

**AI VIETNAM**  
**All-in-One Course**

# **Insight into Logistic Regression**

**Quang-Vinh Dinh**  
**Ph.D. in Computer Science**

# Outline

- Vectorization
- Optimiztion for 1+ samples
- Logistic Regression – Mini-batch
- Logistic Regression – Batch
- BCE and MSE Loss Functions
- Sigmoid and Tanh Function (Optional)

# Implementation - One Sample

Feature	Label
Petal_Length	Label
1.4	0
1.5	0
3	1
4.1	1

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

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

$$z = wx + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = (-y\log\hat{y} - (1-y)\log(1-\hat{y}))$$

4) Compute derivative

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

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

5) Update parameters

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

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

```
def sigmoid_function(z):
    return 1 / (1 + math.exp(-z))

def predict(x, w, b):
    z = w*x + b
    y_hat = sigmoid_function(z)

    return y_hat

def loss_function(y_hat, y):
    return -y*math.log(y_hat) - (1 - y)*math.log(1 - y_hat)

def compute_gradient(x, y_hat, y):
    dw = x*(y_hat - y)
    db = (y_hat - y)

    return dw, db

def update(w, b, dw, db, lr):
    w = w - lr*dw
    b = b - lr*db

    return w, b
```

if #features changes, which functions are affected?

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

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

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

$$z = w_1x_1 + w_2x_2 + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = (-y \log \hat{y} - (1 - y) \log (1 - \hat{y}))$$

4) Compute derivative

$$\frac{\partial L}{\partial w_i} = x_i(\hat{y} - y) \quad \frac{\partial L}{\partial b} = (\hat{y} - y)$$

5) Update parameters

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

```
def sigmoid_function(z):
    return 1 / (1 + np.exp(-z))
```

```
def predict(x1, x2, b, w1, w2):
    z = x1*w1 + x2*w2 + b
    y_hat = sigmoid_function(z)

    return y_hat
```

```
def loss_function(y_hat, y):
    return -y*np.log(y_hat) - (1 - y)*np.log(1 - y_hat)
```

```
def compute_gradient(x1, x2, y_hat, y):
    db = (y_hat - y)
    dw1 = x1*(y_hat - y)
    dw2 = x2*(y_hat - y)

    return (db, dw1, dw2)
```

How to solve the problem?

```
def update(b, w1, w2, lr, db, dw1, dw2):
    b = b - lr*db
    w1 = w1 - lr*dw1
    w2 = w2 - lr*dw2

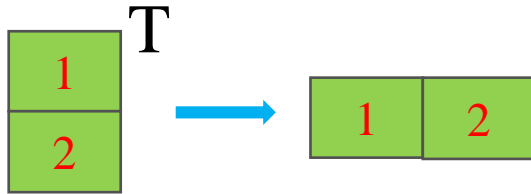
    return (b, w1, w2)
```

# Vector/Matrix Operations

## Transpose

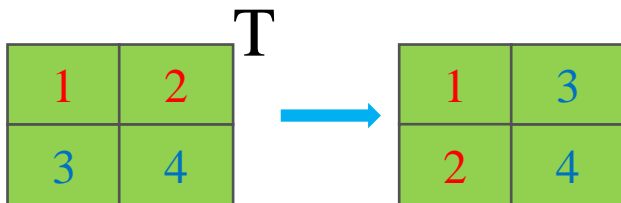
$$\vec{v} = \begin{bmatrix} v_1 \\ \dots \\ v_n \end{bmatrix}$$

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



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

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

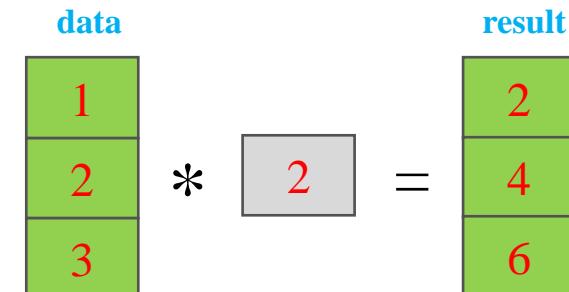


```
2 import numpy as np
3
4 # create data
5 data = np.array([1,2,3])
6 factor = 2
7
8 # broadcasting
9 result_multiplication = data*factor
```

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

## Multiply with a number

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

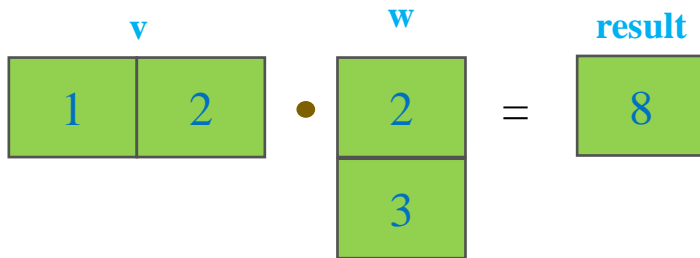


# Vector/Matrix Operations

## Dot product

$$\vec{v} = \begin{bmatrix} v_1 \\ \dots \\ v_n \end{bmatrix} \quad \vec{u} = \begin{bmatrix} u_1 \\ \dots \\ u_n \end{bmatrix}$$

$$\vec{v} \cdot \vec{u} = v_1 \times u_1 + \dots + v_n \times u_n$$



```
1 def dot_product(vector1, vector2):
2     '''
3     Compute dot product between two vectors
4     Output is a floating-point number
5     '''
6
7     return sum([v1*v2 for v1, v2 in zip(vector1, vector2)])
8
9 # test case
10 vector1 = [1, 2, 3]
11 vector2 = [2, 3, 4]
12
13 ouptut = dot_product(vector1, vector2)
14 print(ouptut)
```

20

```
2 import numpy as np
3
4 v = np.array([1, 2])
5 w = np.array([2, 3])
6
7 # Tính inner product giữa v và w
8 print('method 1 \n', v.dot(w))
9 print('method 2 \n', np.dot(v, w))
```

method 1  
8

method 2  
8

# Vectorization

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  $\hat{y}$

$$z = wx + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Traditional

3) Compute loss

$$L(\hat{y}, y) = (-y \log \hat{y} - (1 - y) \log (1 - \hat{y}))$$

4) Compute derivative

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

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

5) Update parameters

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

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

$\eta$  is learning rate

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

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

$$z = wx + b1 = [b \quad w] \begin{bmatrix} 1 \\ x \end{bmatrix} = \theta^T x$$

dot product

# Vectorization

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

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

$$z = wx + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = \underline{(-y \log \hat{y} - (1 - y) \log(1 - \hat{y}))}$$

4) Compute derivative

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

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

5) Update parameters

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

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

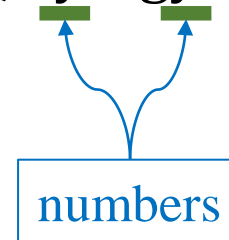
$\eta$  is learning rate

Traditional

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

$$z = \theta^T x \quad \hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$L(\hat{y}, y) = (-y \log \hat{y} - (1 - y) \log(1 - \hat{y}))$$



What will we do?



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

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

$$z = wx + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = (-y \log \hat{y} - (1-y) \log(1-\hat{y}))$$

4) Compute derivative

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

5) Update parameters

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

$$\begin{bmatrix} (\hat{y} - y) \times 1 \\ (\hat{y} - y) \times x \end{bmatrix} = \underbrace{(\hat{y} - y)}_{\text{common factor}} \begin{bmatrix} 1 \\ x \end{bmatrix} = (\hat{y} - y) \mathbf{x} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix} = \nabla_{\theta} L \quad \rightarrow \quad \nabla_{\theta} L = \mathbf{x}(\hat{y} - y)$$

Traditional

# Vectorization

$$z = \mathbf{w} \mathbf{x} + b \qquad \mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \qquad \boldsymbol{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\begin{cases} \frac{\partial L}{\partial b} = (\hat{y} - y) = (\hat{y} - y) \times 1 \\ \frac{\partial L}{\partial w} = x(\hat{y} - y) = (\hat{y} - y) \times x \end{cases}$$

# Vectorization

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

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

$$z = wx + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = (-y \log \hat{y} - (1-y) \log(1-\hat{y}))$$

4) Compute derivative

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

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

5) Update parameters

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

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

$\eta$  is learning rate

Traditional

$$z = \theta^T x$$

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

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

$$\nabla_{\theta} L = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{bmatrix}$$

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

$$\rightarrow \theta = \theta - \eta \nabla_{\theta} L$$

# Vectorization

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

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

$$z = wx + b \quad \hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = (-y \log \hat{y} - (1-y) \log(1-\hat{y}))$$

4) Compute derivative

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

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

5) Update parameters

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

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

$\eta$  is learning rate

Traditional

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

2) Compute output  $\hat{y}$

$$z = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta} \quad \hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

$$L(\hat{y}, y) = (-y \log \hat{y} - (1-y) \log(1-\hat{y}))$$

4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{y} - y)$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

$\eta$  is learning rate

Vectorized

# Vectorization

## ❖ Implementation (using Numpy)

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



2) Compute output  $\hat{y}$



$$z = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta} \quad \hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss



$$L(\hat{y}, y) = (-y \log \hat{y} - (1-y) \log (1-\hat{y}))$$

4) Compute derivative



$$\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{y} - y)$$

→ 5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

$\eta$  is learning rate

```
def sigmoid_function(z):  
    return 1 / (1 + np.exp(-z))  
  
def predict(X, theta):  
    return sigmoid_function( np.dot(X.T, theta) )  
  
def loss_function(y_hat, y):  
    return -y*np.log(y_hat) - (1 - y)*np.log(1 - y_hat)  
  
def compute_gradient(X, y_hat, y):  
    return X*(y_hat - y)  
  
def update(theta, lr, gradient):  
    return theta - lr*gradient
```

# compute output

y\_hat = predict(X, theta)

# compute loss

loss = loss\_function(y\_hat, y)

# compute mean of gradient

gradient = compute\_gradient(X, y\_hat, y)

# update

theta = update(theta, lr, gradient)

# Given X and y

Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

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

↓

2) Compute output  $\hat{y}$

↓  $z = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta} \qquad \hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$

3) Compute loss

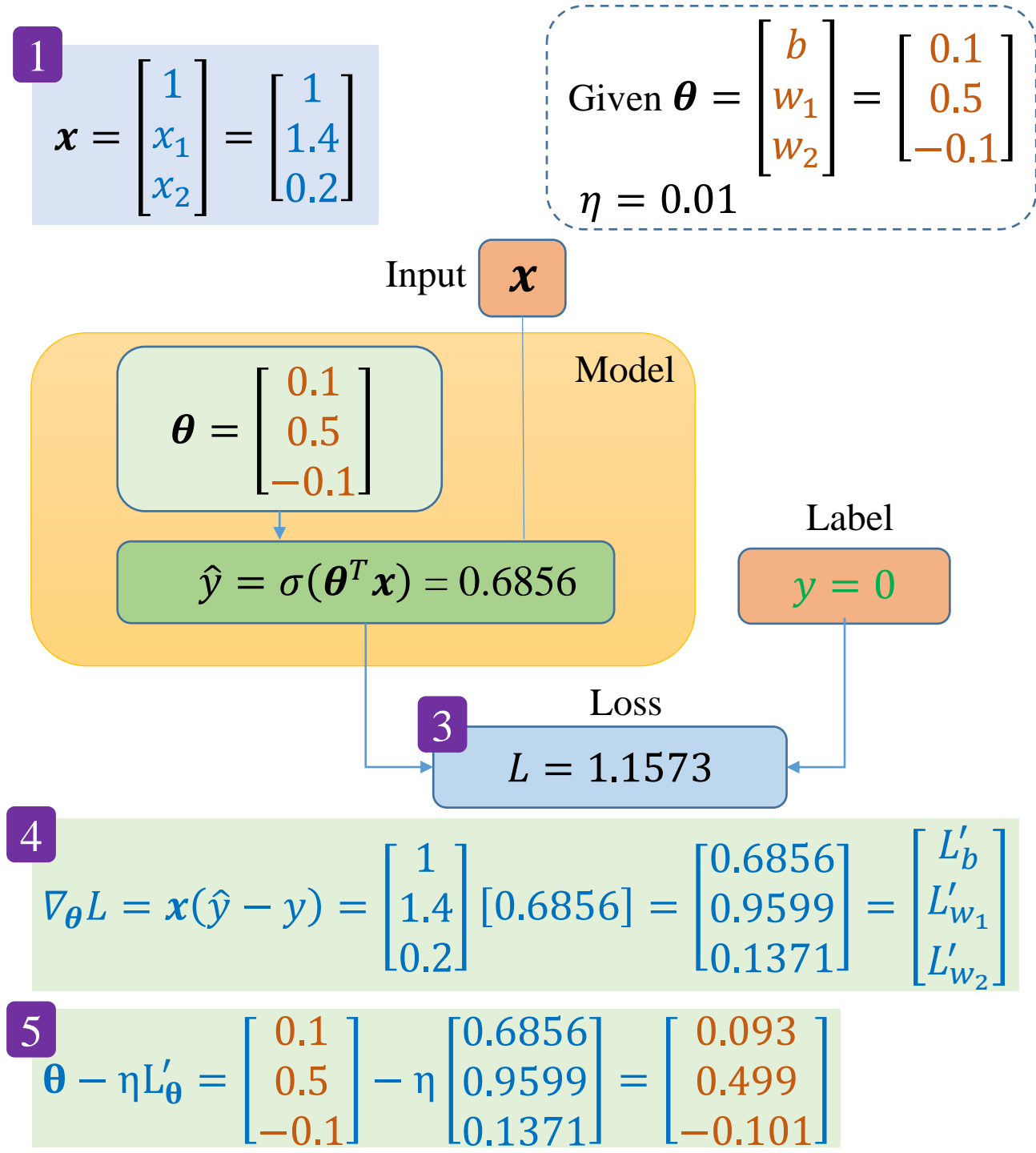
↓  $L(\hat{y}, y) = (-y \log \hat{y} - (1-y) \log(1-\hat{y}))$

