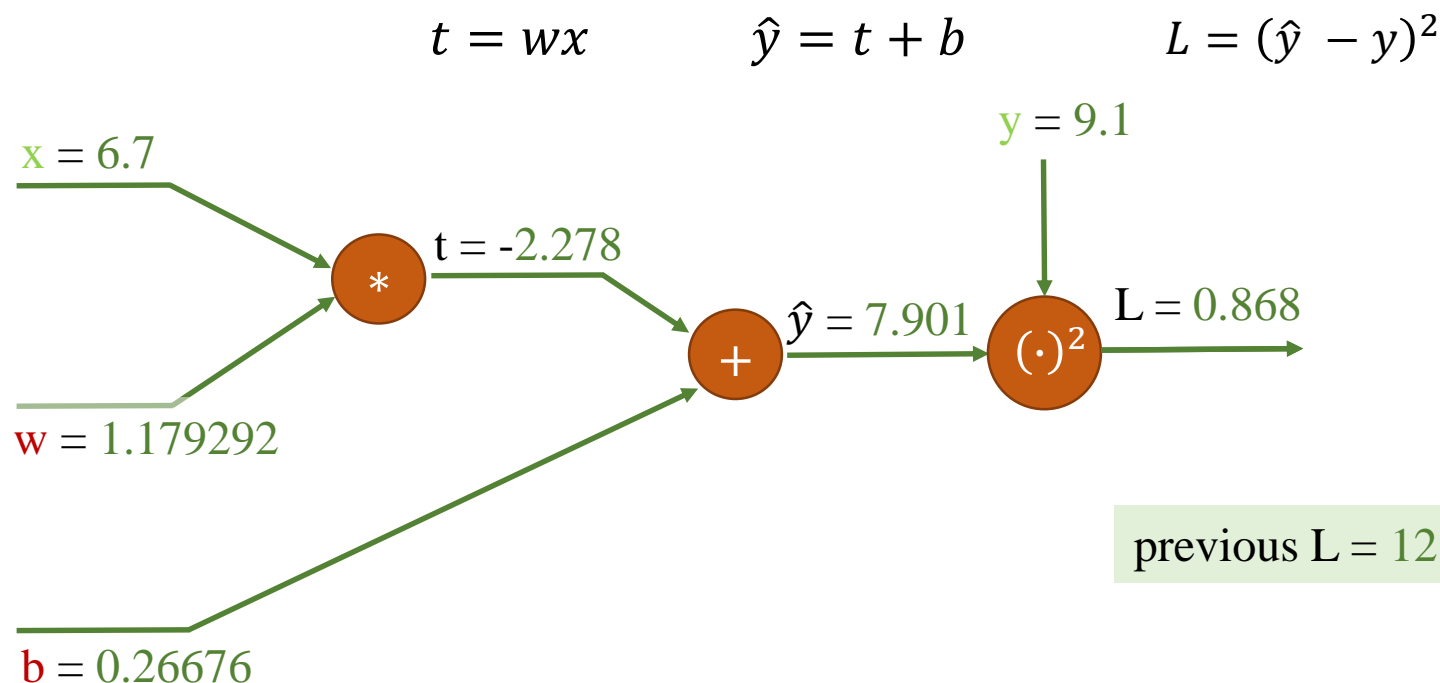


Feature		Label
area		price
6.7	9.1	
4.6	5.9	
3.5	4.6	
5.5	6.7	

Computational graph

❖ House price prediction

❖ One-sample training



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

Updated a and b values help to reduce the L value

Implementation

One-sample training

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

column index=0
column index=1

```
1 # data preparation
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 def get_column(data, index):
6     result = [row[index] for row in data]
7     return result
8
9 data = np.genfromtxt('data.csv',
10                      delimiter=',').tolist()
11
12 x_data = get_column(data, 0)
13 y_data = get_column(data, 1)
14 N = len(x_data)
15
16 print(f'areas: {x_data}')
17 print(f'prices: {y_data}')
18 print(f'data_size: {N}')
```

```
areas: [6.7, 4.6, 3.5, 5.5]
prices: [9.1, 5.9, 4.6, 6.7]
data_size: 4
```

Implementation

❖ House price prediction

❖ One-sample training

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

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

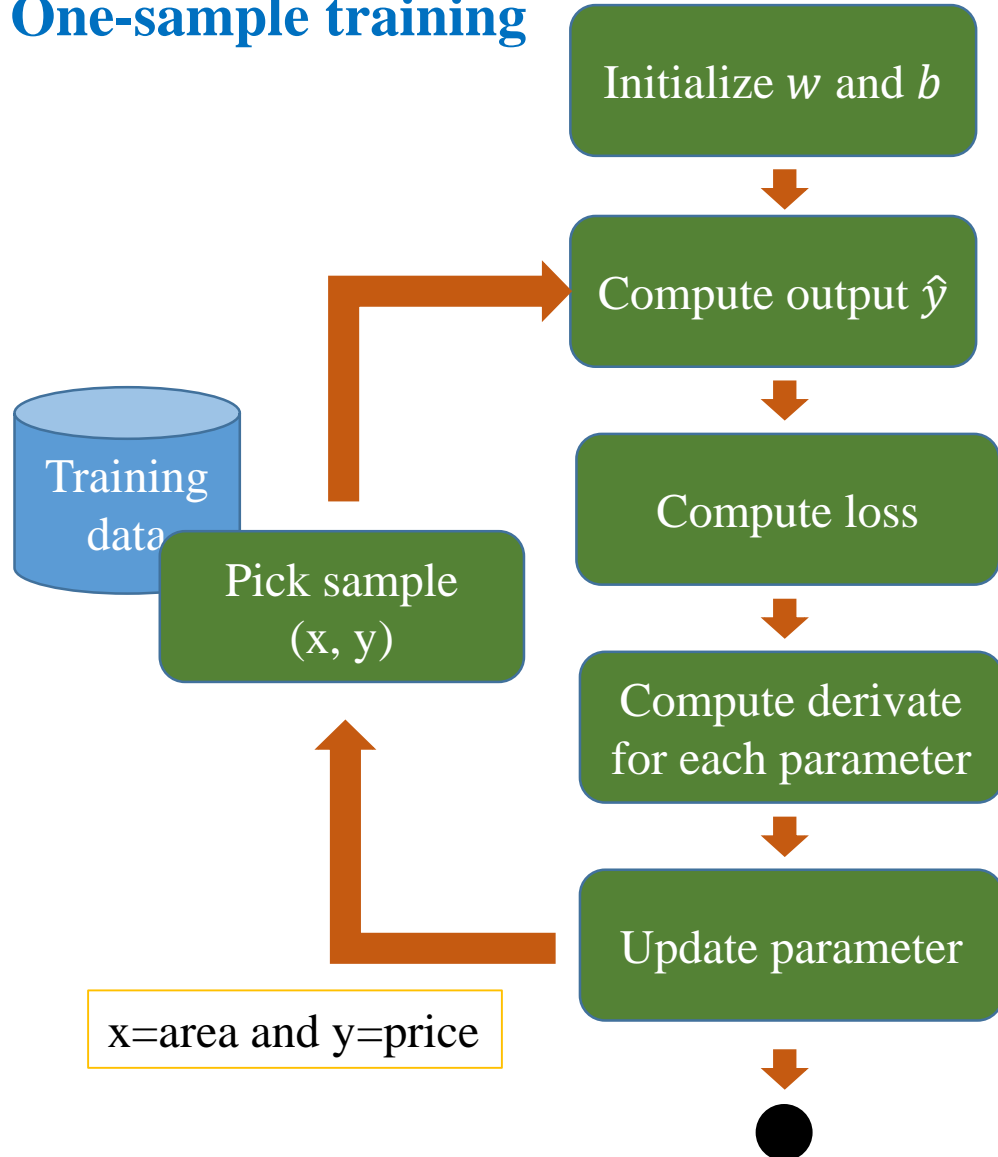
5) Update parameters

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

```
1 # forward
2 def predict(x, w, b):
3     return x*w + b
4
5 # compute gradient
6 def gradient(y_hat, y, x):
7     dw = 2*x*(y_hat-y)
8     db = 2*(y_hat-y)
9
10    return (dw, db)
11
12 # update weights
13 def update_weight(w, b, lr, dw, db):
14     w_new = w - lr*dw
15     b_new = b - lr*db
16
17    return (w_new, b_new)
```