4) Compute derivative

↓  $\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{y} - y)$

5) Update parameters

$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$



# Logistic Regression-Stochastic

Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1.4 \\ 0.2 \end{bmatrix} \quad \mathbf{y} = [0]$$

Demo

- 1) Pick a sample  $(\mathbf{x}, y)$  from training data
- 2) Compute output  $\hat{y}$

$$z = \boldsymbol{\theta}^T \mathbf{x}$$
$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

- 3) Compute loss

$$L(\boldsymbol{\theta}) = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

- 4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{y} - y)$$

- 5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

$\eta$  is learning rate

# Outline

- Vectorization
- Optimiztion for 1+ samples
- Logistic Regression – Mini-batch
- Logistic Regression – Batch
- BCE and MSE Loss Functions
- Sigmoid and Tanh Function (Optional)

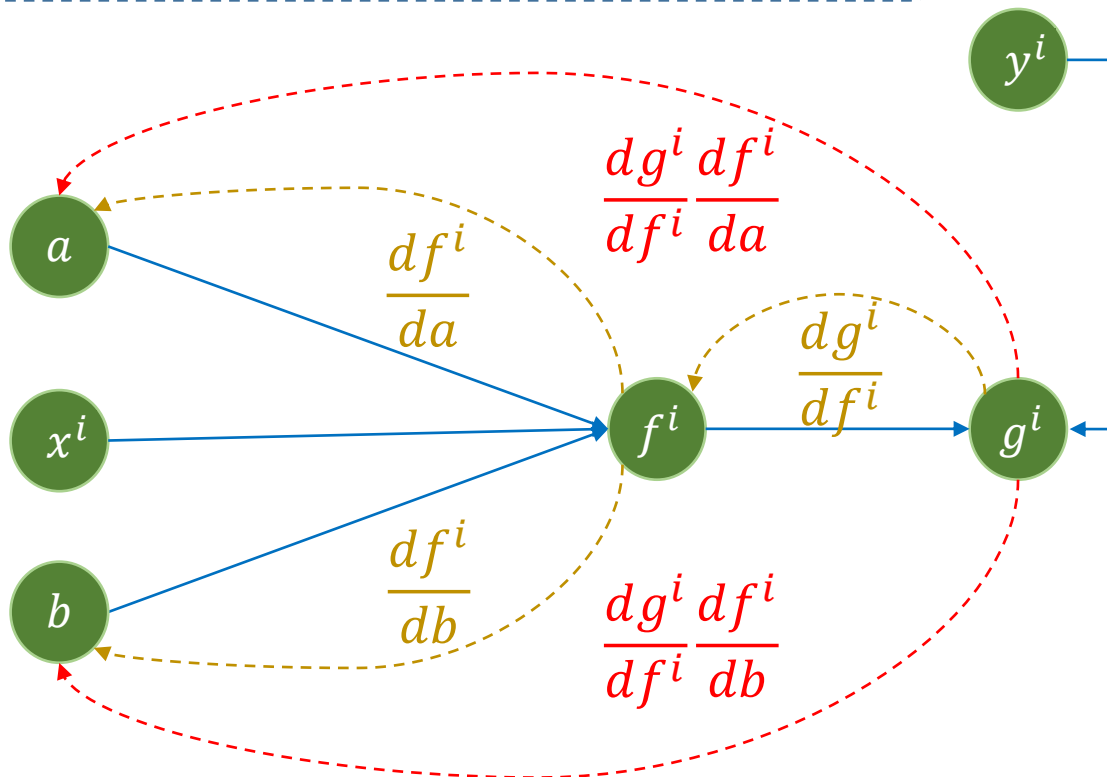
# Optimization for One+ Samples

## ❖ Equations for partial gradients

$$f(x^i) = ax^i + b \quad (x^1=1, y^1=5)$$

$$g(f^i) = (f^i - y^i)^2 \quad (x^2=2, y^2=7)$$

illustration



$$\frac{df}{da} = x$$

$$\frac{df}{db} = 1$$

$$\frac{dg}{df} = 2(f - y)$$

$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$

$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$

During looking for optimal  $a$  and  $b$ , at a given time,  $a$  and  $b$  have concrete values



## ❖ Optimization for a composite function

Find a and b so that  $g(f(x))$  is minimum

$$f(x^i) = ax^i + b \quad (x^1=1, y^1=5)$$

$$g(f^i) = (f^i - y^i)^2 \quad (x^2=2, y^2=7)$$

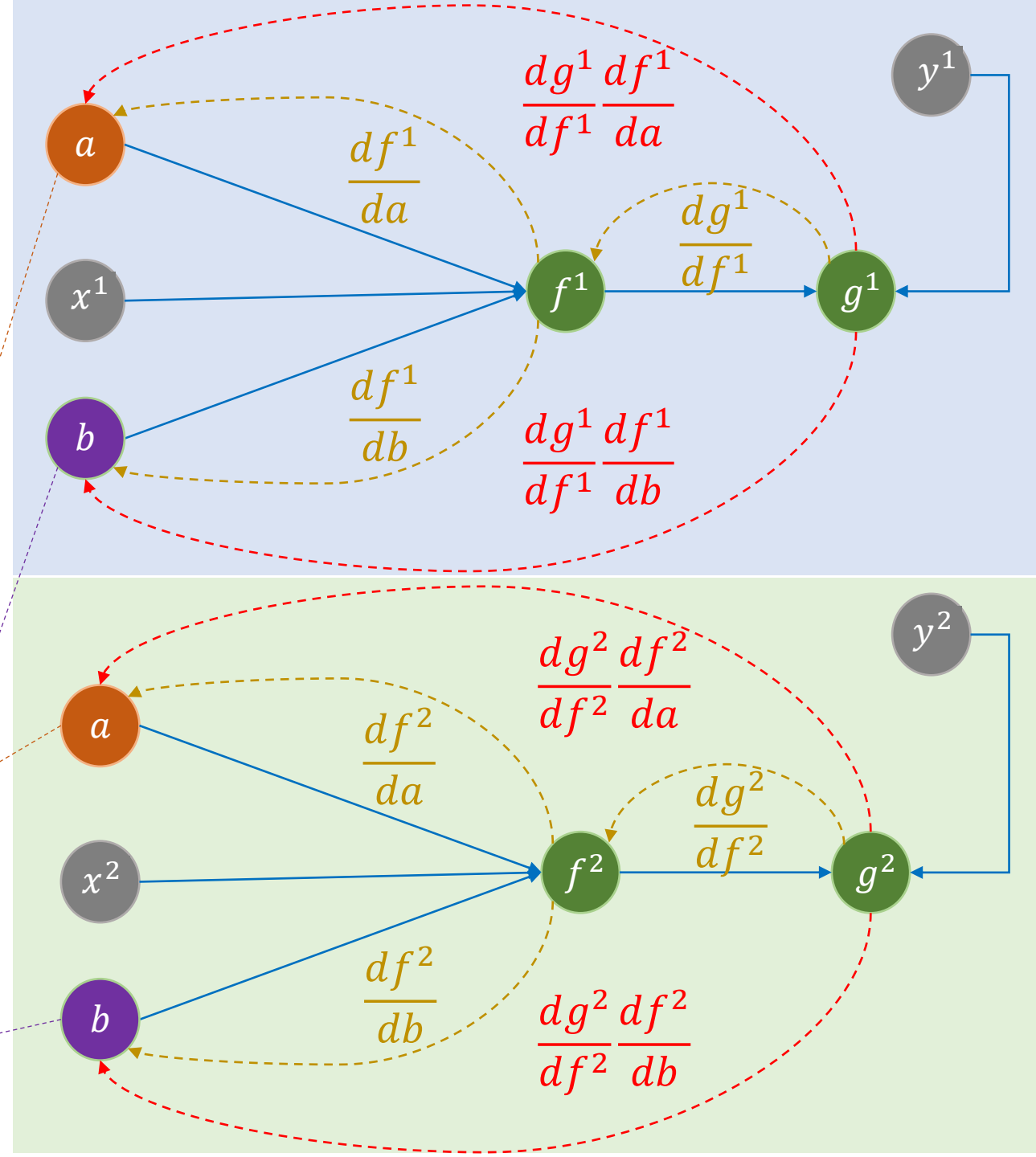
Partial derivative functions

$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$

$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$

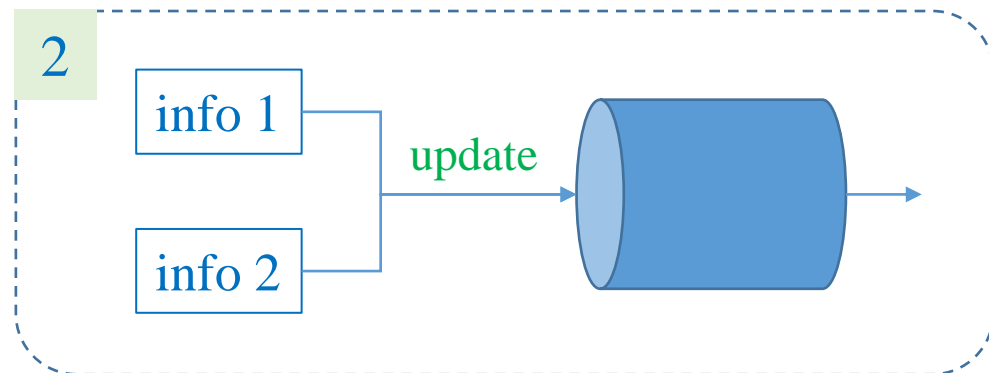
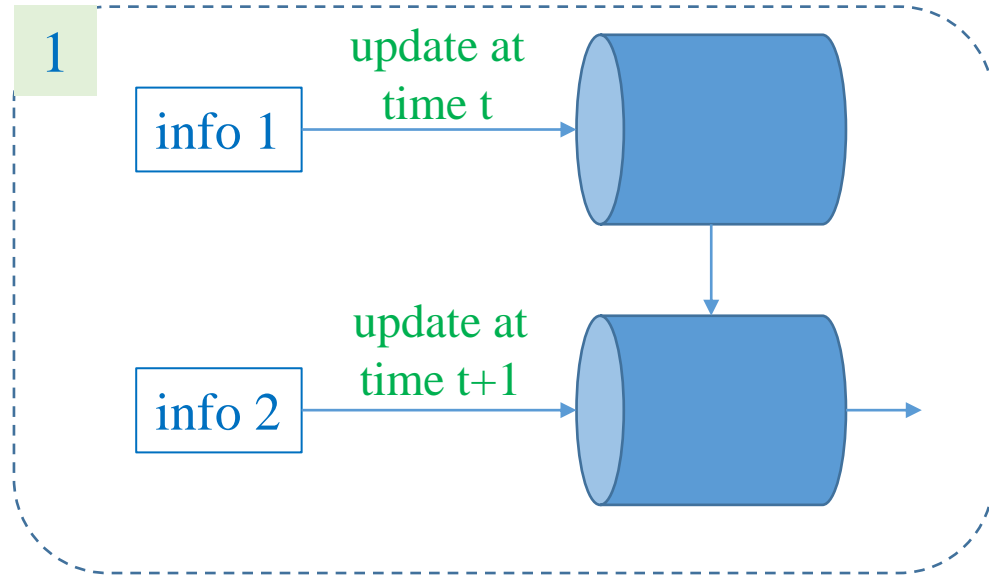
$$\sum_i \frac{dg^i}{da} = \frac{dg^1}{df^1} \frac{df^1}{da} + \frac{dg^2}{df^2} \frac{df^2}{da}$$

$$\sum_i \frac{dg_i}{db} = \frac{dg^1}{df^1} \frac{df^1}{db} + \frac{dg^2}{df^2} \frac{df^2}{db}$$

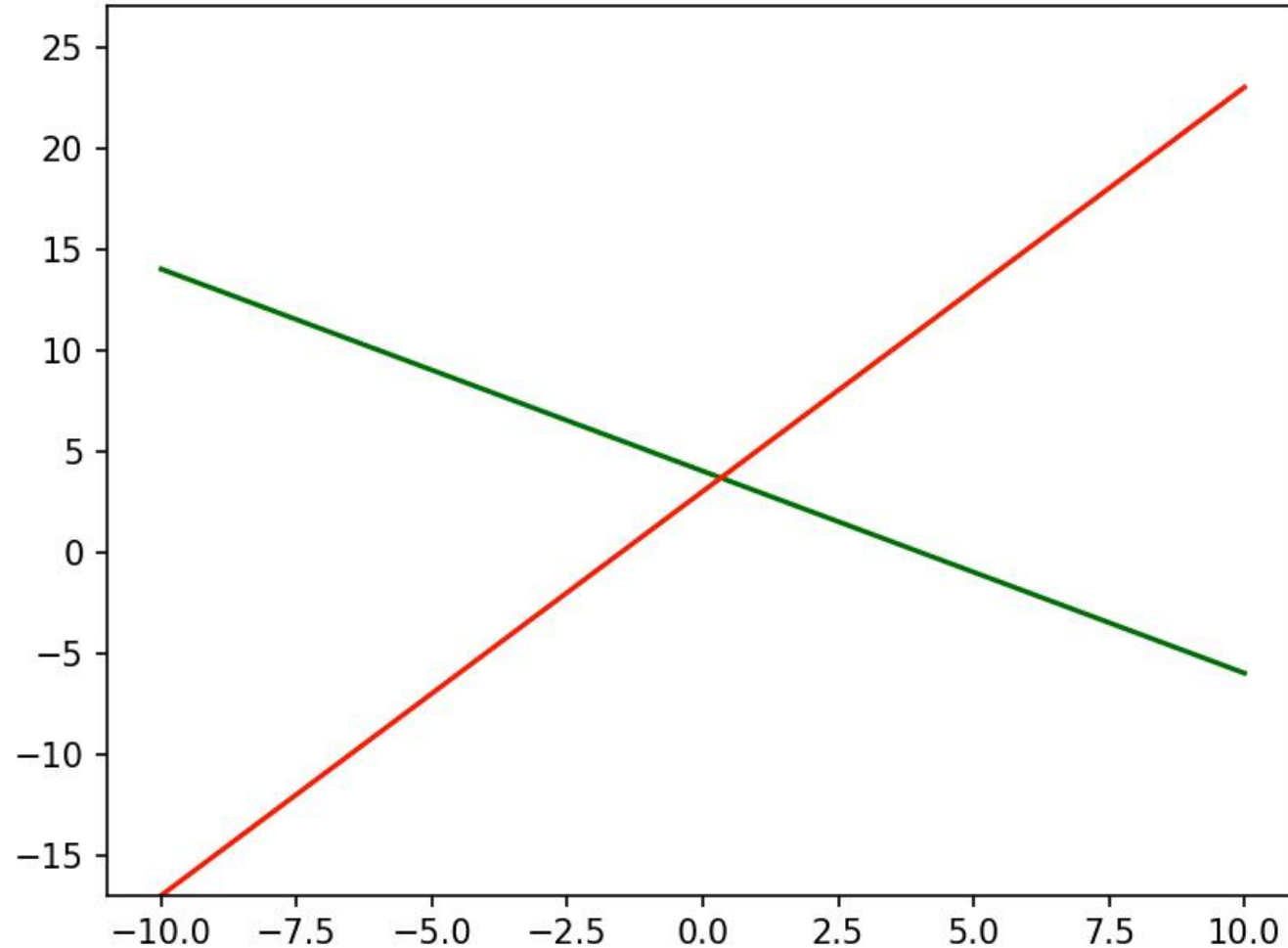
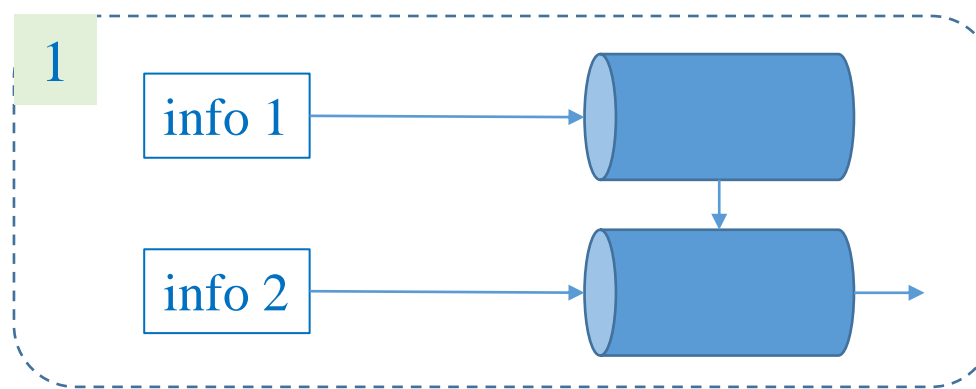


# Optimization

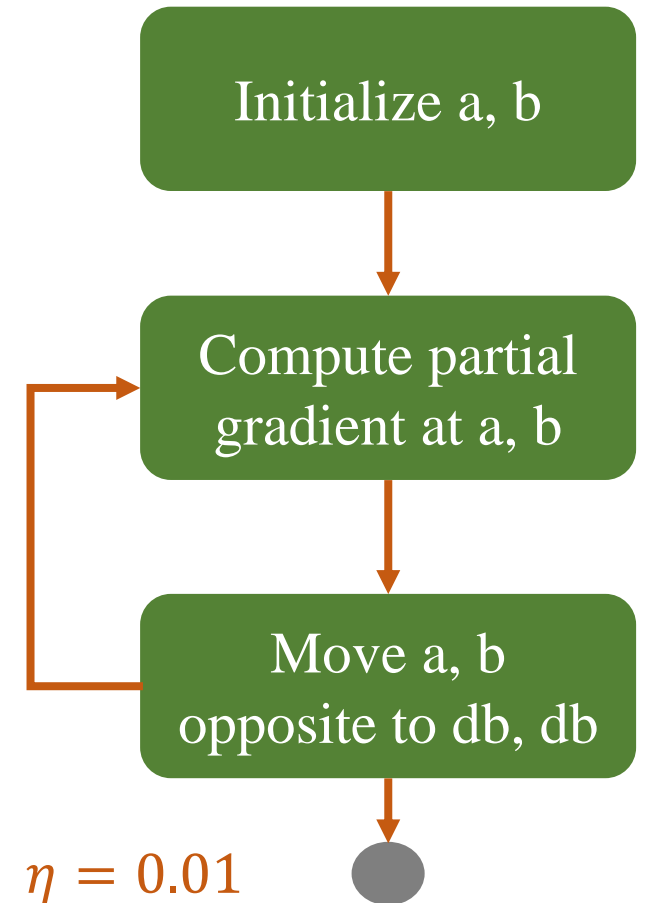
## ❖ How to use gradient information



# Summary

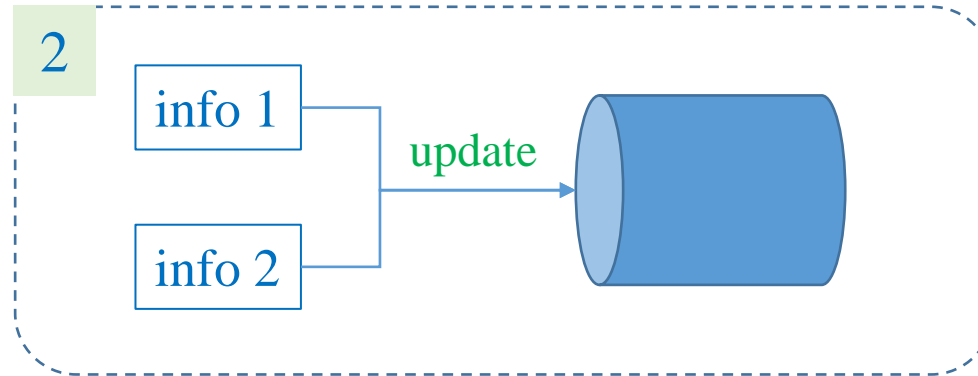


$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$

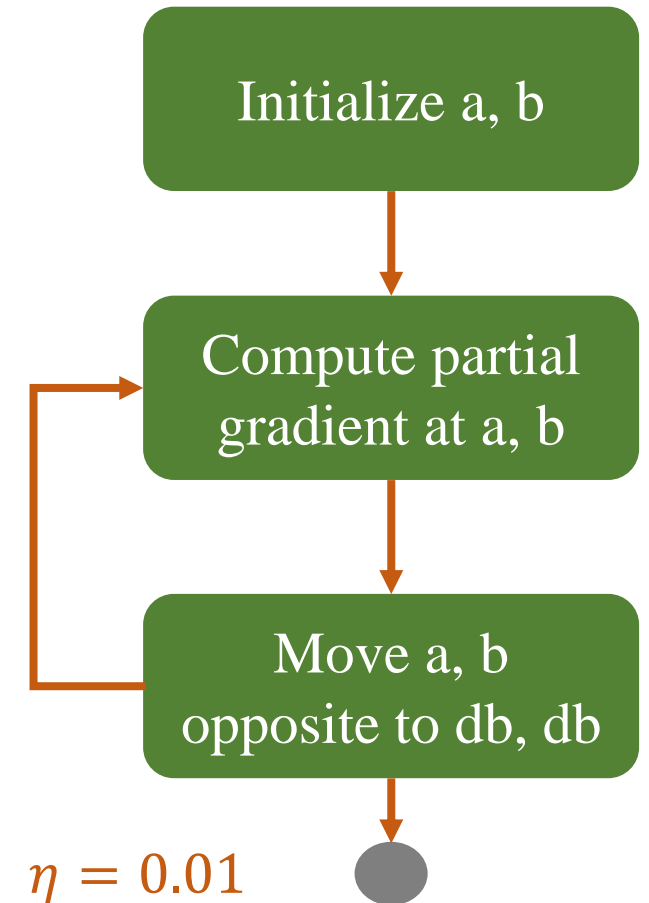
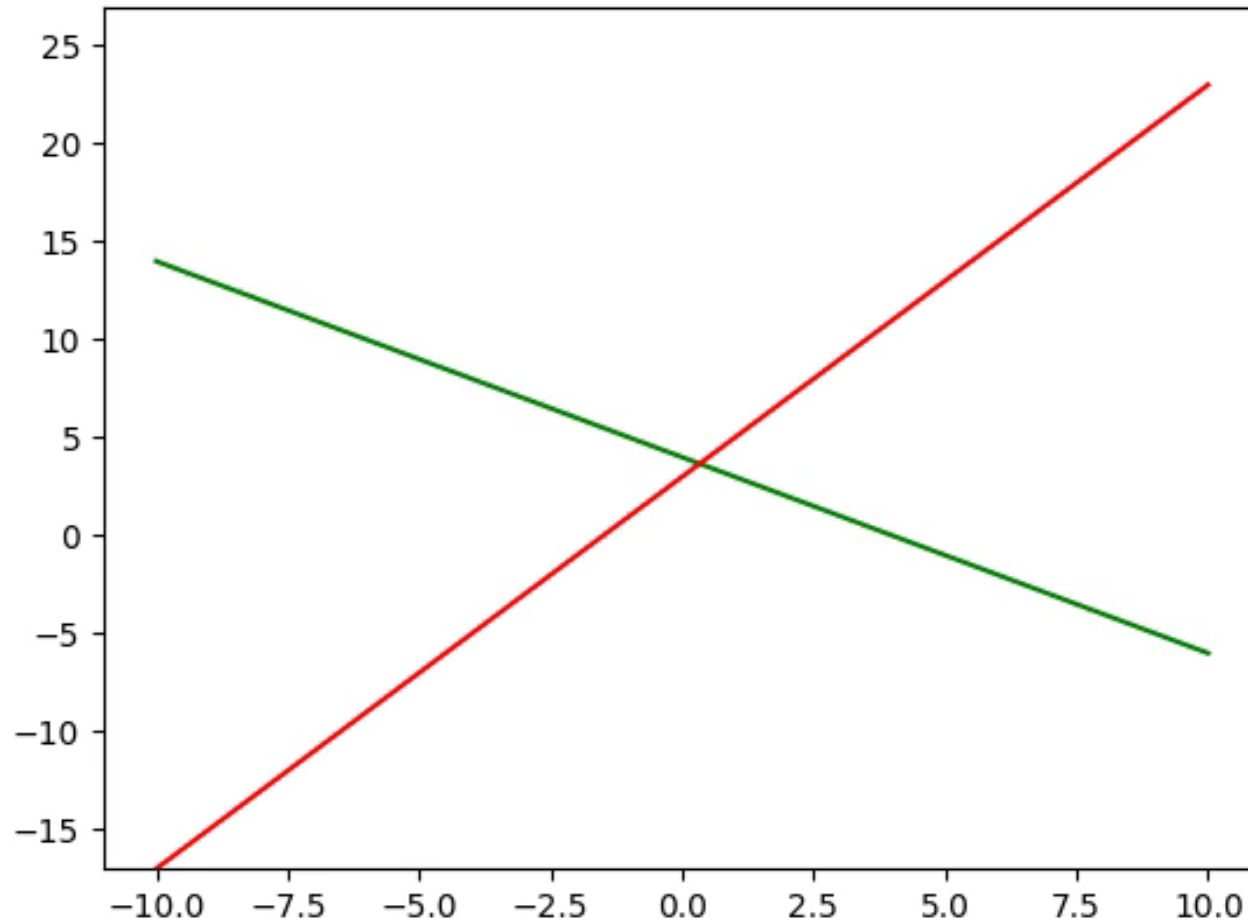