Implementation

One-sample training

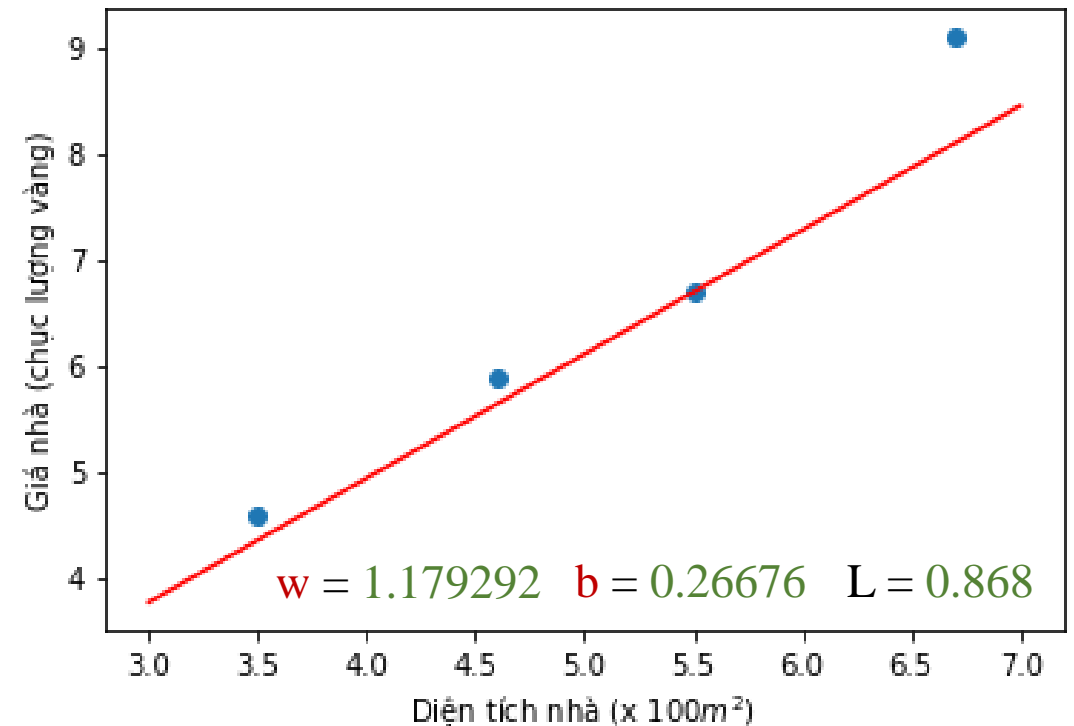
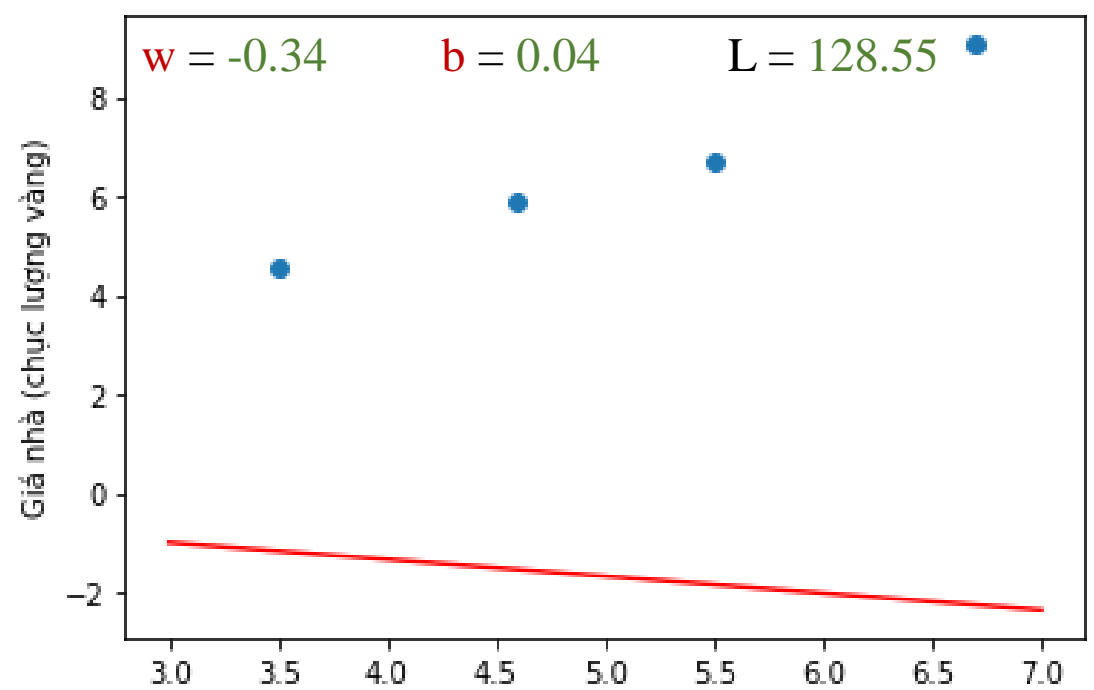


```
1 # init weights
2 b = 0.04
3 w = -0.34
4 lr = 0.01
5
6 # how long
7 epoch_max = 10
8
9 for epoch in range(epoch_max):
10     for i in range(data_size):
11         # get a sample
12         # ...
13
14         # predict y_hat
15         # ...
16
17         # compute loss
18         # ...
19
20         # compute gradient
21         # ...
22
23         # update weights
24         # ...
25
```

```

1 # init weights
2 b = 0.04
3 w = -0.34
4 lr = 0.01
5
6 # how long
7 epoch_max = 10
8 data_size = 4
9
10 for epoch in range(epoch_max):
11     for i in range(data_size):
12         # get a sample
13         x = areas[i]
14         y = prices[i]
15
16         # predict y_hat
17         y_hat = predict(x, w, b)
18
19         # compute loss
20         loss = (y_hat-y)*(y_hat-y)
21
22         # compute gradient
23         (dw, db) = gradient(y_hat, y, x)
24
25         # update weights
26         (w, b) = update_weight(w, b, lr, dw, db)

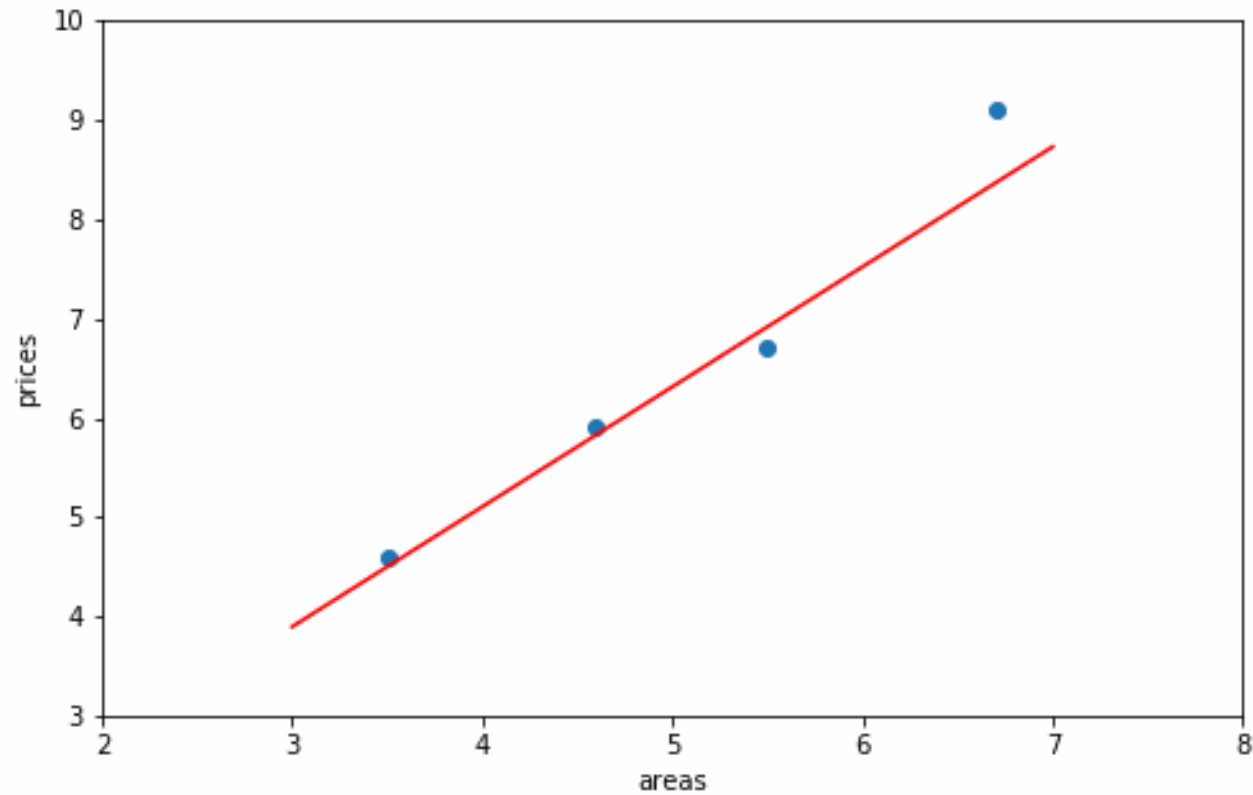
```



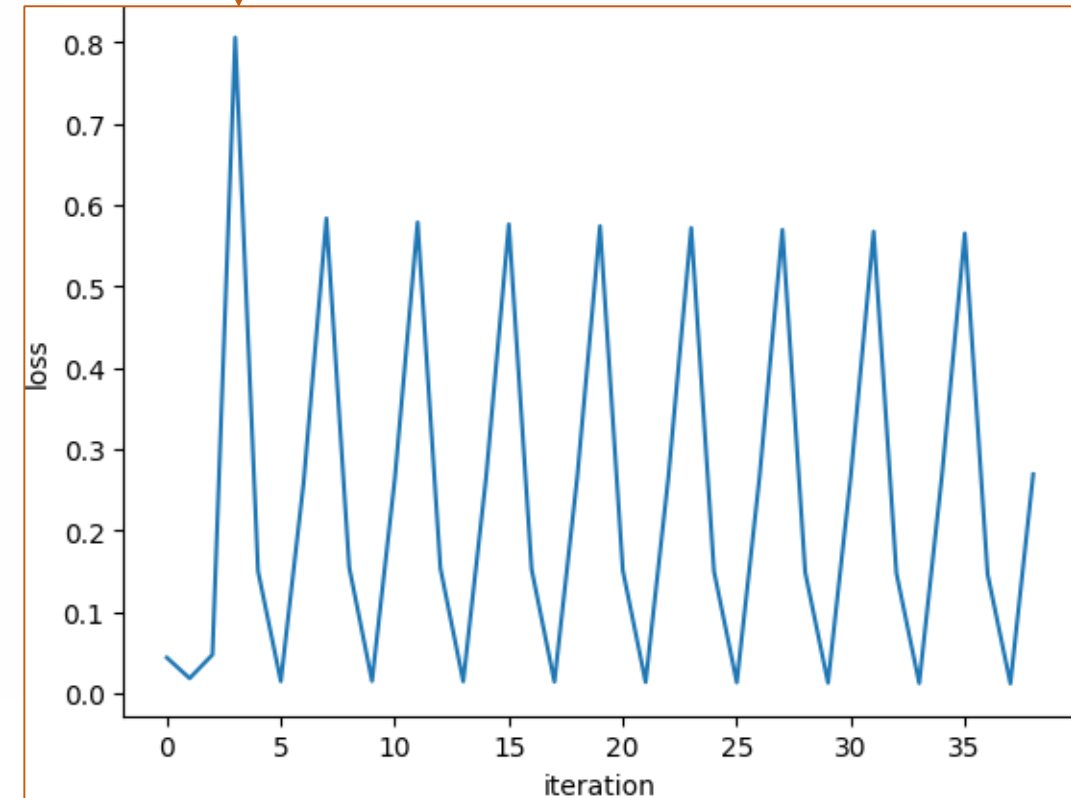
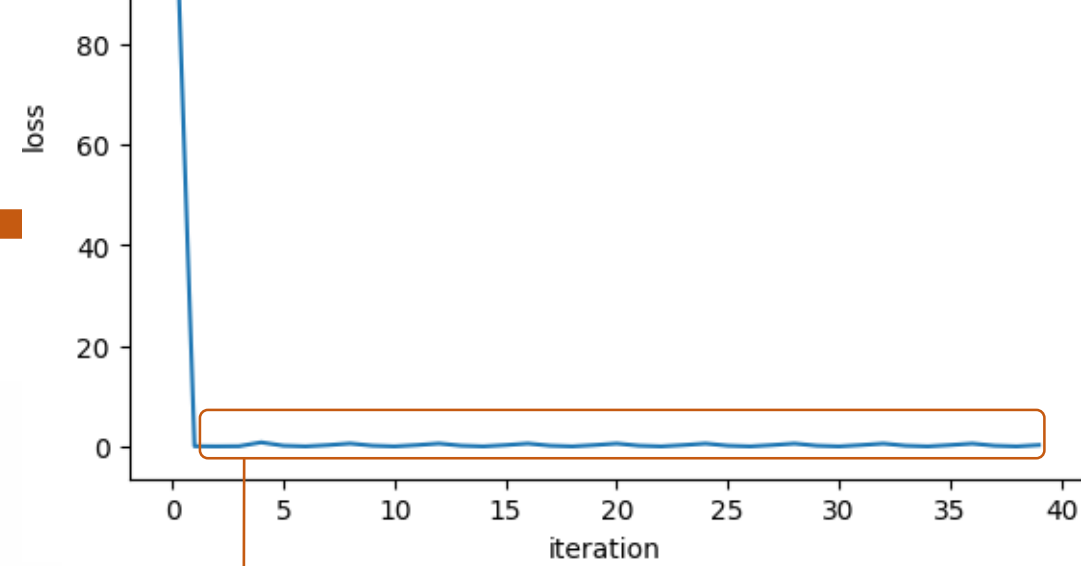
Computational graph

❖ House price prediction

❖ One-sample training



Model after training



Quiz 1: Is it OK to use the following loss function?

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

Quiz 2: if so, construct formulas and change somehow so that the two models have the same output.

Quiz 3: What about the following loss function?

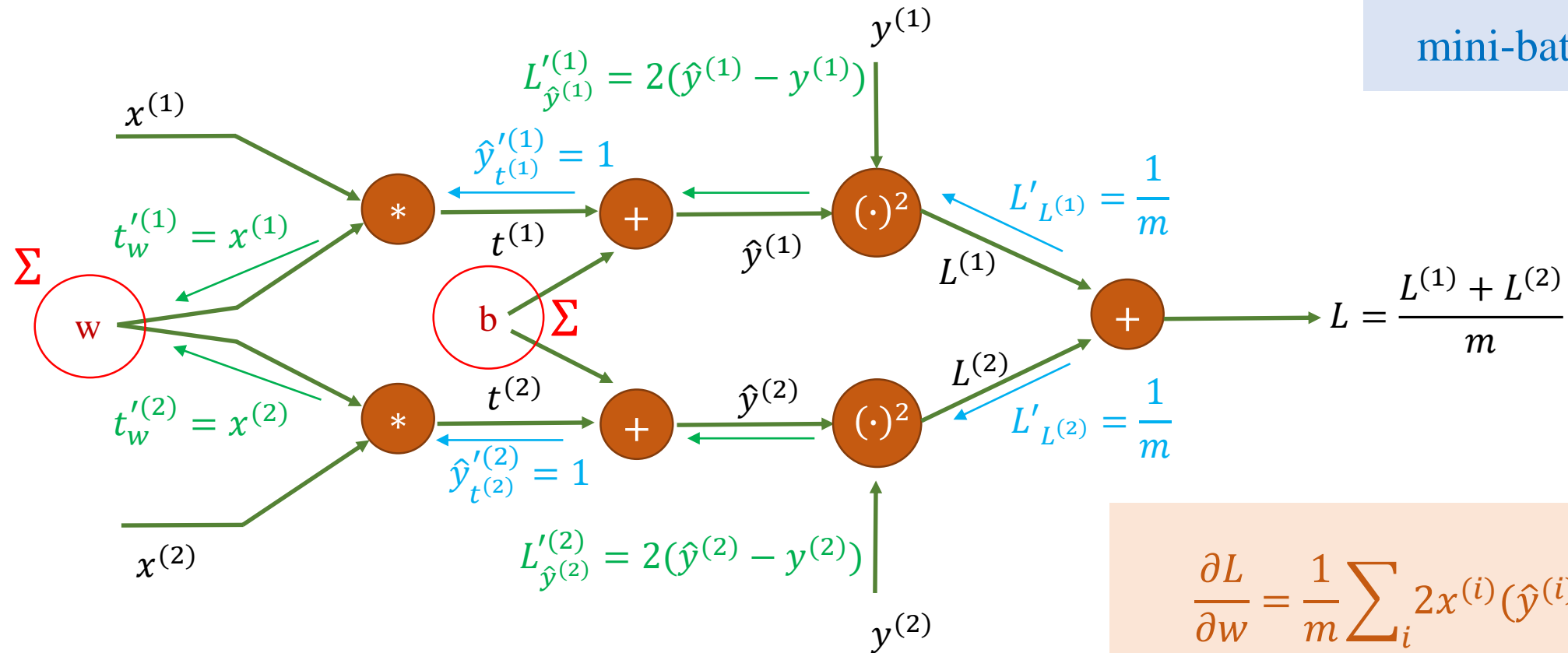
$$L = \frac{1}{2} (y - \hat{y})^2$$

Outline

- **Simple version of Linear Regression**
- **Computational Graph**
- **Mini-batch Training**
- **Batch Training**
- **Generalization of Linear Regression**
- **Loss Functions**

Computational graph

❖ Compute derivate for w and b



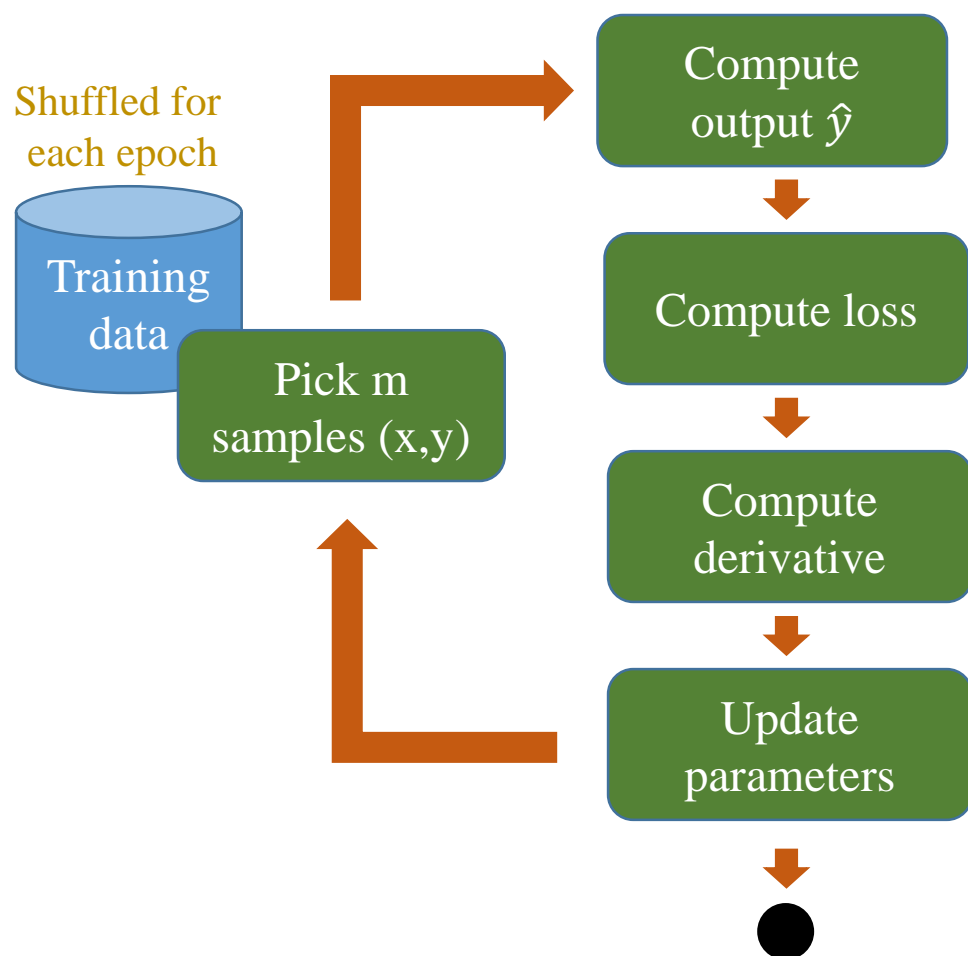
$$\frac{\partial L}{\partial w} = \frac{1}{m} \sum_i 2x^{(i)} (\hat{y}^{(i)} - y^{(i)})$$

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

Computational graph

❖ House price prediction

❖ m-sample training ($1 < m < N$)



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

2) Tính output $\hat{y}^{(i)}$

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

3) Tính loss

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

4) Tính đạo hàm

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

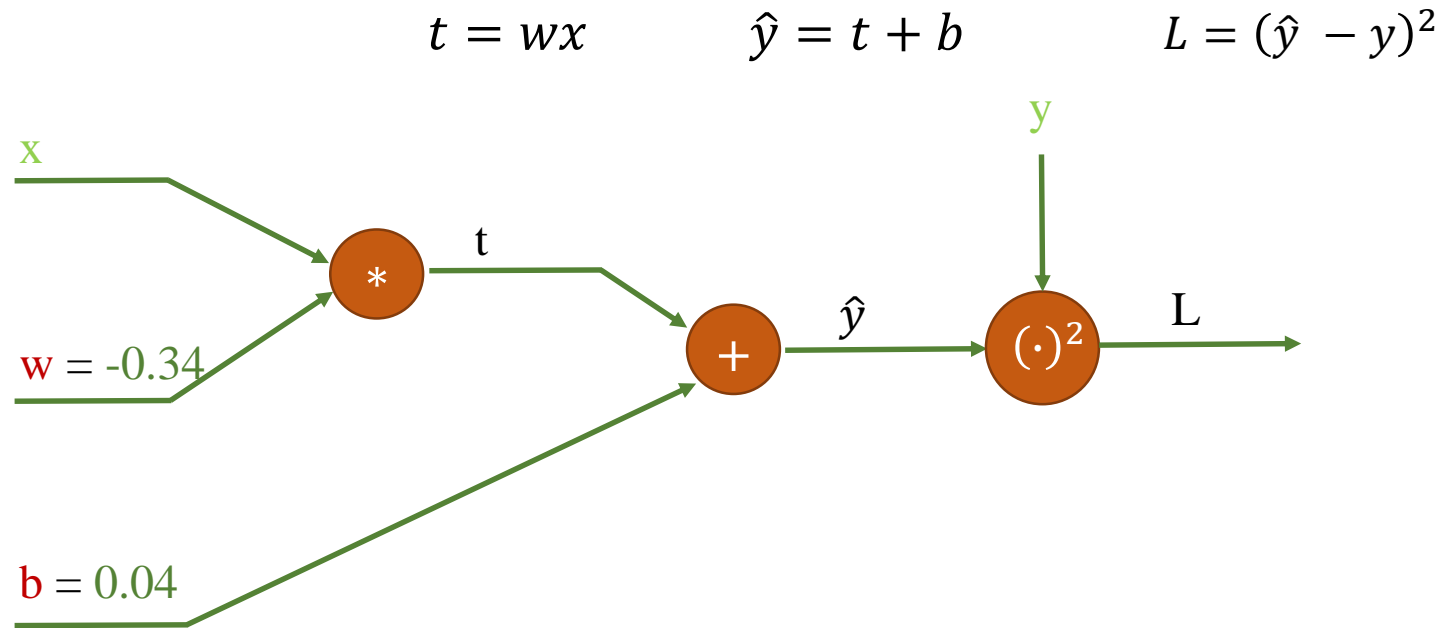
5) Cập nhật tham số

$$w = w - \eta \frac{\sum_i \frac{\partial L^{(i)}}{\partial w}}{m} \quad b = b - \eta \frac{\sum_i \frac{\partial L^{(i)}}{\partial b}}{m}$$

Computational graph

❖ House price prediction

❖ m-sample training ($1 < m < N$)



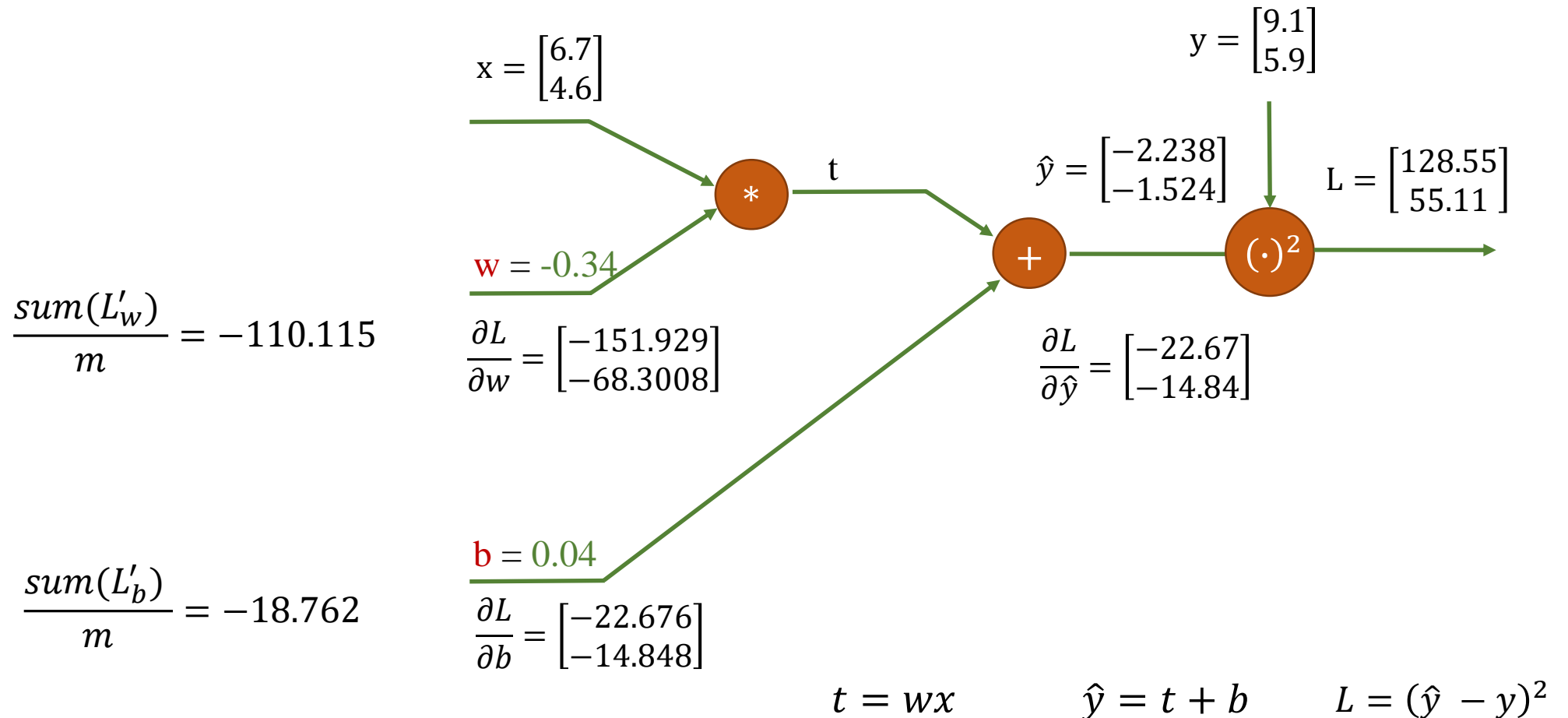
Computational graph

❖ House price prediction

❖ m-sample training ($1 < m < N$)

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

$m = 2$



$$\frac{\text{sum}(L'_w)}{m} = -110.115$$

$$\frac{\text{sum}(L'_b)}{m} = -18.762$$

Computational graph

❖ House price prediction

❖ m-sample training ($1 < m < N$)