# Summary

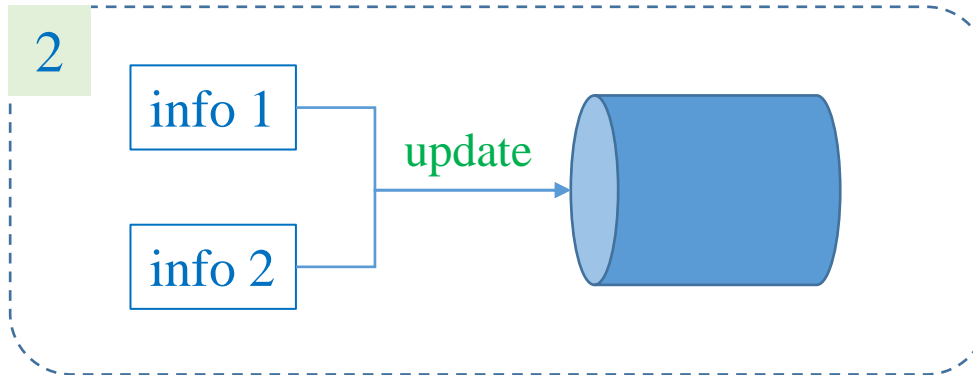
2



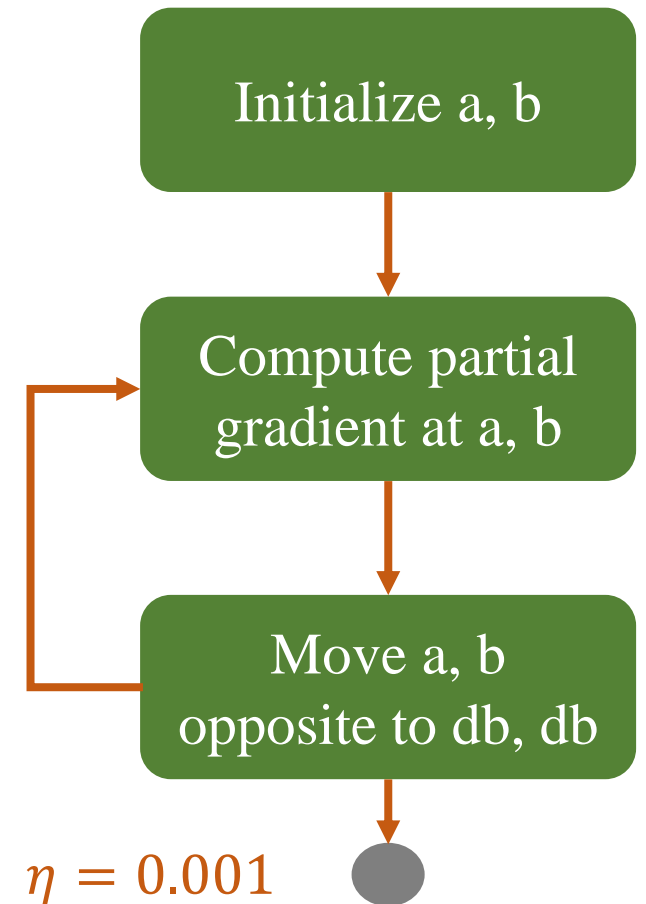
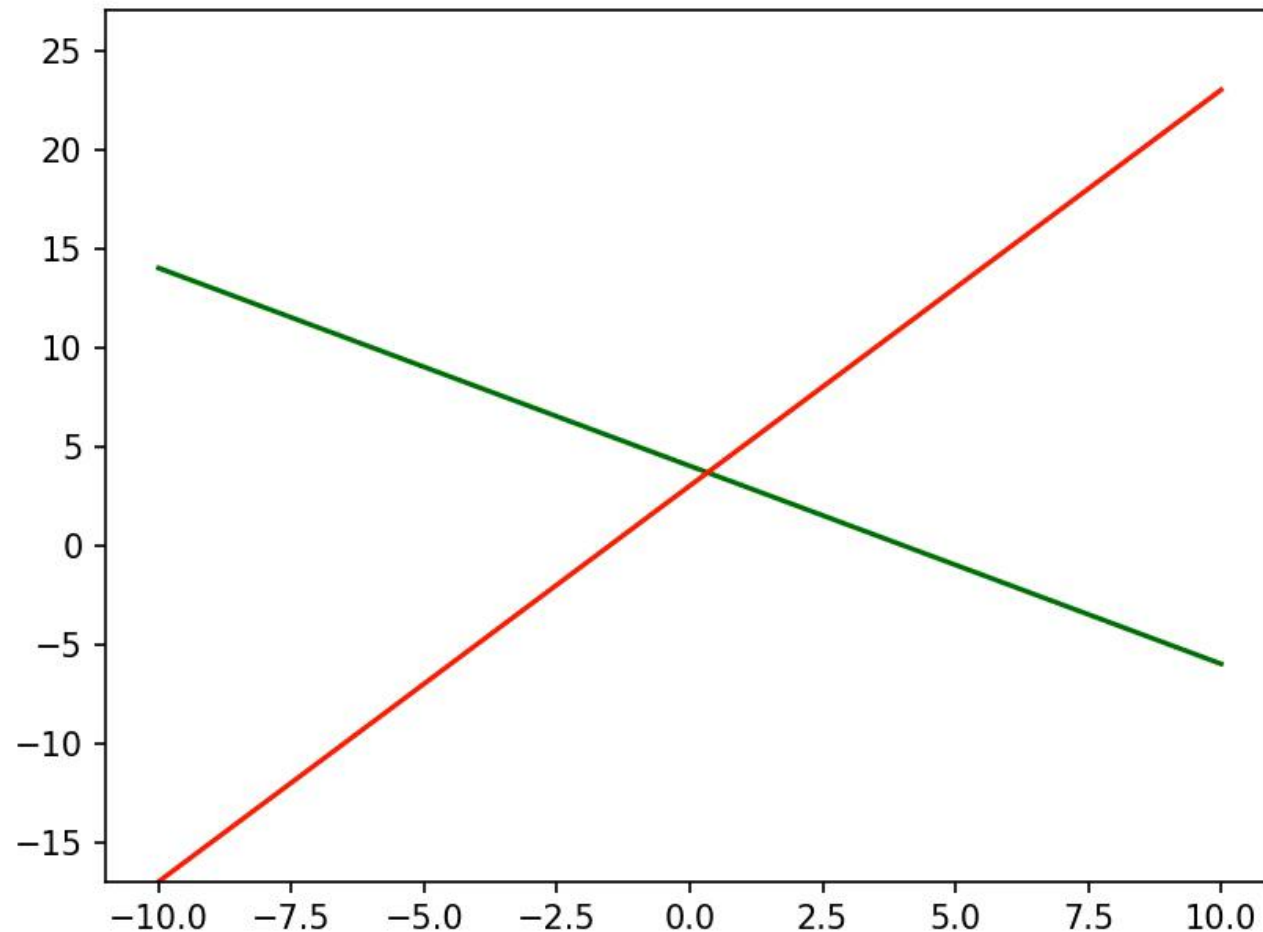
$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$



# Summary



$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$



# Outline

- Vectorization
- Optimiztion for 1+ samples
- Logistic Regression – Mini-batch
- Logistic Regression – Batch
- BCE and MSE Loss Functions
- Sigmoid and Tanh Function (Optional)

# Linear Regression (m-samples)

## ❖ Construct formulas

Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

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

$$\mathbf{x} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

$$\mathbf{z} = \begin{bmatrix} z^{(1)} \\ z^{(2)} \end{bmatrix} = \begin{bmatrix} w_1 x_1^{(1)} + w_2 x_2^{(1)} + b \\ w_1 x_1^{(2)} + w_2 x_2^{(2)} + b \end{bmatrix}$$

$$= \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \mathbf{x}\boldsymbol{\theta} = \begin{bmatrix} 0.83 \\ 2.02 \end{bmatrix}$$

$$\mathbf{z} = \boldsymbol{\theta}^T \mathbf{x}$$

$$\hat{y} = \sigma(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

# Linear Regression (m-samples)

## ❖ Construct formulas

Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

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

$$\mathbf{x} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

$$\mathbf{z} = \mathbf{x}\boldsymbol{\theta} = \begin{bmatrix} 0.83 \\ 2.02 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \sigma(\mathbf{z}) = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \end{bmatrix} = \begin{bmatrix} \frac{1}{1 + e^{-z^{(1)}}} \\ \frac{1}{1 + e^{-z^{(2)}}} \end{bmatrix} = \frac{1}{1 + e^{-z}} = \begin{bmatrix} 0.69 \\ 0.88 \end{bmatrix}$$

$$\mathbf{z} = \boldsymbol{\theta}^T \mathbf{x}$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Numpy perspective



# Linear Regression (m-samples)

## ❖ Construct formulas

Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

### 3) Compute loss

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

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \frac{L^{(1)}(\hat{y}^{(1)}, y^{(1)}) + L^{(2)}(\hat{y}^{(2)}, y^{(2)})}{m}$$

$$L^{(1)}(\hat{y}^{(1)}, y^{(1)}) = -y^{(1)} \log \hat{y}^{(1)} - (1-y^{(1)}) \log(1-\hat{y}^{(1)})$$

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

+

$$\mathbf{y}^T \log \hat{\mathbf{y}}$$

$$(1-\mathbf{y})^T \log(1-\hat{\mathbf{y}})$$

## 4) Compute derivative

sample 1

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

$$\frac{\partial L^{(1)}}{\partial w_1} = x_1^{(1)} (\hat{y}^{(1)} - y^{(1)})$$

$$\frac{\partial L^{(1)}}{\partial w_2} = x_2^{(1)} (\hat{y}^{(1)} - y^{(1)})$$

sample 2

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

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

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

$$\begin{aligned} \frac{\partial L}{\partial b} &= \frac{\partial L^{(1)}}{\partial b} + \frac{\partial L^{(2)}}{\partial b} = \frac{(\hat{y}^{(1)} - y^{(1)}) + (\hat{y}^{(2)} - y^{(2)})}{m} \\ &= \frac{1}{m} [1 * (\hat{y}^{(1)} - y^{(1)}) + 1 * (\hat{y}^{(2)} - y^{(2)})] \\ &= \frac{1}{m} \begin{bmatrix} x_0^{(1)} & x_0^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix} \end{aligned}$$

$$x_0^{(1)} = 1$$

$$x_0^{(2)} = 1$$

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= \frac{\partial L^{(1)}}{\partial w_1} + \frac{\partial L^{(2)}}{\partial w_1} = \frac{x_1^{(1)} (\hat{y}^{(1)} - y^{(1)}) + x_1^{(2)} (\hat{y}^{(2)} - y^{(2)})}{m} \\ &= \frac{1}{m} \begin{bmatrix} x_1^{(1)} & x_1^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial w_2} &= \frac{\partial L^{(1)}}{\partial w_2} + \frac{\partial L^{(2)}}{\partial w_2} = \frac{x_2^{(1)} (\hat{y}^{(1)} - y^{(1)}) + x_2^{(2)} (\hat{y}^{(2)} - y^{(2)})}{m} \\ &= \frac{1}{m} \begin{bmatrix} x_2^{(1)} & x_2^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix} \end{aligned}$$

## 4) Compute derivative

sample 1

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

$$\frac{\partial L^{(1)}}{\partial w_1} = x_1^{(1)} (\hat{y}^{(1)} - y^{(1)})$$

$$\frac{\partial L^{(1)}}{\partial w_2} = x_2^{(1)} (\hat{y}^{(1)} - y^{(1)})$$

sample 2

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

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

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

$$\frac{\partial L}{\partial b} = \frac{1}{m} \begin{bmatrix} x_0^{(1)} & x_0^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_1} = \frac{1}{m} \begin{bmatrix} x_1^{(1)} & x_1^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_2} = \frac{1}{m} \begin{bmatrix} x_2^{(1)} & x_2^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix}$$

$$\nabla_{\theta} L = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = \frac{1}{m} \begin{bmatrix} x_0^{(1)} & x_0^{(2)} \\ x_1^{(1)} & x_1^{(2)} \\ x_2^{(1)} & x_2^{(2)} \end{bmatrix} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \end{bmatrix}$$

$(\mathbf{x}^{(1)})^T$        $(\mathbf{x}^{(2)})^T$        $\hat{\mathbf{y}}$        $\mathbf{y}$