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

$m = 2$

Update w and b

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

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

Learning rate $\eta = 0.01$

replace

replace

$$x = \begin{bmatrix} 6.7 \\ 4.6 \end{bmatrix}$$

$$w = -0.34$$

$$\frac{\partial L}{\partial w} = -110.115$$

$$b = 0.04$$

$$\frac{\partial L}{\partial b} = -18.762$$

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

*

t

+

$(\cdot)^2$

$$t = wx$$

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

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

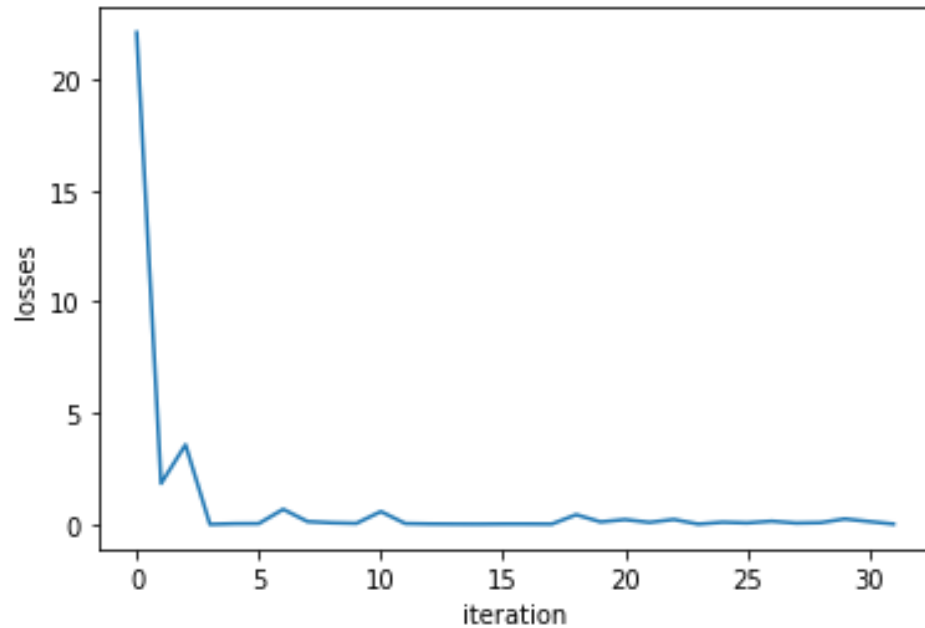
$$w = -0.34 - (0.01 * (-110.115)) = 0.761$$

$$b = 0.04 - (0.01 * (-18.762)) = 0.227$$

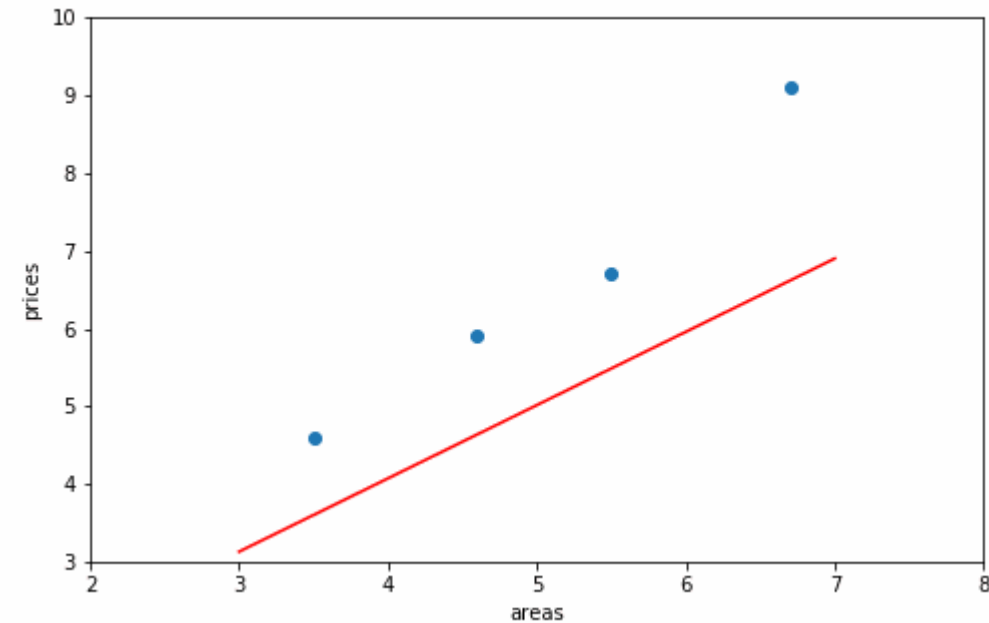
Computational graph

❖ House price prediction

❖ m-sample training ($1 < m < N$)



Losses for 30 iterations

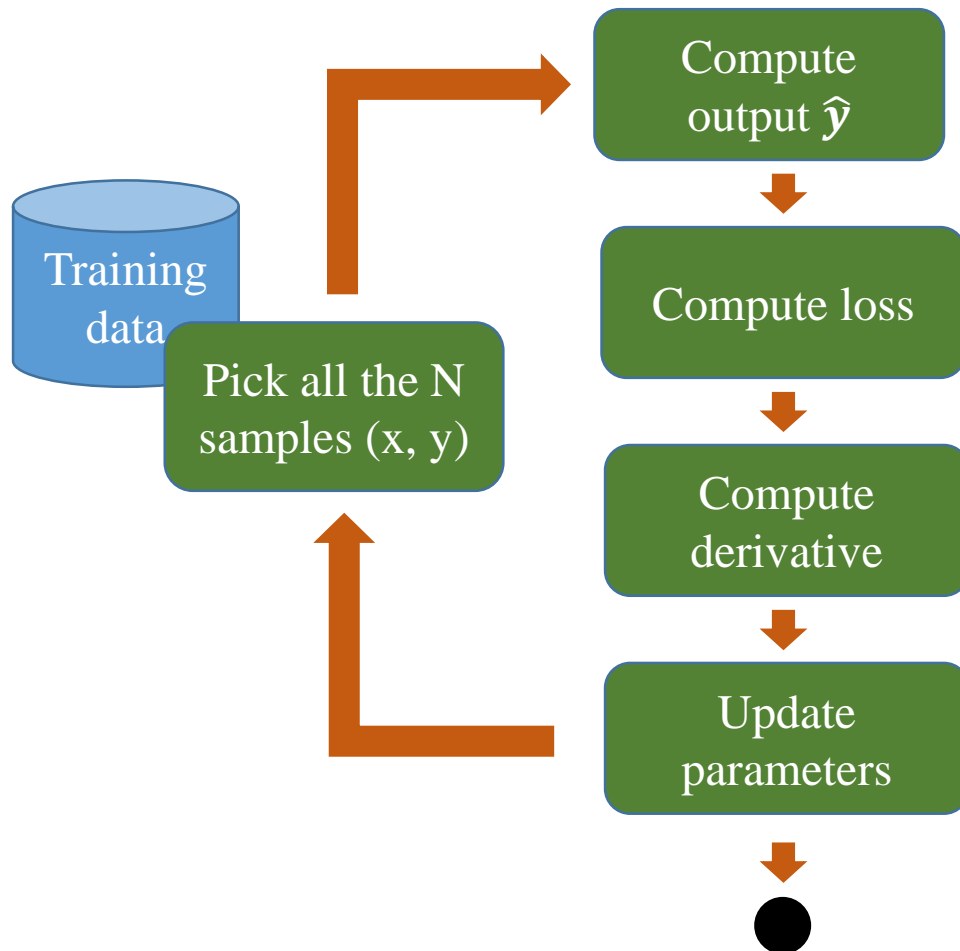


Model updating for different iterations

Computational graph

❖ House price prediction

❖ N-sample training



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

2) Tính output $\hat{y}^{(i)}$

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

3) Tính loss

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

4) Tính đạo hàm

$$\begin{aligned} \frac{\partial L^{(i)}}{\partial w} &= 2x^{(i)}(\hat{y}^{(i)} - y^{(i)}) \\ \frac{\partial L^{(i)}}{\partial b} &= 2(\hat{y}^{(i)} - y^{(i)}) \end{aligned} \quad \text{for } 0 \leq i < N$$

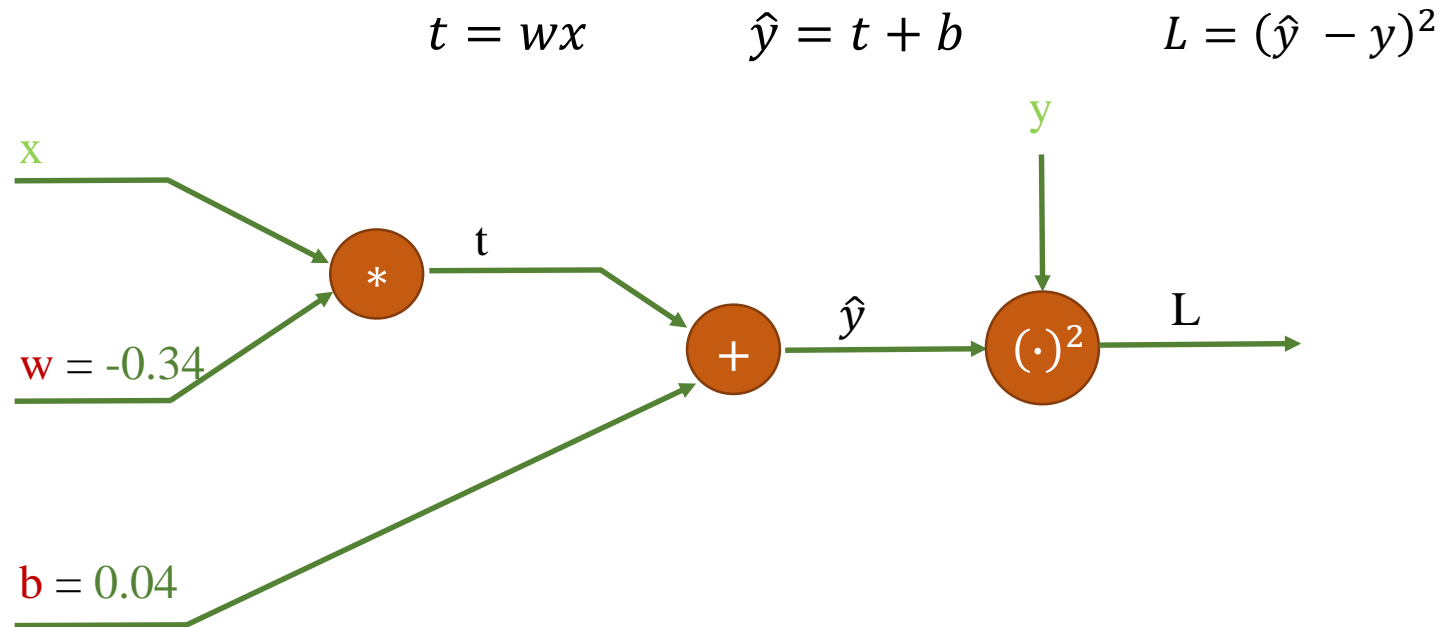
5) Cập nhật tham số

$$w = w - \eta \frac{\sum_i \frac{\partial L^{(i)}}{\partial w}}{N} \quad b = b - \eta \frac{\sum_i \frac{\partial L^{(i)}}{\partial b}}{N}$$

Computational graph

❖ House price prediction

❖ N-sample training



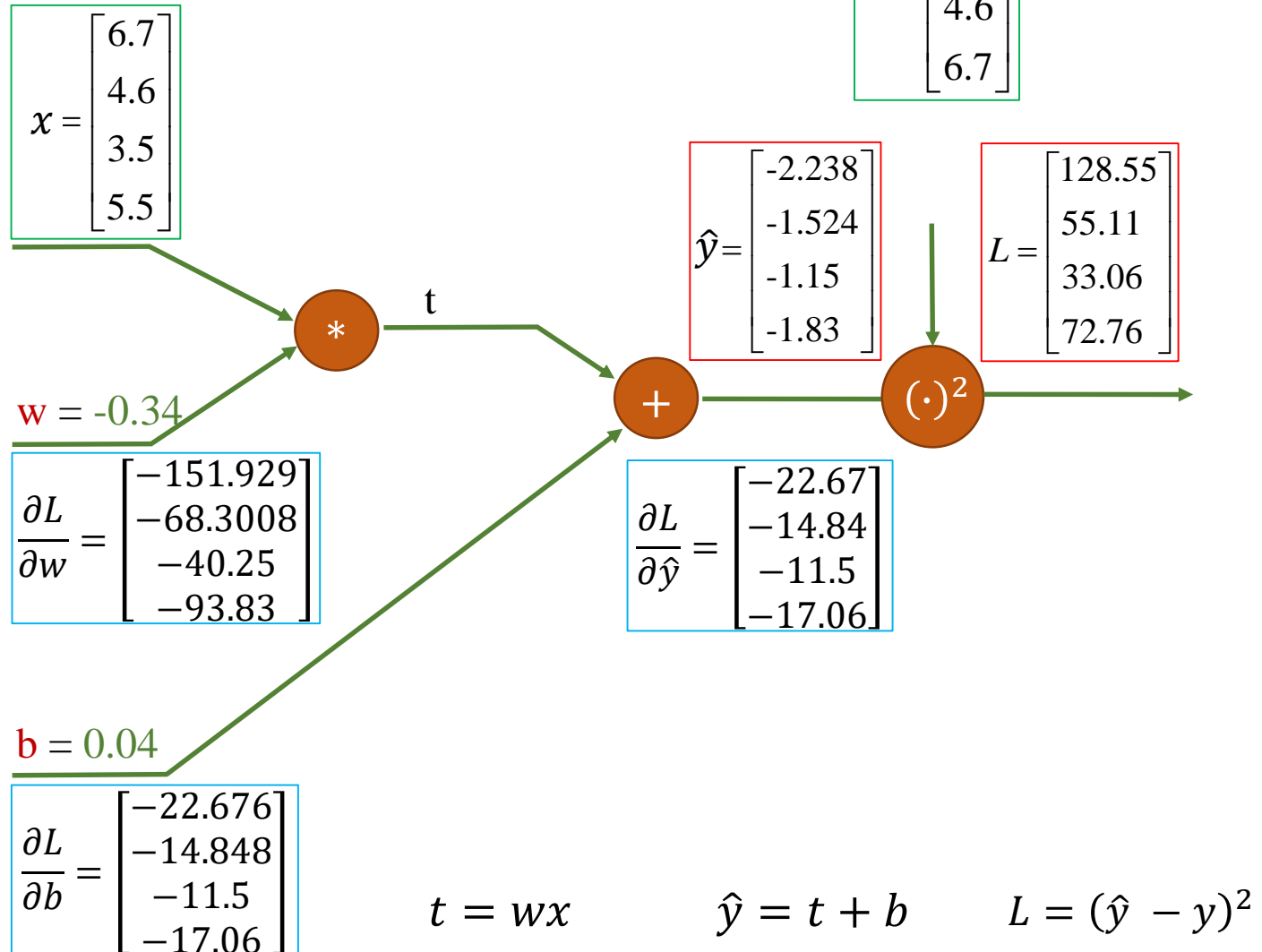
Computational graph

❖ House price prediction

❖ N-sample training

$$\frac{\text{sum}(\frac{\partial L}{\partial w})}{4} = -88.5775$$

$$\frac{\text{sum}(\frac{\partial L}{\partial b})}{4} = -16.521$$



Computational graph

❖ House price prediction

❖ N-sample training

Update w and b

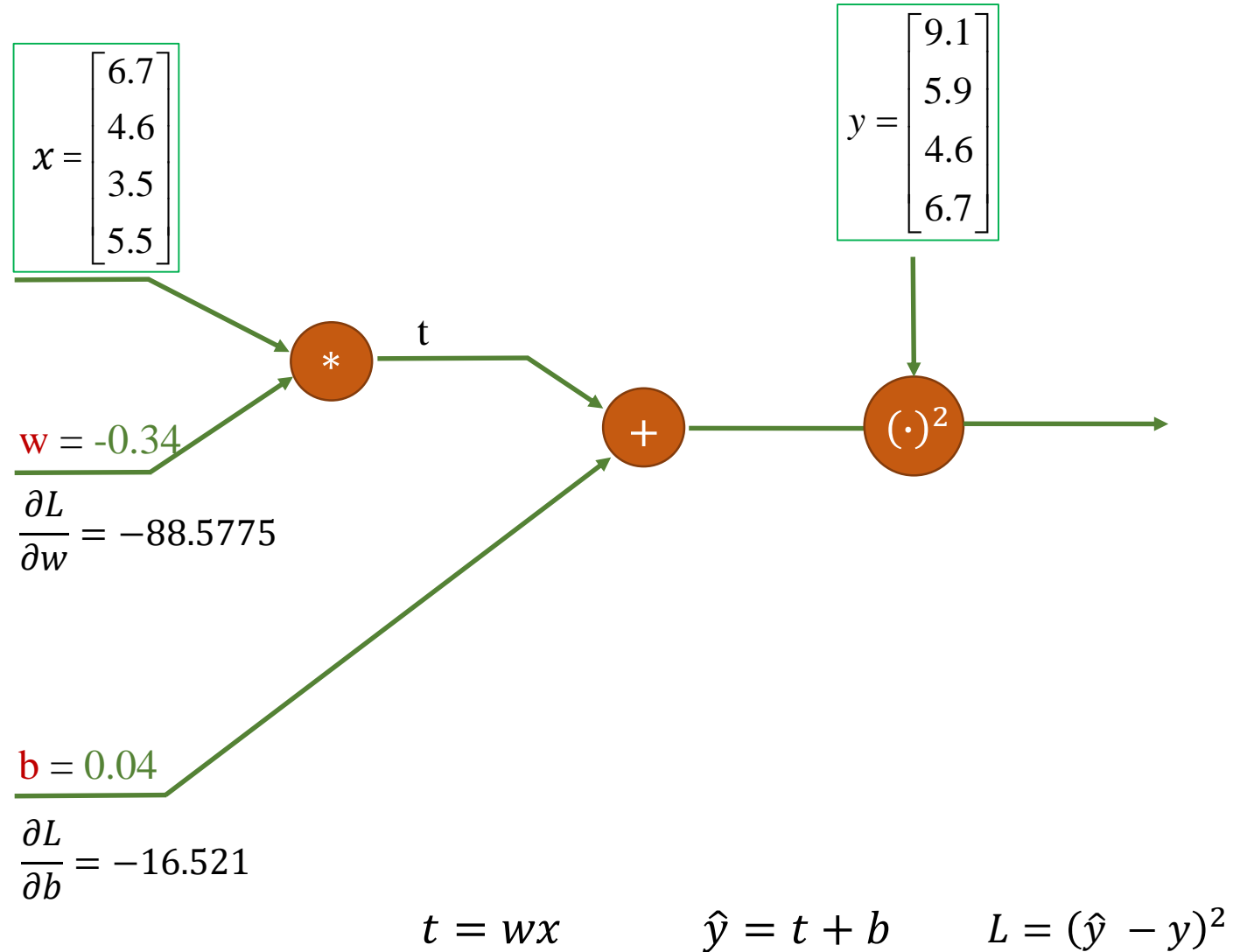
$$w = w - \eta * \frac{\partial L}{\partial w}$$

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

Learning rate $\eta = 0.01$

$$w = -0.34 - (0.01 * (-88.5775)) = 0.54$$

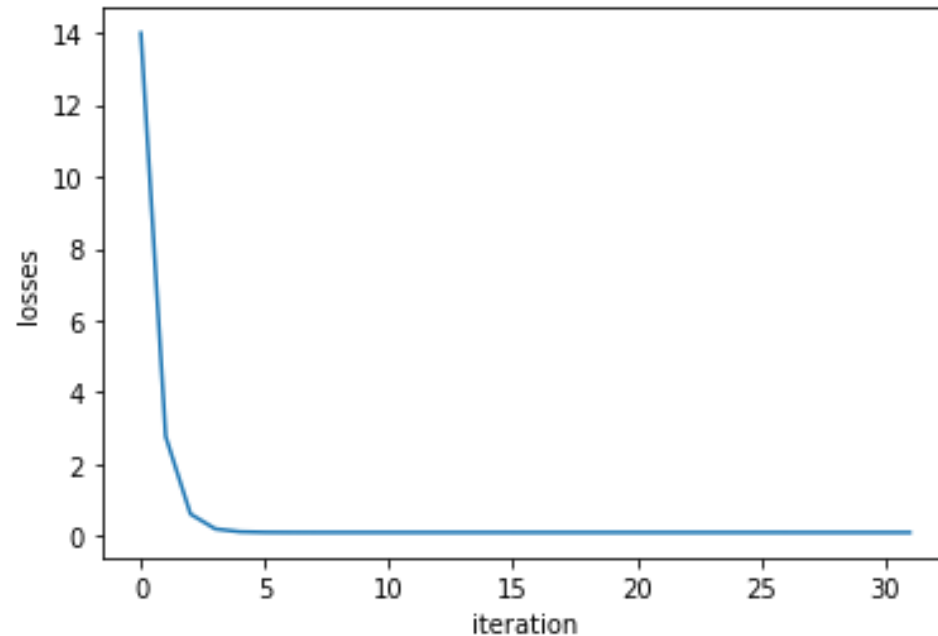
$$b = 0.04 - (0.01 * (-16.521)) = 0.205$$



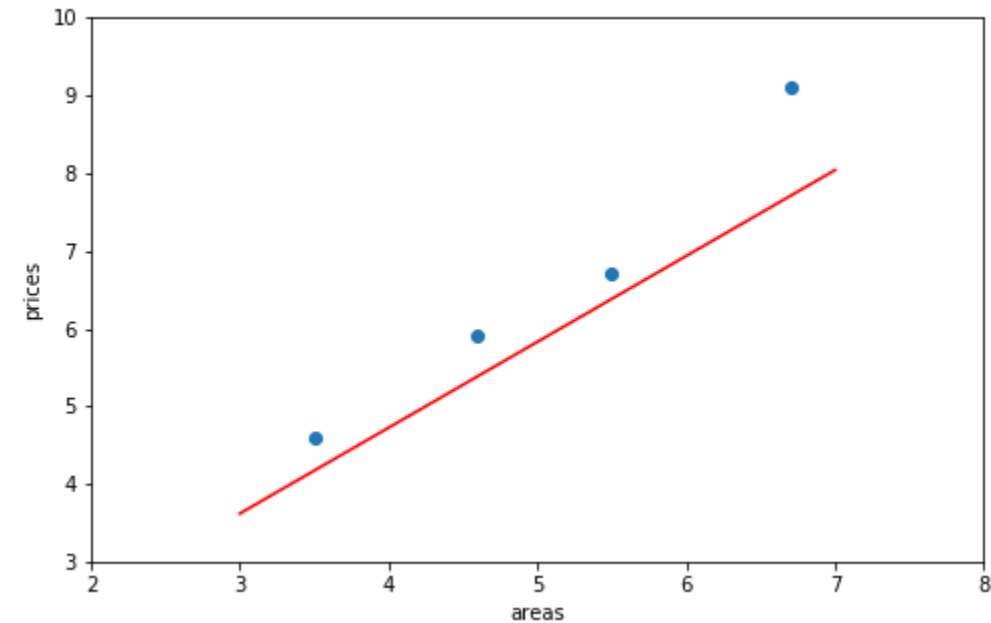
Computational graph

❖ House price prediction

❖ N-sample training



Losses for 30 iterations



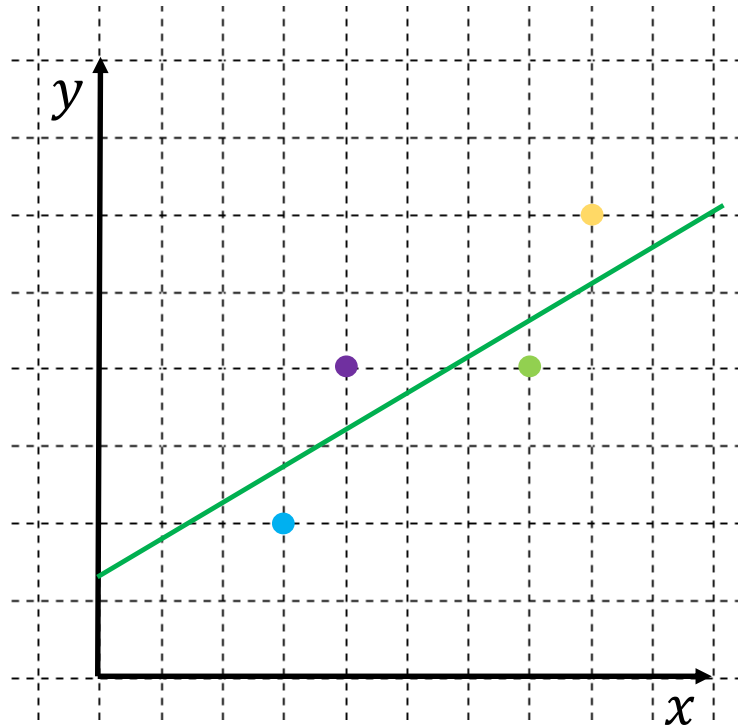
Model updating for different iterations

Outline

- **Simple version of Linear Regression**
- **Computational Graph**
- **Mini-batch Training**
- **Batch Training**
- **Generalization of Linear Regression**
- **Loss Functions**

Linear Regression

❖ General formula



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

House price data

Linear regression models ← Linear equations

Linear equation = $w_1x_1 + w_2x_2 + \dots + w_nx_n + b$

Linear Regression

❖ General formula

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

❖ General formula

Boston House
Price Data

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	

$$\text{medv} = w_1 * x_1 + \cdots + w_{13} * x_{13} + b$$

Linear Regression

❖ Generalized formula

House price data

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

Model

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

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

Model (vectorization)

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} \quad \text{where} \quad \boldsymbol{\theta}^T = [b \quad w]^T$$

$$\mathbf{x} = [x_0 \quad \text{area}]^T$$

$$x_0 = 1$$

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

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

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

Model (vectorization)