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

## 5) Update parameters

**Dataset**

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$x = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\theta = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

$$\nabla_{\theta} L = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = \frac{1}{m} x^T (\hat{y} - y)$$

$$\begin{array}{l} b \\ w_1 \\ w_2 \\ \theta \end{array} = \begin{array}{l} b \\ w_1 \\ w_2 \\ \theta \end{array} - \eta \begin{array}{l} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \nabla_{\theta} L \end{array}$$



$$\theta = \theta - \eta \nabla_{\theta} L$$

# Logistic Regression - Minibatch

1) Pick  $m$  samples from training data

2) Compute output  $\hat{y}$

$$\mathbf{z} = \mathbf{x}\boldsymbol{\theta}$$

$$\hat{y} = \sigma(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

3) Compute loss

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

4) Compute derivative

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

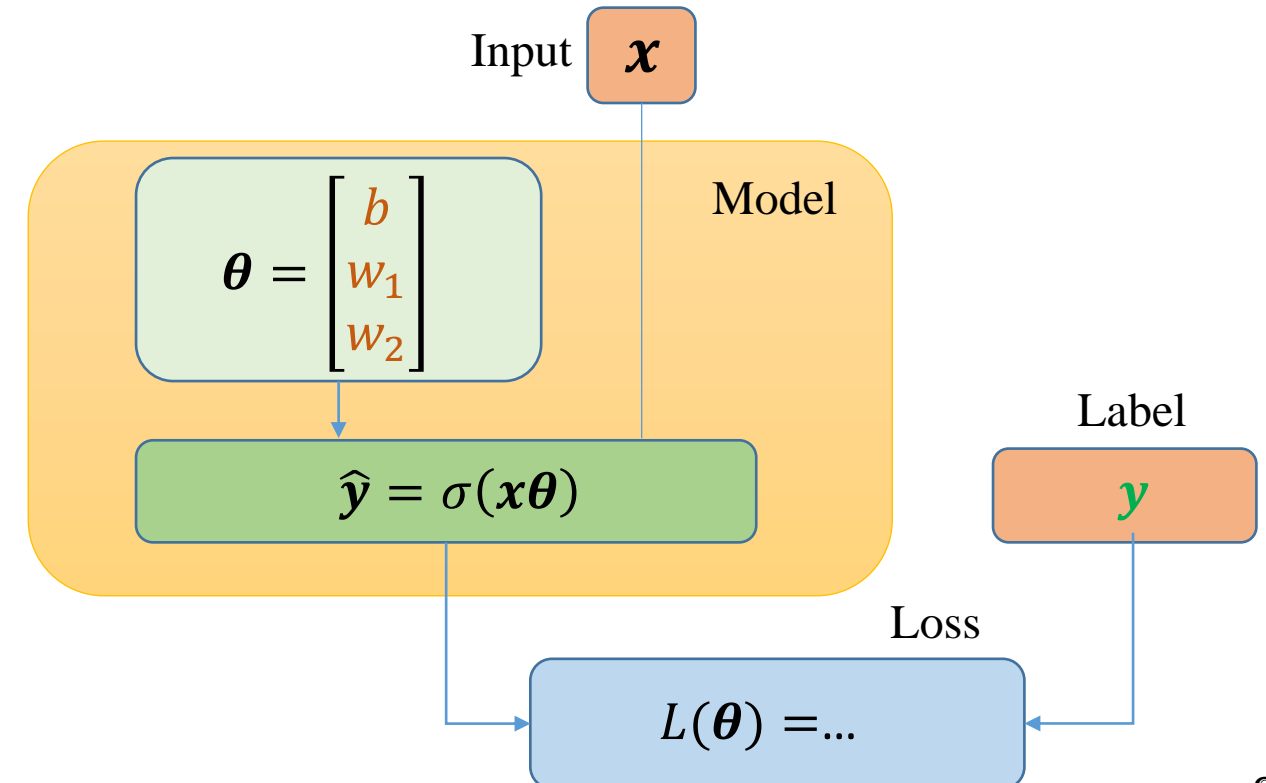
5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

$\eta$  is learning rate

Mini-batch  $m=2$

$$\mathbf{x} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$



## Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

1) Pick m samples from training data

2) Compute output  $\hat{y}$

Mini-batch m=2

$$z = x\theta$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

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

4) Compute derivative

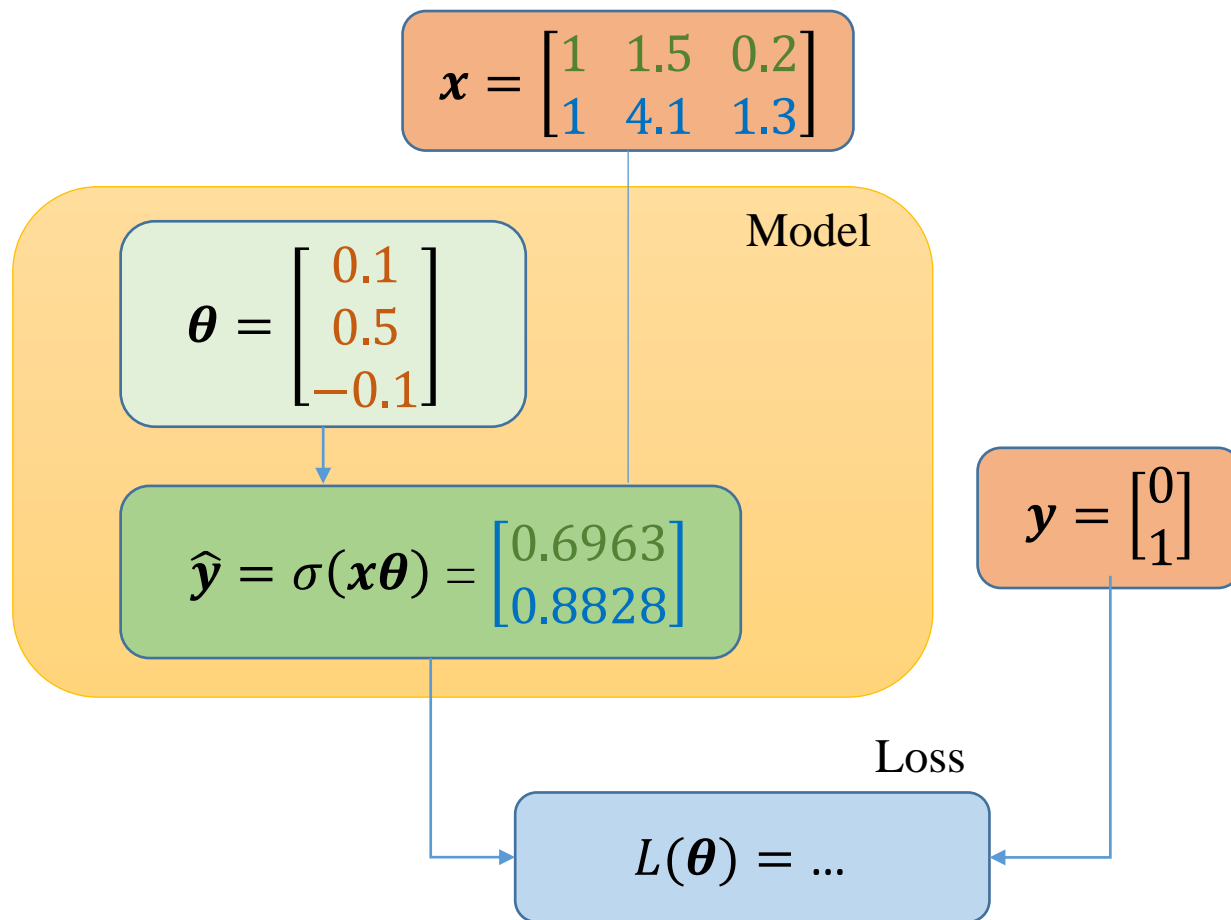
$$\nabla_{\theta} L = \frac{1}{m} x^T (\hat{y} - y)$$

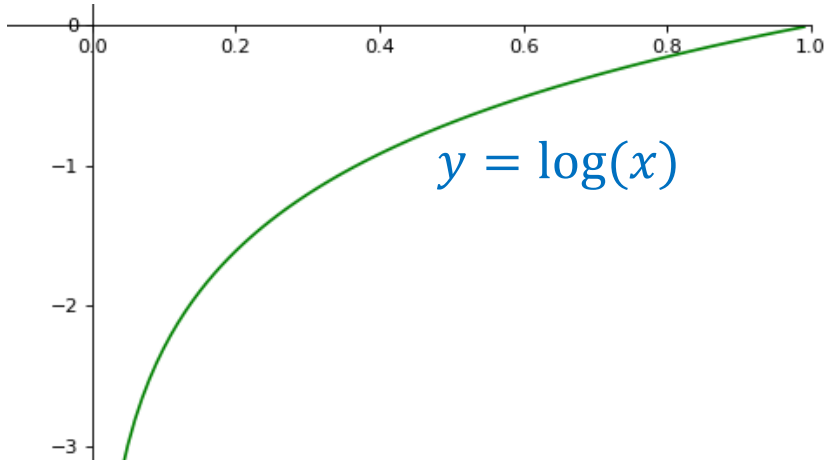
5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$

2

$$\begin{aligned} \hat{y} = \sigma(x\theta) &= \sigma \left( \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix} \right) \\ &= \sigma \left( \begin{bmatrix} 0.83 \\ 2.02 \end{bmatrix} \right) = \begin{bmatrix} 0.6963 \\ 0.8828 \end{bmatrix} \end{aligned}$$





3

$$\begin{aligned}
 L(\theta) &= \frac{1}{m} \left( -[0 \quad 1] \begin{bmatrix} \log 0.6963 \\ \log 0.8828 \end{bmatrix} - [1 \quad 0] \begin{bmatrix} \log(1 - 0.6963) \\ \log(1 - 0.8828) \end{bmatrix} \right) \\
 &= \frac{1}{m} (-\log 0.8828 - \log(1 - 0.6963)) \\
 &= \frac{0.1246 + 1.1917}{m} = 0.65815
 \end{aligned}$$

1) Pick  $m$  samples from training data

2) Compute output  $\hat{y}$

Mini-batch  $m=2$

$$z = x\theta$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

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

4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$

$$x = \begin{bmatrix} 1 & 1.5 & 0.2 \\ 1 & 4.1 & 1.3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix}$$

Model

$$\hat{y} = \sigma(x\theta) = \begin{bmatrix} 0.6963 \\ 0.8828 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Loss

$$L(\theta) = 0.65815$$

## Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

1) Pick m samples from training data

2) Compute output  $\hat{y}$

Mini-batch m=2

$$z = x\theta$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

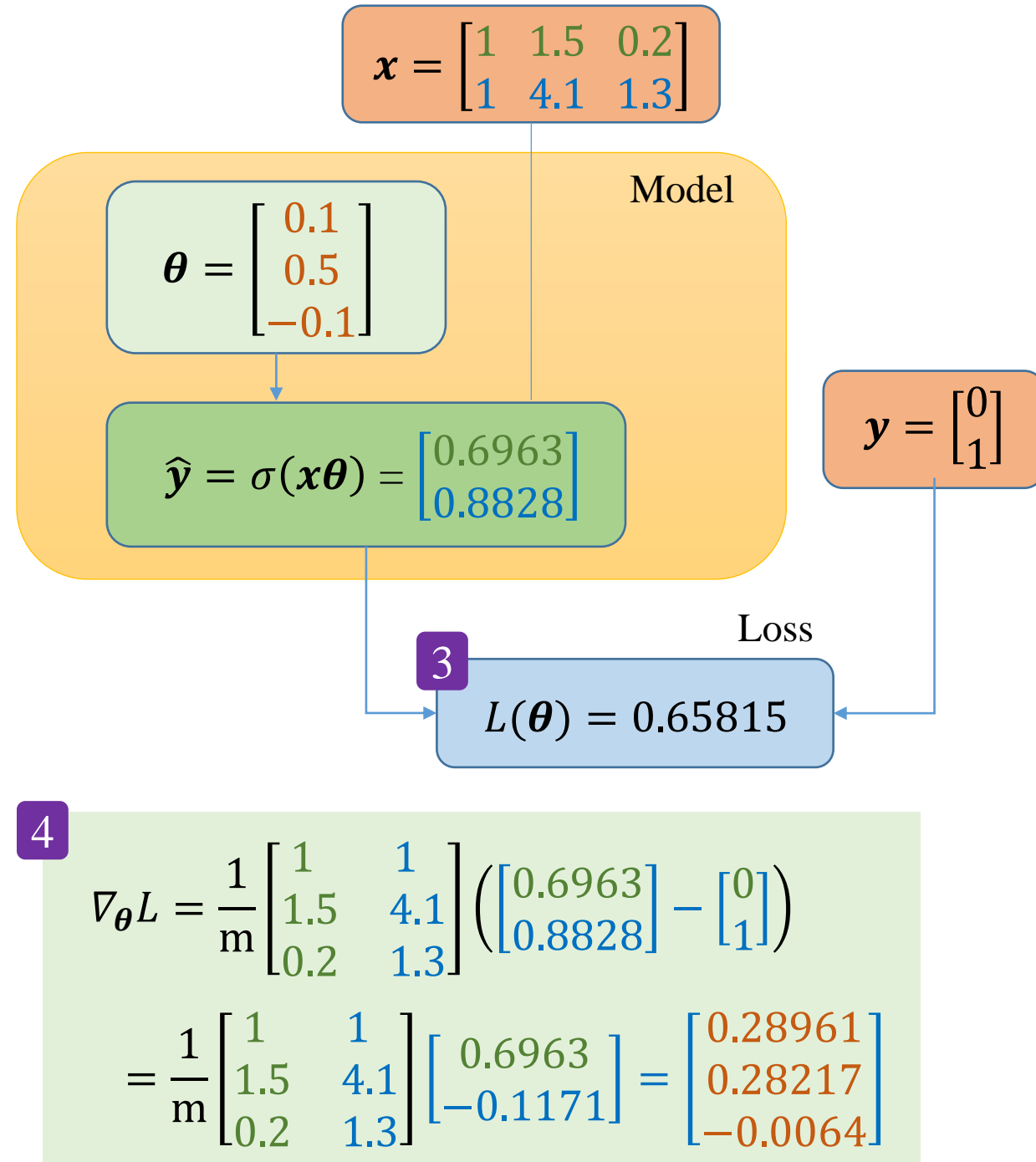
$$L(\hat{y}, y) = \frac{1}{m} (-y^T \log \hat{y} - (1-y)^T \log(1-\hat{y}))$$

4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$





## Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

1) Pick m samples from training data

2) Compute output  $\hat{y}$

Mini-batch m=2

$$z = x\theta$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

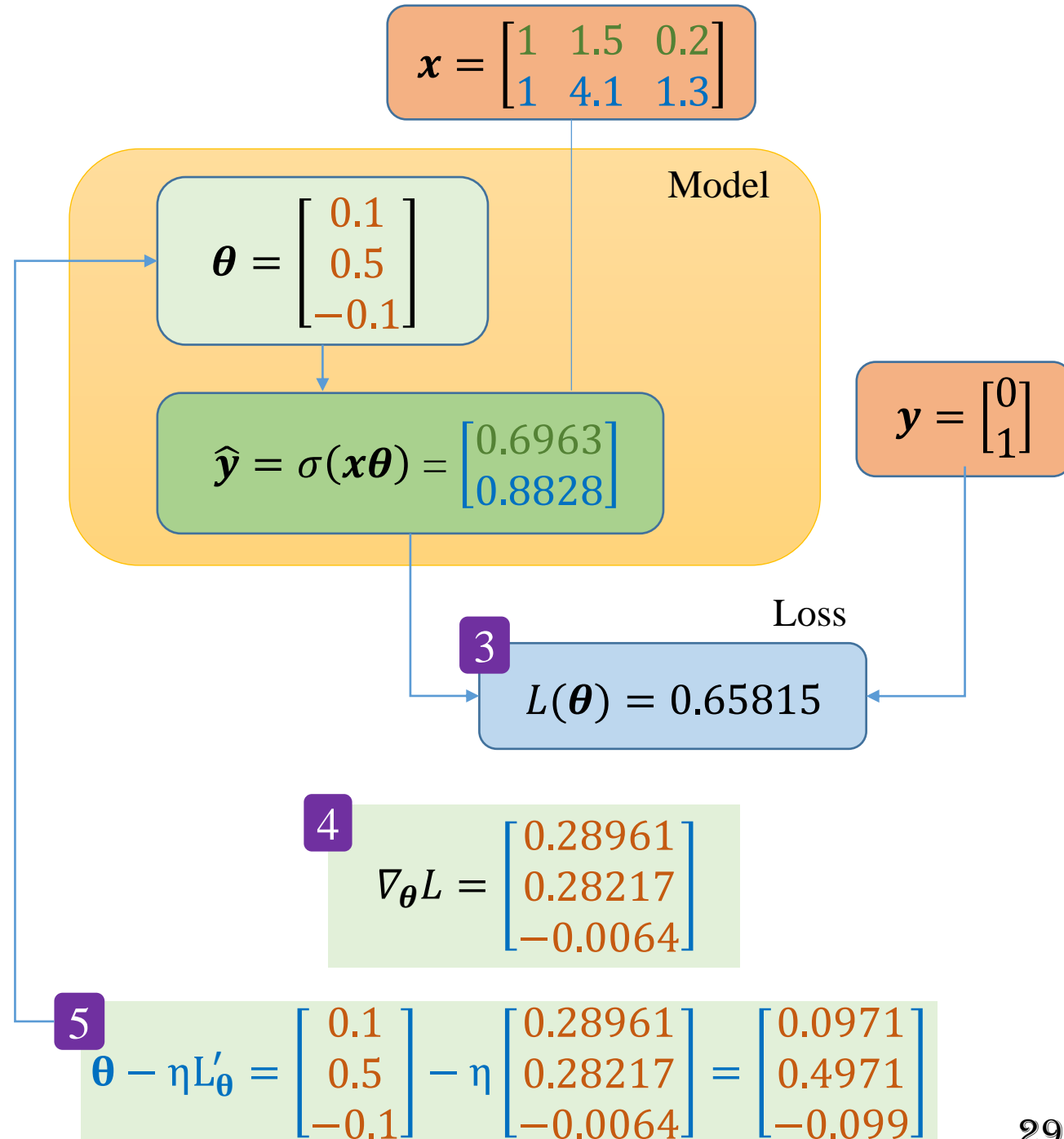
$$L(\hat{y}, y) = \frac{1}{m} (-y^T \log \hat{y} - (1-y)^T \log(1-\hat{y}))$$

4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$



# Outline

- Vectorization
- Optimiztion for 1+ samples
- Logistic Regression – Mini-batch
- Logistic Regression – Batch
- BCE and MSE Loss Functions
- Sigmoid and Tanh Function (Optional)

# Logistic Regression - Batch

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{x} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} \\ 1 & x_1^{(4)} & x_2^{(4)} \end{bmatrix}$$

1) Pick all the samples from training data

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

$$\mathbf{z} = \mathbf{x}\boldsymbol{\theta}$$
$$\hat{\mathbf{y}} = \sigma(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

3) Compute loss

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

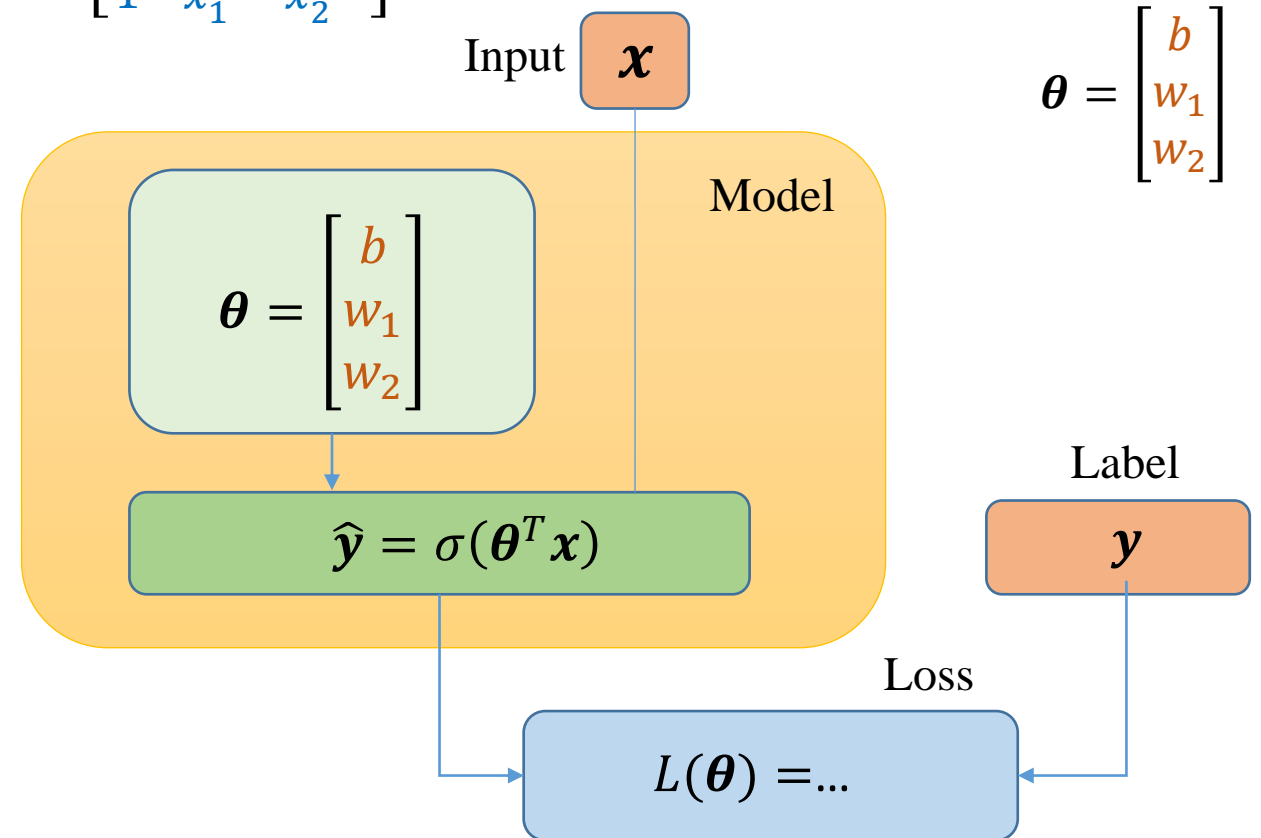
4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \mathbf{x}^T (\hat{\mathbf{y}} - \mathbf{y})$$

5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

$\eta$  is learning rate



## Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

1) Pick all the samples from training data

2) Compute output  $\hat{y}$

$$z = x\theta$$
$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

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

4) Compute derivative

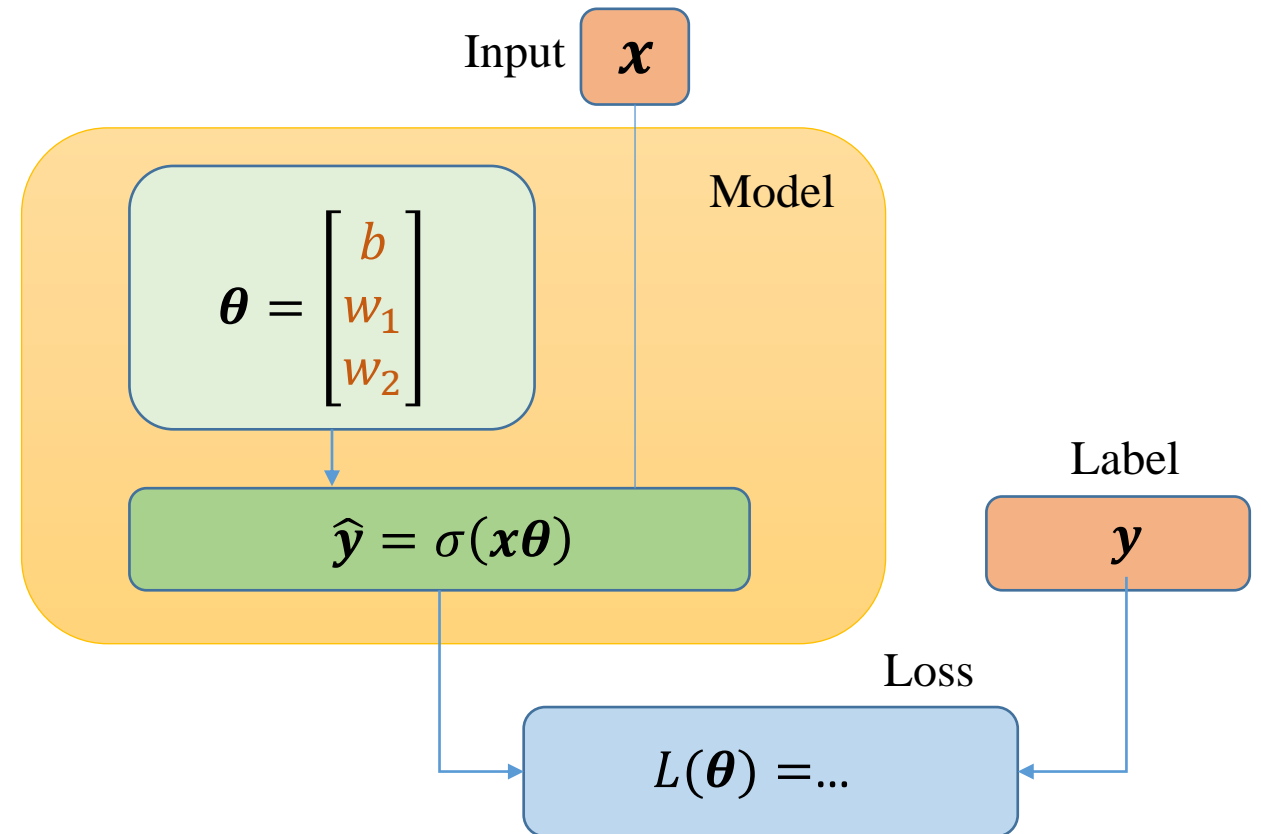
$$\nabla_{\theta} L = \frac{1}{N} x^T (\hat{y} - y)$$

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$

# Logistic Regression - Batch

$$x = \begin{bmatrix} 1 & 1.4 & 0.2 \\ 1 & 1.5 & 0.2 \\ 1 & 3.0 & 1.1 \\ 1 & 4.1 & 1.3 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$



## Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

$$\mathbf{x} = \begin{bmatrix} 1 & 1.4 & 0.2 \\ 1 & 1.5 & 0.2 \\ 1 & 3.0 & 1.1 \\ 1 & 4.1 & 1.3 \end{bmatrix}$$

1) Pick all the samples from training data

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

$$\mathbf{z} = \mathbf{x}\boldsymbol{\theta}$$

$$\hat{\mathbf{y}} = \sigma(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

3) Compute loss

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

4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \mathbf{x}^T (\hat{\mathbf{y}} - \mathbf{y})$$

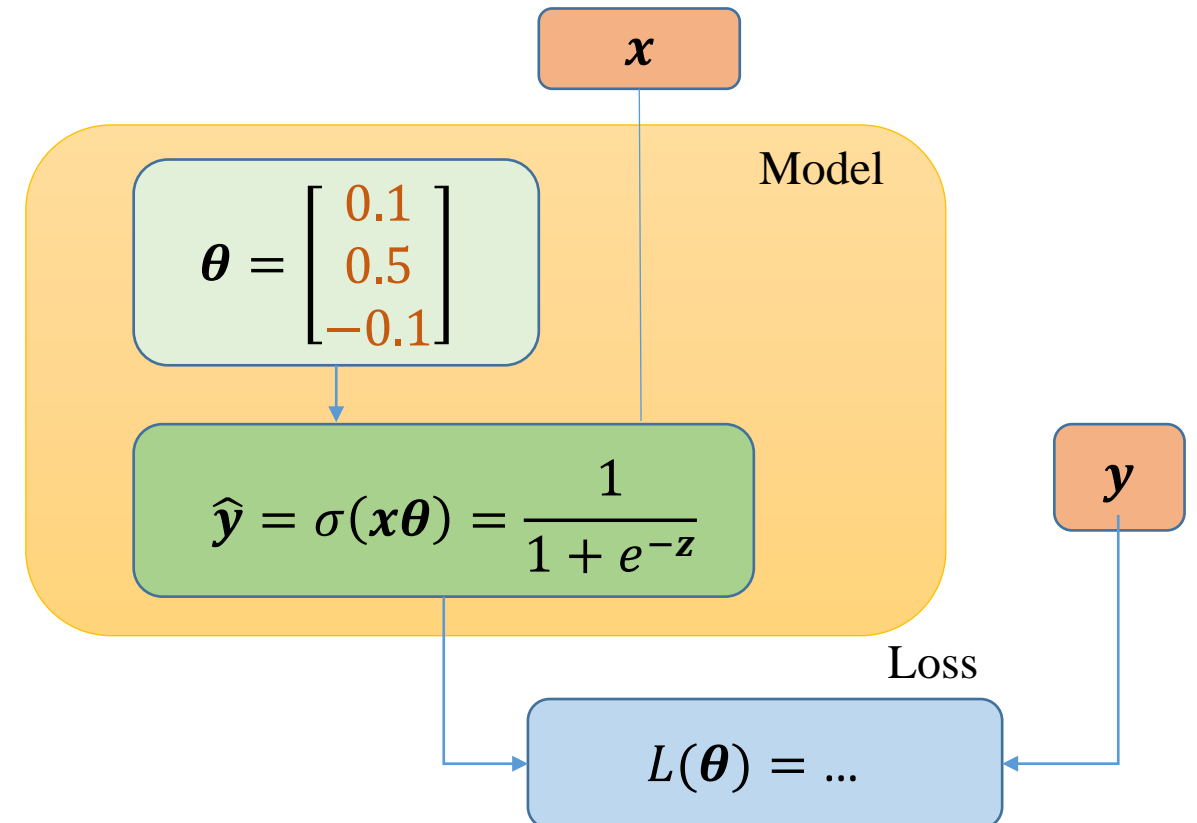
5) Update parameters

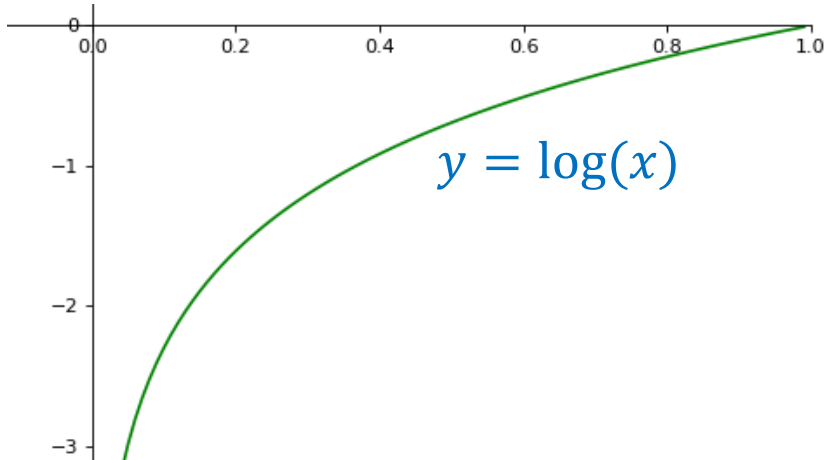
$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$

2

$$\hat{\mathbf{y}} = \sigma(\mathbf{x}\boldsymbol{\theta}) = \sigma \left( \begin{bmatrix} 1 & 1.4 & 0.2 \\ 1 & 1.5 & 0.2 \\ 1 & 3.0 & 1.1 \\ 1 & 4.1 & 1.3 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.78 \\ 0.83 \\ 1.49 \\ 2.02 \end{bmatrix} \right) = \begin{bmatrix} 0.6856 \\ 0.6963 \\ 0.8160 \\ 0.8828 \end{bmatrix}$$





3

$$L(\theta) = \frac{1}{N} \left\{ - \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}^T \log \begin{bmatrix} 0.6856 \\ 0.6963 \\ 0.8160 \\ 0.8828 \end{bmatrix} - \left( \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \right)^T \log \left( 1 - \begin{bmatrix} 0.6856 \\ 0.6963 \\ 0.8160 \\ 0.8828 \end{bmatrix} \right) \right\}$$

1) Pick all the samples from training data

2) Compute output  $\hat{y}$

$$z = x\theta$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

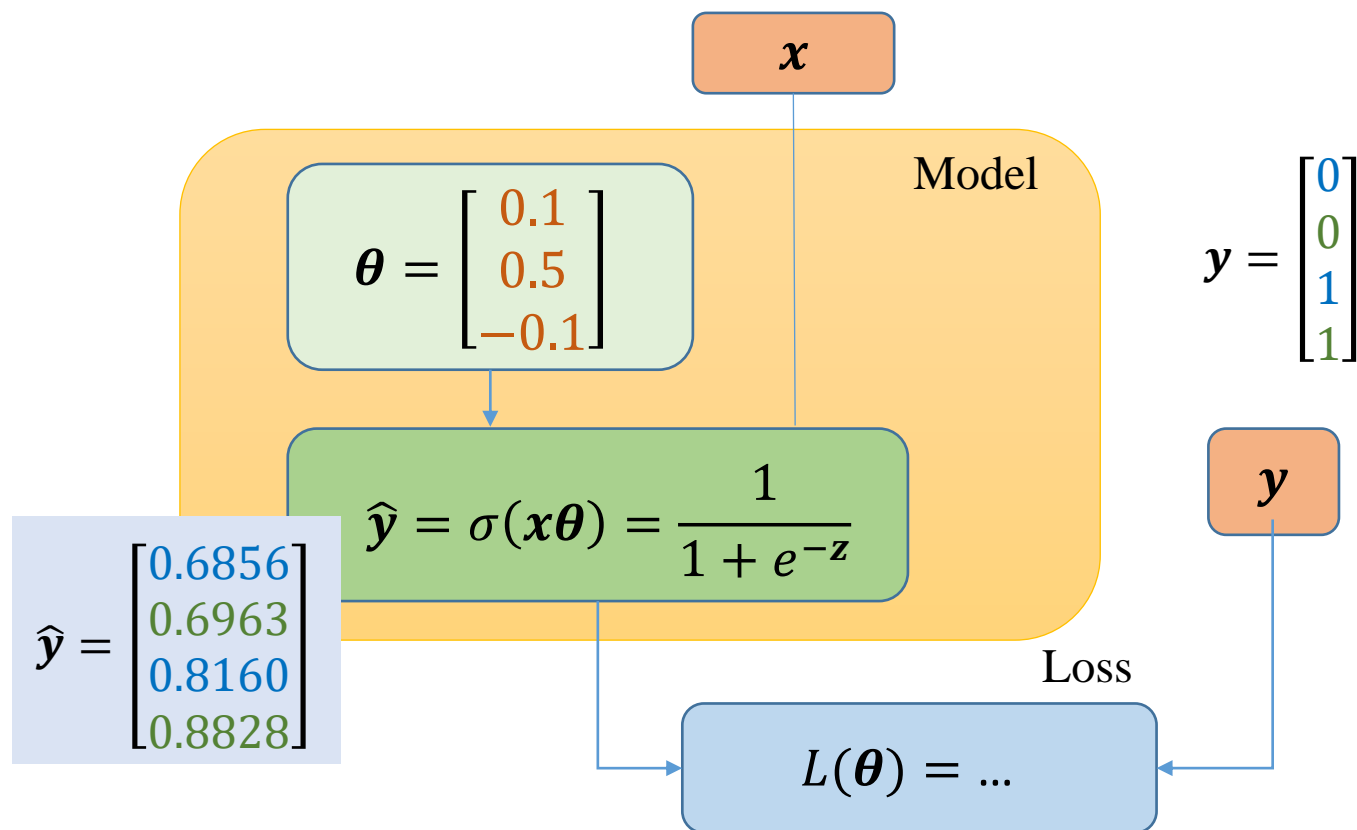
$$L(\hat{y}, y) = \frac{1}{N} (-y^T \log \hat{y} - (1-y)^T \log(1-\hat{y}))$$

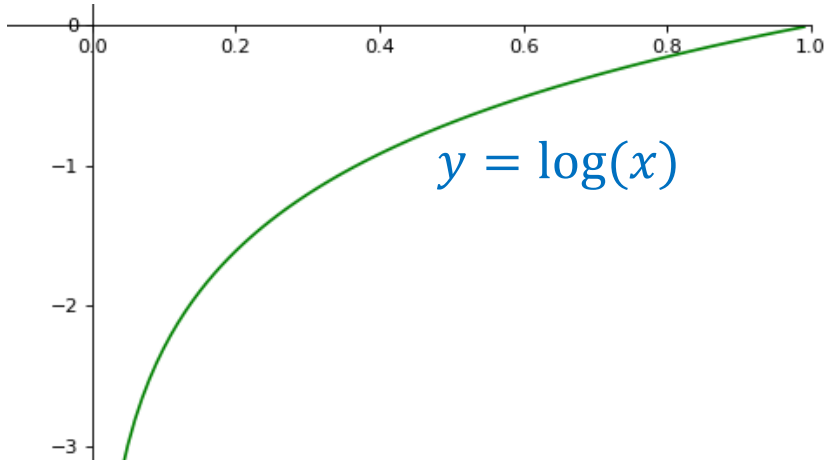
4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$





3

$$\begin{aligned}
 L(\theta) &= \frac{1}{N} \left\{ - \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}^T \log \left( \begin{bmatrix} 0.6856 \\ 0.6963 \\ 0.8160 \\ 0.8828 \end{bmatrix} \right) - \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}^T \log \left( \begin{bmatrix} 0.3144 \\ 0.3037 \\ 0.1840 \\ 0.1172 \end{bmatrix} \right) \right\} \\
 &= \frac{1}{N} (-\log 0.8160 - \log 0.8828 - \log 0.3144 - \log 0.3037) \\
 &= 0.6691
 \end{aligned}$$

1) Pick all the samples from training data

2) Compute output  $\hat{y}$

$$\begin{aligned}
 z &= x\theta \\
 \hat{y} &= \sigma(z) = \frac{1}{1 + e^{-z}}
 \end{aligned}$$

3) Compute loss

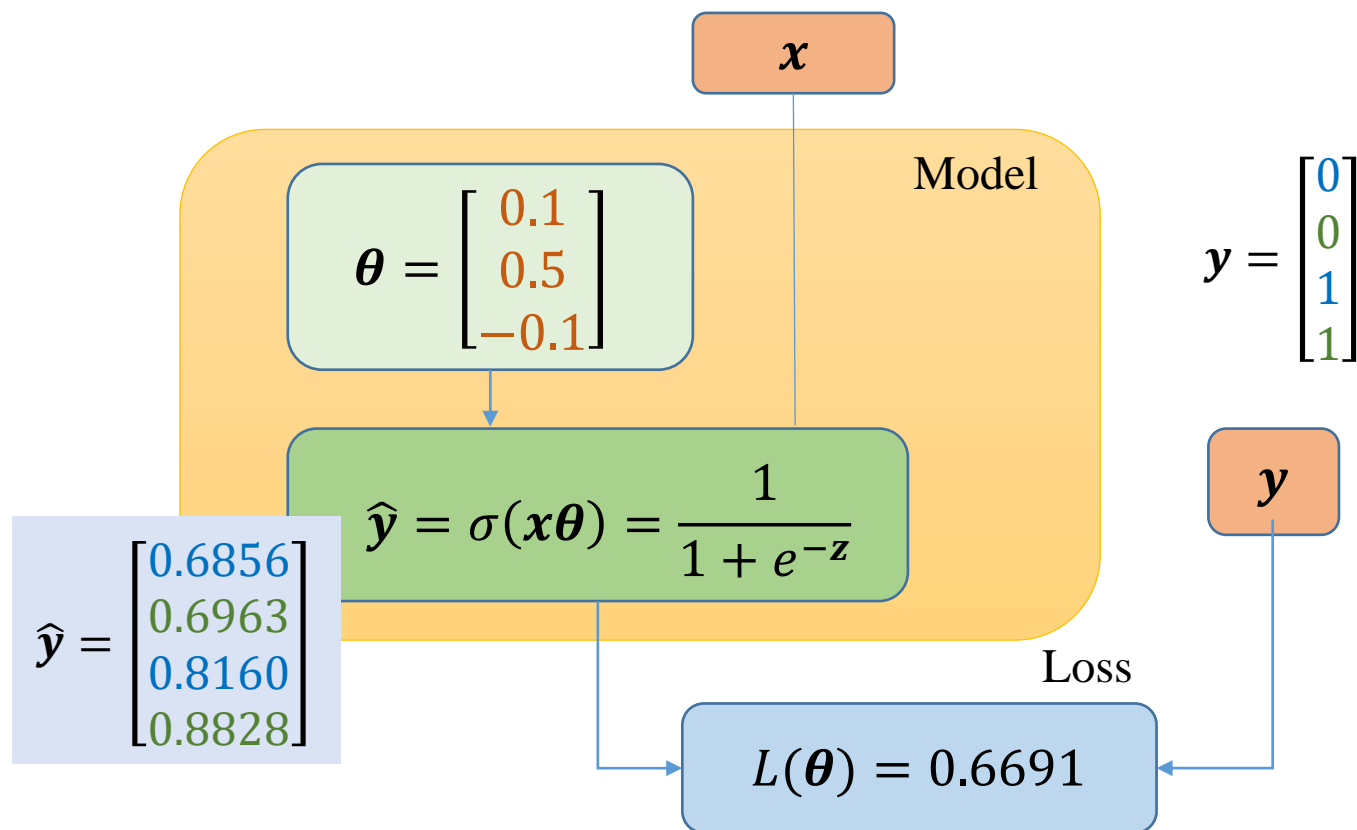
$$L(\hat{y}, y) = \frac{1}{N} (-y^T \log \hat{y} - (1-y)^T \log(1-\hat{y}))$$

4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$



$$\mathbf{x} = \begin{bmatrix} 1 & 1.4 & 0.2 \\ 1 & 1.5 & 0.2 \\ 1 & 3.0 & 1.1 \\ 1 & 4.1 & 1.3 \end{bmatrix} \quad \hat{\mathbf{y}} = \begin{bmatrix} 0.6856 \\ 0.6963 \\ 0.8160 \\ 0.8828 \end{bmatrix}$$

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

$$\mathbf{z} = \mathbf{x}\boldsymbol{\theta}$$

$$\hat{\mathbf{y}} = \sigma(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

- 3) Compute loss

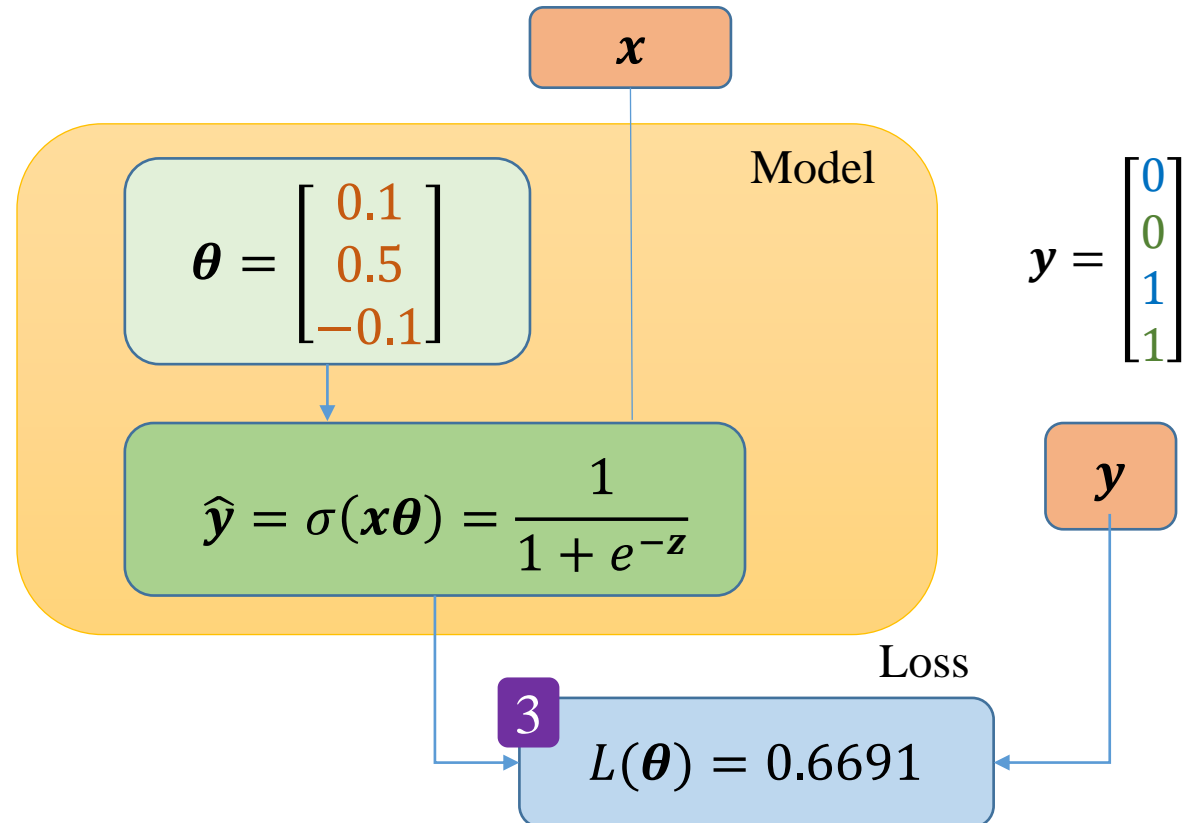
$$L(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{N} (-\mathbf{y}^T \log \hat{\mathbf{y}} - (1 - \mathbf{y})^T \log(1 - \hat{\mathbf{y}}))$$

- 4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \mathbf{x}^T (\hat{\mathbf{y}} - \mathbf{y})$$

- 5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L$$



$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1.4 & 1.5 & 3.0 & 4.1 \\ 0.2 & 0.2 & 1.1 & 1.3 \end{bmatrix} \left( \begin{bmatrix} 0.6856 \\ 0.6963 \\ 0.8160 \\ 0.8828 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \right)$$

$$= \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1.4 & 1.5 & 3.0 & 4.1 \\ 0.2 & 0.2 & 1.1 & 1.3 \end{bmatrix} \begin{bmatrix} 0.6856 \\ 0.6963 \\ -0.184 \\ -0.117 \end{bmatrix} = \begin{bmatrix} 0.2702 \\ 0.2431 \\ -0.019 \end{bmatrix}$$



## Dataset

Petal_Length	Petal_Width	Label
1.4	0.2	0
1.5	0.2	0
3	1.1	1
4.1	1.3	1

1) Pick all the samples from training data

2) Compute output  $\hat{y}$

$$z = x\theta$$
$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3) Compute loss

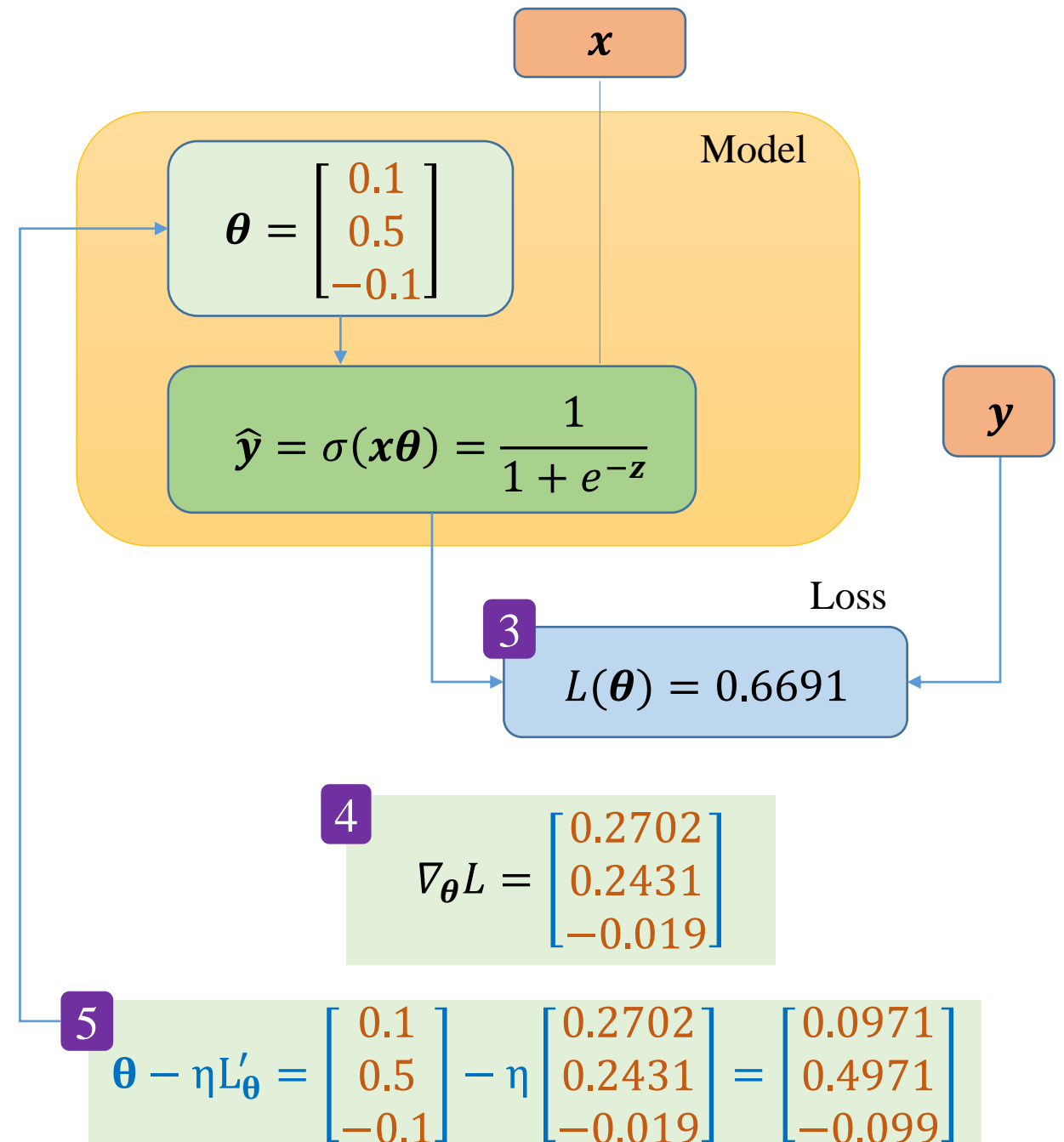
$$L(\hat{y}, y) = \frac{1}{N} (-y^T \log \hat{y} - (1-y)^T \log(1-\hat{y}))$$

4) Compute derivative

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

5) Update parameters

$$\theta = \theta - \eta \nabla_{\theta} L$$



# Outline

- Vectorization
- Optimiztion for 1+ samples
- Logistic Regression – Mini-batch
- Logistic Regression – Batch
- BCE and MSE Loss Functions
- Sigmoid and Tanh Function (Optional)

# Hessian Matrices

## ❖ Definition

The Hessian matrix or Hessian is a square matrix of second-order partial derivatives of a scalar-valued function

[https://en.wikipedia.org/wiki/Hessian\\_matrix](https://en.wikipedia.org/wiki/Hessian_matrix)

Given  $f(x, y)$

$f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$$H_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

Given  $f(x, y) = x^2 + 2x^2y + y^3$

$$\frac{\partial f}{\partial x} = 2x + 4xy$$

$$\frac{\partial f}{\partial y} = 2x^2 + 3y^2$$

$$H_f = \begin{bmatrix} 2 + 4y & 4x \\ 4x & 6y \end{bmatrix}$$

# Binary Cross-entropy

## ❖ Convex function

$$z = \theta^T x$$

Model and