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} \quad \text{where} \quad \boldsymbol{\theta}^T = [b \quad w_1 \quad w_2 \quad w_3]^T$$

$$\mathbf{x} = [x_0 \quad TV \quad Radio \quad Newspaper]^T$$

$$x_0 = 1$$

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
```

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

```

1 def initialize_params():
2     w1 = random.gauss(mu=0.0, sigma=0.01)
3     w2 = random.gauss(mu=0.0, sigma=0.01)
4     w3 = random.gauss(mu=0.0, sigma=0.01)
5     b = 0
6
7     return w1, w2, w3, b
8
9 # initialize model's parameters
10 w1, w2, w3, b = initialize_params()
11 print(w1, w2, w3, b)

```

0.01609506469549467 0.00607778501208891 0.0023344573891806507 0

```

1 import numpy as np
2 import random
3
4 def get_column(data, index):
5     result = [row[index] for row in data]
6     return result
7
8 data = np.genfromtxt('advertising.csv',
9                     delimiter=',',
10                    skip_header=1).tolist()
11
12 # get tv (index=0)
13 tv_data = get_column(data, 0)
14
15 # get radio (index=1)
16 radio_data = get_column(data, 1)
17
18 # get newspaper (index=2)
19 newspaper_data = get_column(data, 2)
20
21 # get sales (index=3)
22 sales_data = get_column(data, 3)

```

Unnormalized data $\eta = 10^{-5}$

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

0.01609506469549467 0.00607778501208891 0.0023344573891806507 0

x1: 230.1

x2: 37.8

x1: 69.2

y: 22.1

y_hat: 4.0947591112215855

d1_dw1: -8286.011857015827

d1_dw2: -1361.1962111916482

d1_dw3: -2491.925339006933

d1_db: -36.01048177755683

w1: 0.09895518326565295

w2: 0.019689747124005393

w3: 0.027253710779249984

b: 0.0003601048177755684

```
1 for epoch in range(epoch_max):
2     for i in range(N):
3         # get a sample
4         x1 = tv_data[i]
5         x2 = radio_data[i]
6         x3 = newspaper_data[i]
7         y = sales_data[i]
8
9         # compute output
10        y_hat = predict(x1, x2, x3, w1, w2, w3, b)
11
12        # compute gradient w1, w2, w3, b
13        dl_dw1 = compute_gradient_wi(x1, y, y_hat)
14        dl_dw2 = compute_gradient_wi(x2, y, y_hat)
15        dl_dw3 = compute_gradient_wi(x3, y, y_hat)
16        dl_db = compute_gradient_b(y, y_hat)
17
18        # update parameters
19        w1 = update_weight_wi(w1, dl_dw1, lr)
20        w2 = update_weight_wi(w2, dl_dw2, lr)
21        w3 = update_weight_wi(w3, dl_dw3, lr)
22        b = update_weight_b(b, dl_db, lr)
```

Normalized data

$$\eta = 10^{-2}$$

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

0.01609506469549467 0.00607778501208891 0.0023344573891806507 0

x1: 0.5504267881241568
 x2: -0.09835863697705782
 x1: 0.007579284750337614
 y: 22.1
 y_hat: 0.008279045632653116

d1_dw1: -24.319750018095103
 d1_dw2: 4.345823123098159
 d1_dw3: -0.33487888747630074
 d1_db: -44.1834419087347

w1: 0.2592925648764457
 w2: -0.037380446218892686
 w3: 0.005683246263943658
 b: 0.441834419087347

$$x = \frac{x - x_{mean}}{x_{max} - x_{min}}$$

```

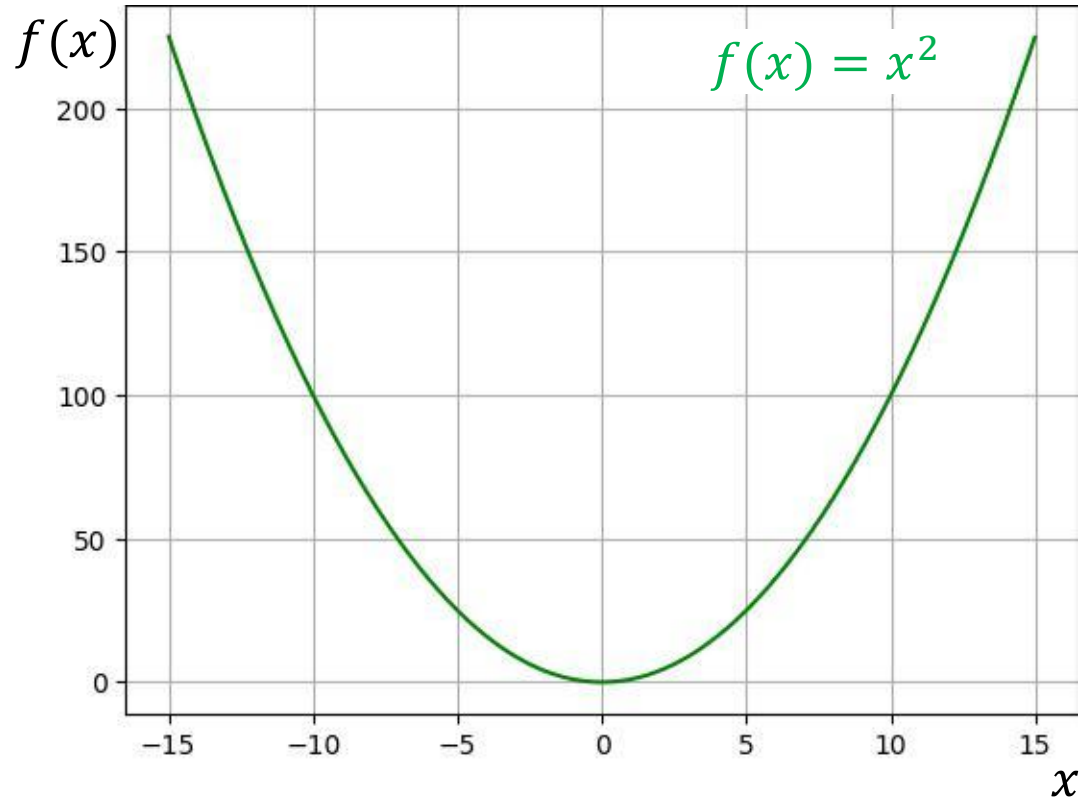
1  for epoch in range(epoch_max):
2      for i in range(N):
3          # get a sample
4          x1 = tv_data[i]
5          x2 = radio_data[i]
6          x3 = newspaper_data[i]
7          y = sales_data[i]
8
9          # compute output
10         y_hat = predict(x1, x2, x3, w1, w2, w3, b)
11
12         # compute gradient w1, w2, w3, b
13         dl_dw1 = compute_gradient_wi(x1, y, y_hat)
14         dl_dw2 = compute_gradient_wi(x2, y, y_hat)
15         dl_dw3 = compute_gradient_wi(x3, y, y_hat)
16         dl_db = compute_gradient_b(y, y_hat)
17
18         # update parameters
19         w1 = update_weight_wi(w1, dl_dw1, lr)
20         w2 = update_weight_wi(w2, dl_dw2, lr)
21         w3 = update_weight_wi(w3, dl_dw3, lr)
22         b = update_weight_b(b, dl_db, lr)
  
```


Outline

- **Simple version of Linear Regression**
- **Computational Graph**
- **Mini-batch Training**
- **Batch Training**
- **Generalization of Linear Regression**
- **Loss Functions**

Loss Functions

❖ Mean Squared Error (MSE)



$$f'(x) = 2x$$

One sample

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

N samples

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

I	y	\hat{y}	L
area	price	prediction	error
6.7	9.1	5.5	12.9
4.6	5.9	3.9	4.41
3.5	4.6	3.1	2.25
5.5	6.7	4.6	4.41

Loss Functions

Feature Label

area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

❖ Mean Squared Error (MSE)

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

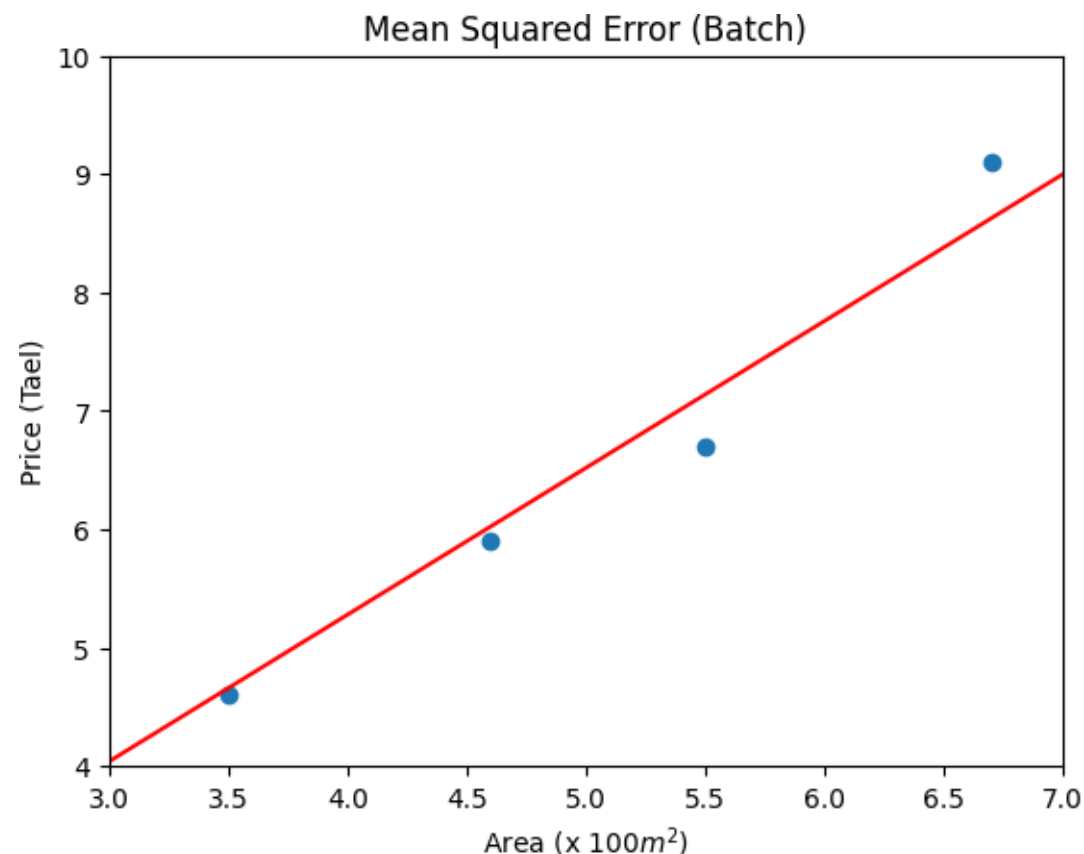
$$\frac{\partial L}{\partial w} = 2x(\hat{y} - y) \quad \frac{\partial L}{\partial b} = 2(\hat{y} - y)$$

5) Update parameters

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

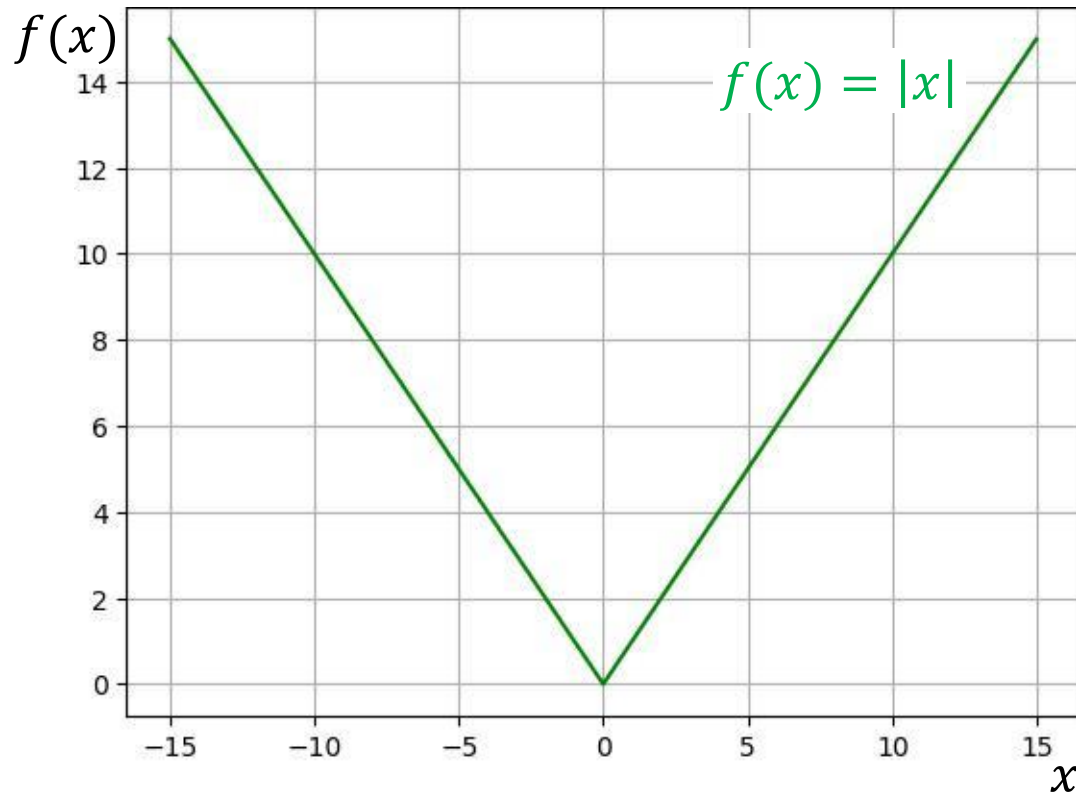
$$w = 1.207$$

$$b = 0.251$$



Loss Functions

❖ Mean Absolute Error (MAE)



$$f'(x) = \frac{x}{|x|} \quad \text{for } x \neq 0$$

One sample

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

N samples

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - y_i|$$

I	y	\hat{y}	L
area	price	prediction	error
6.7	9.1	5.5	3.6
4.6	5.9	3.8	2.1
3.5	4.6	3.1	1.5
5.5	6.7	4.6	2.1

Loss Functions

Feature Label

area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

❖ Mean Absolute Error (MAE)

$$w = 1.185$$

$$b = 0.340$$

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

2) Compute the output \hat{y}

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

3) Compute loss

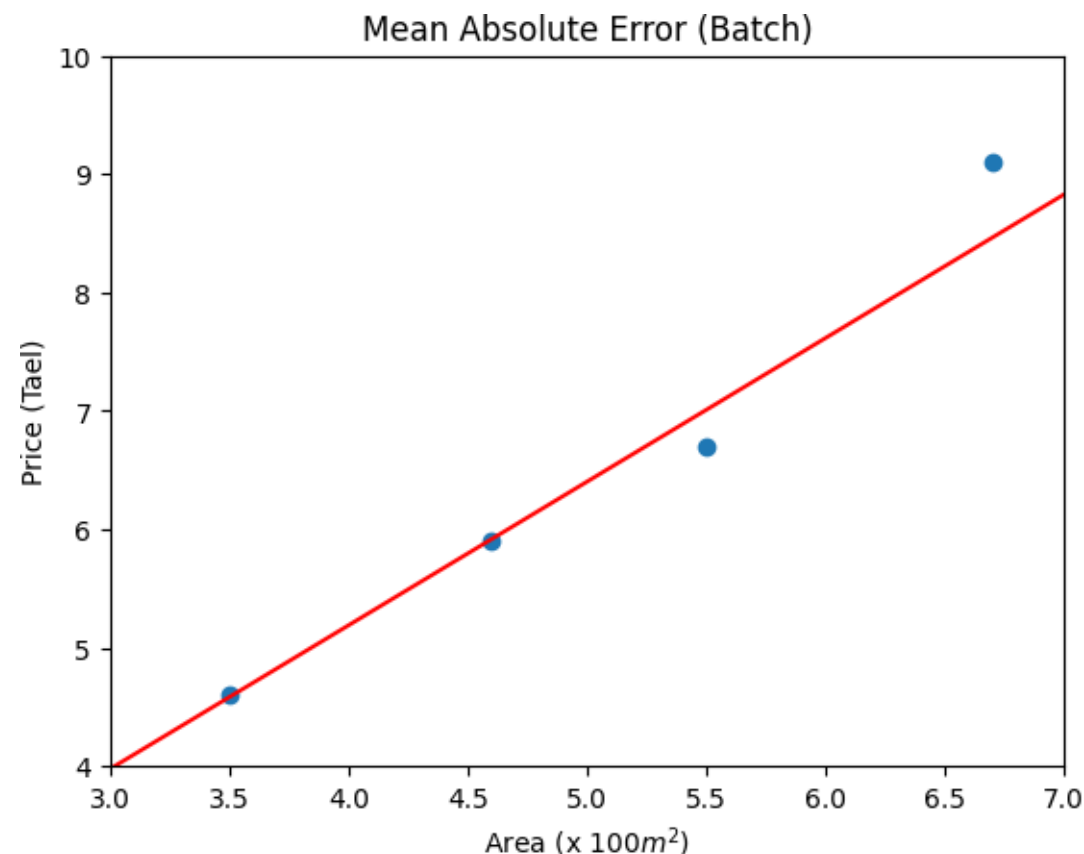
$$L = |\hat{y} - y|$$

4) Compute derivative

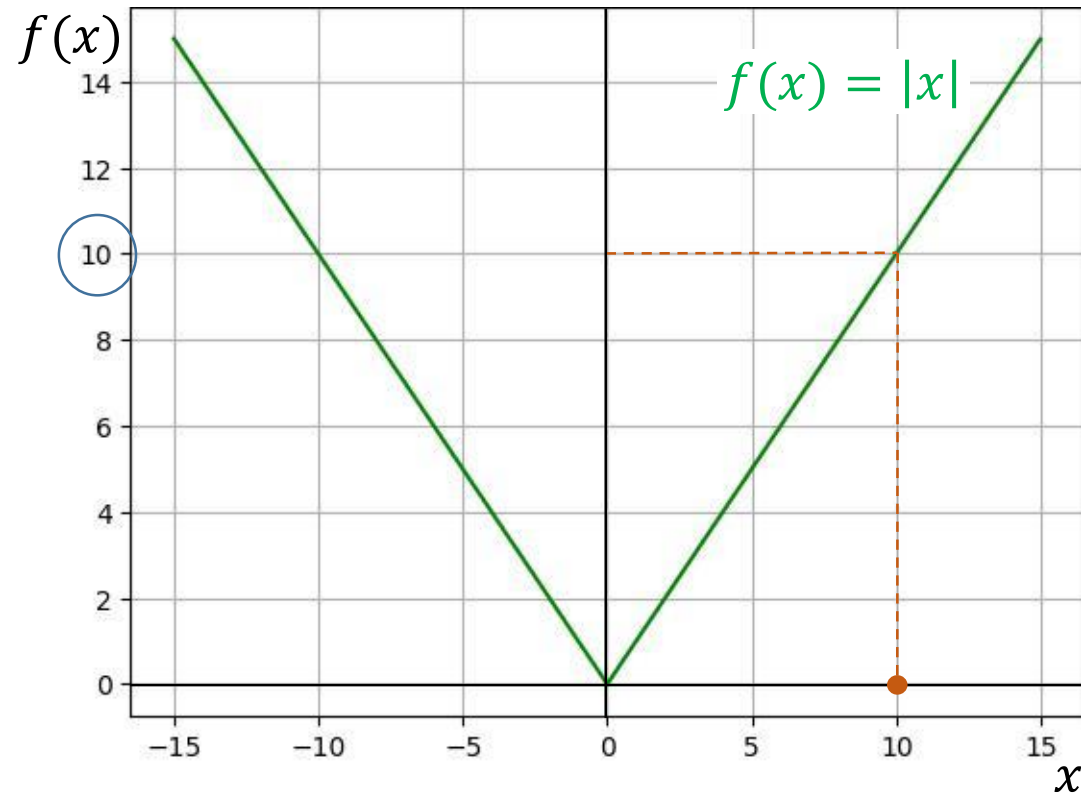
$$\frac{\partial L}{\partial w} = x \frac{(\hat{y} - y)}{|\hat{y} - y|} \quad \frac{\partial L}{\partial b} = \frac{(\hat{y} - y)}{|\hat{y} - y|}$$

5) Update parameters

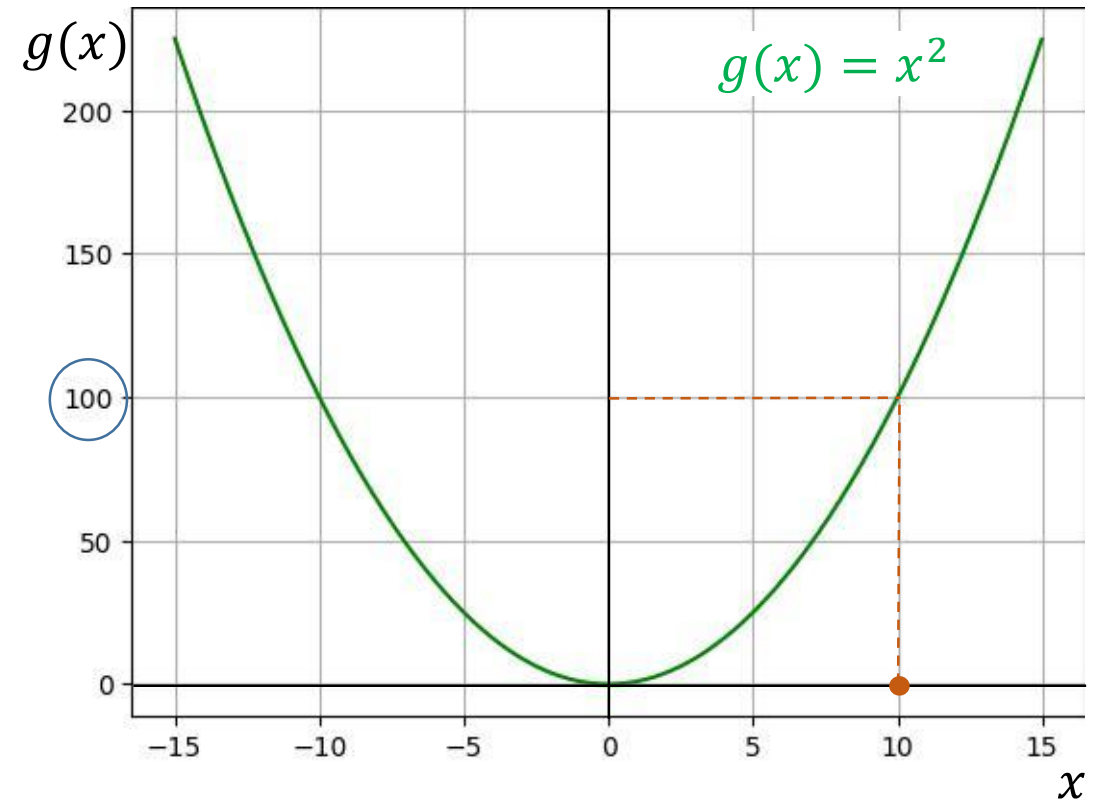
$$w = w - \eta \frac{\partial L}{\partial w} \quad b = b - \eta \frac{\partial L}{\partial b}$$



Quiz 4: The pros and cons of MSE and MAE when data contain outliers?



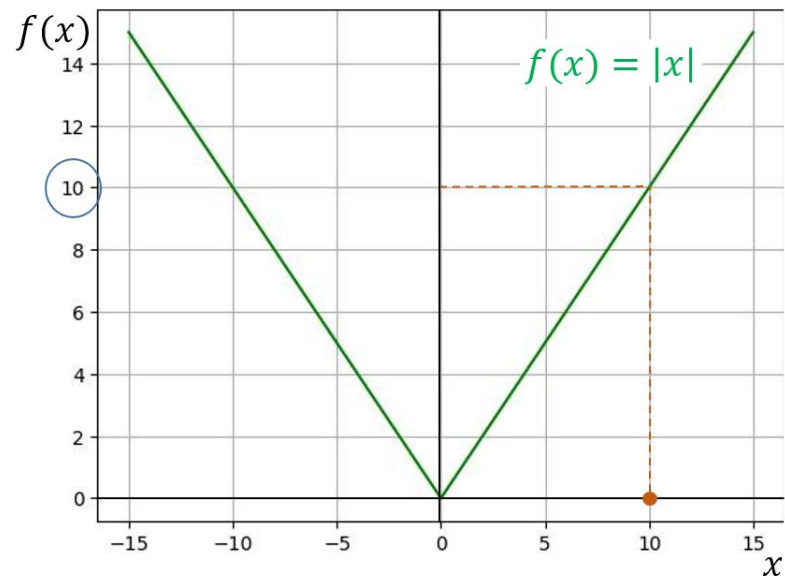
$$g(10) \gg f(10)$$



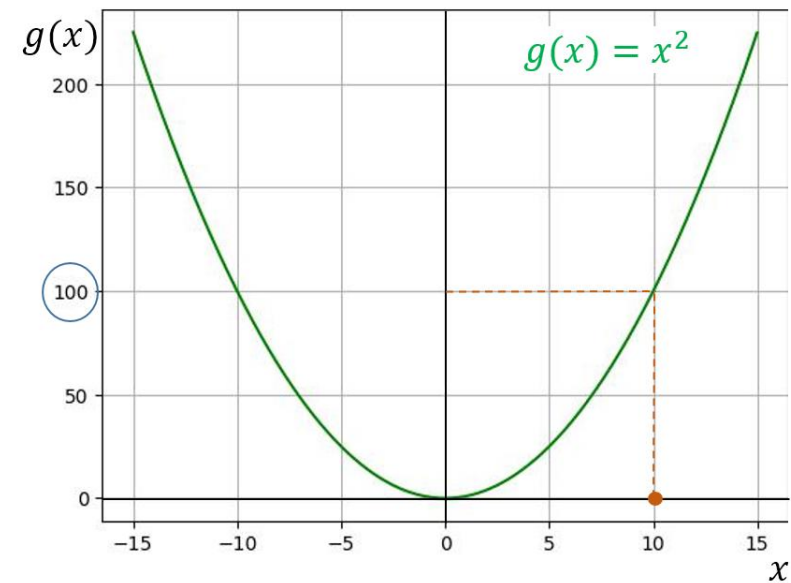
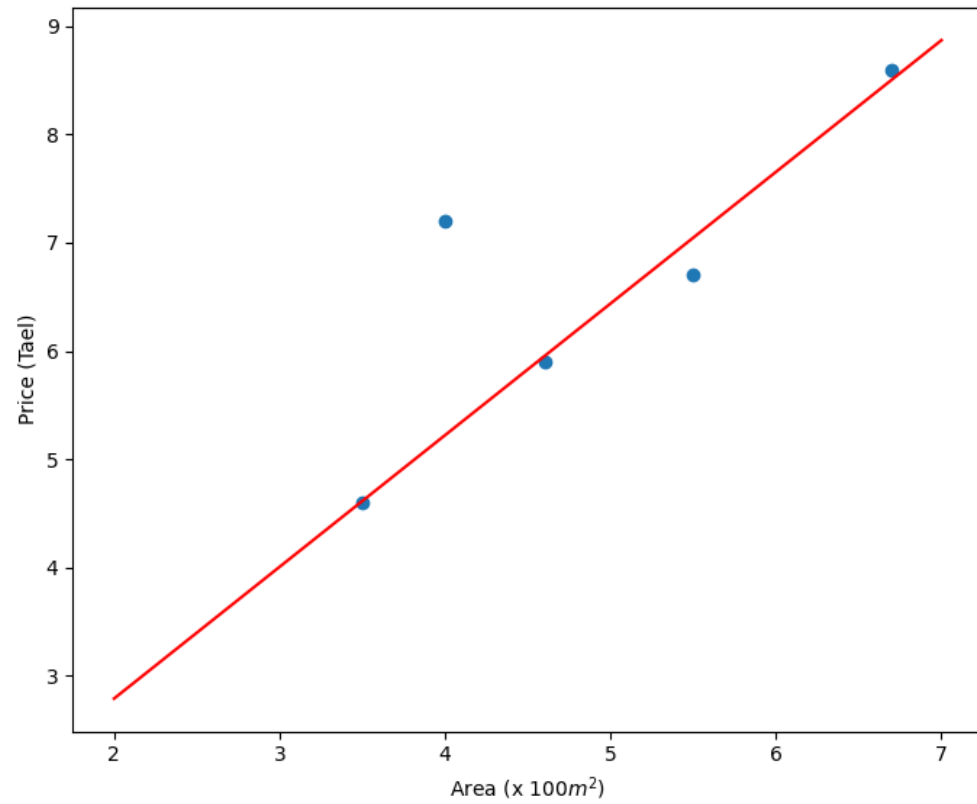
If $x = 10$ is an outlier, $g(10)$ has more negative effect

⇒ MAE is better to tolerate outliers

MAE

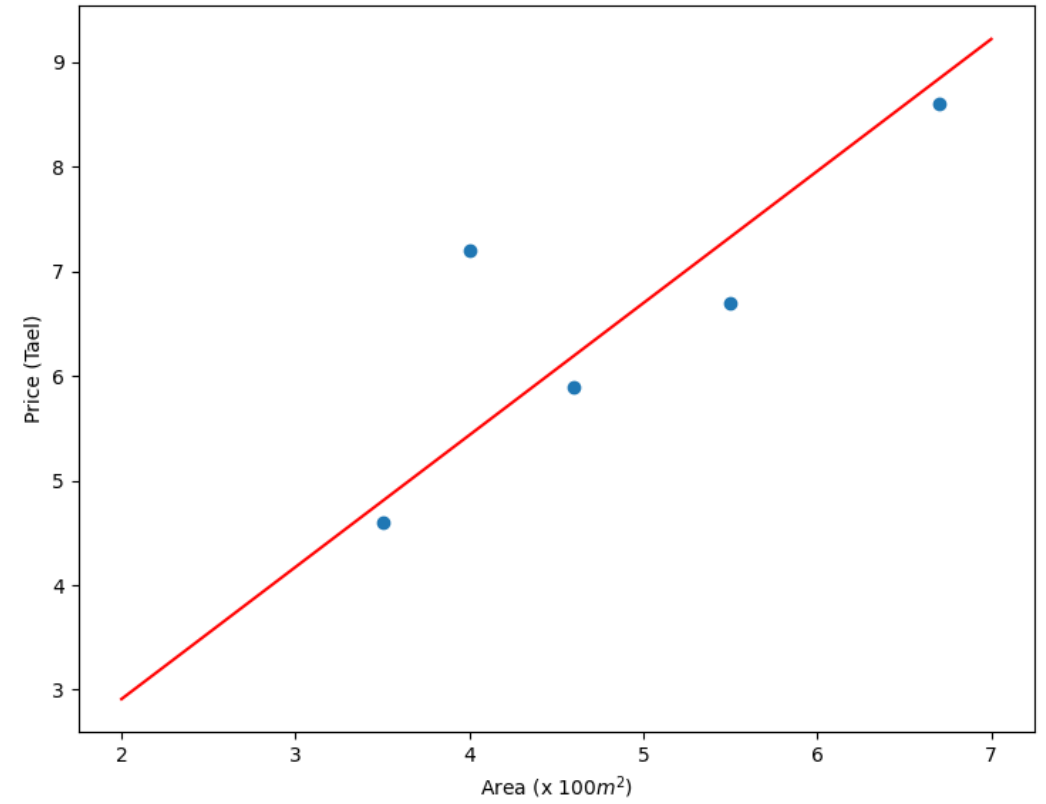


Mean Absolute Error

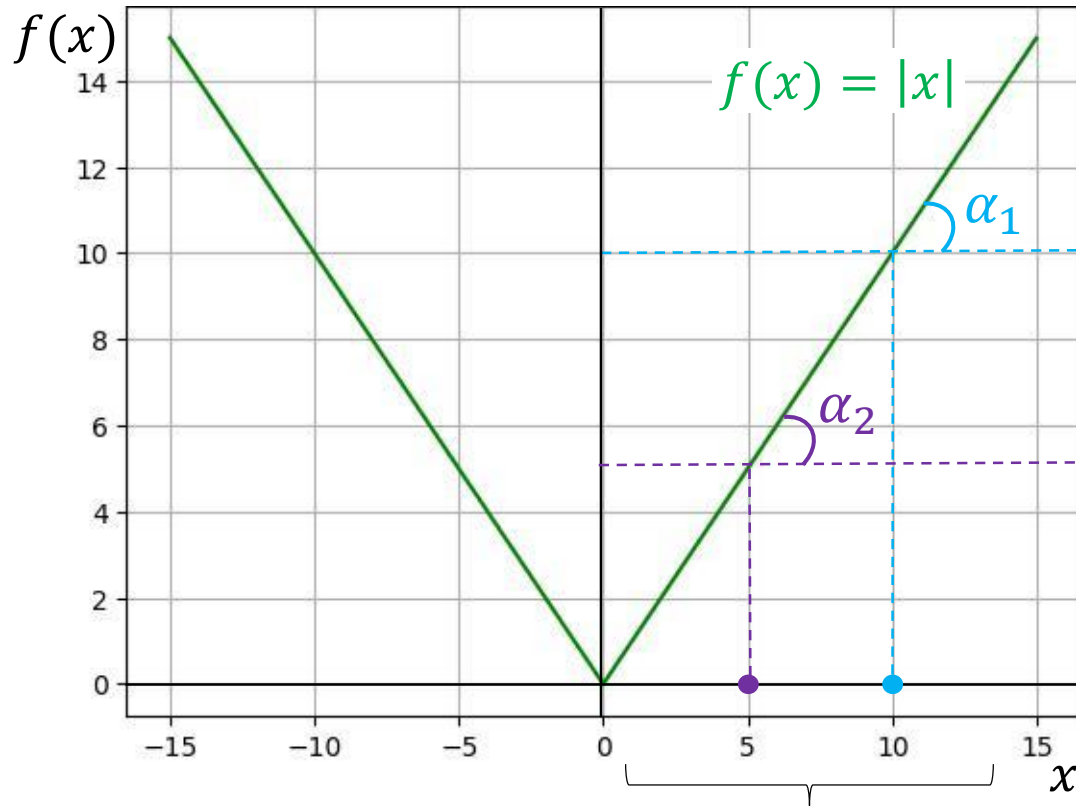


MSE

Mean Squared Error

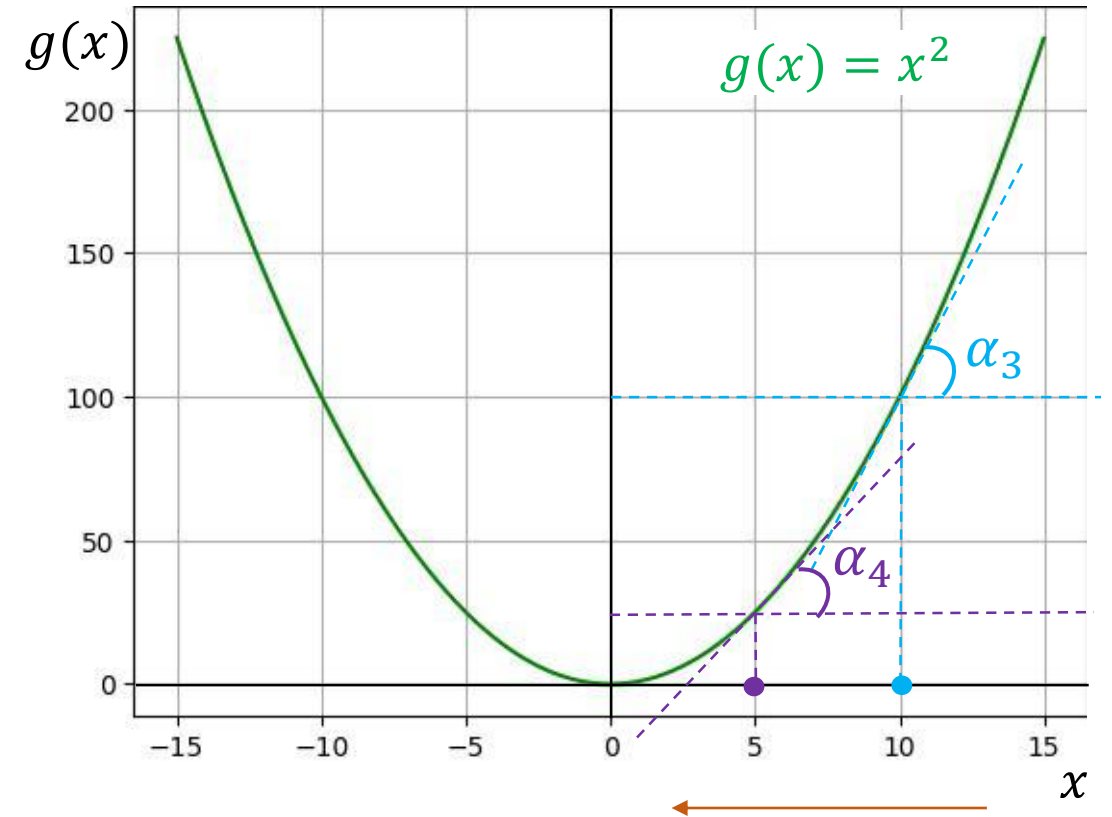


Quiz 5: The pros and cons of MSE and MAE with a fixed learning rate η ?



$$\alpha_1 = \alpha_2$$

$\eta f'(x)$ values are constants



$$\alpha_3 > \alpha_4$$

$\eta f'(x)$ values reduce

\Rightarrow MSE is better when working with a fixed learning rate

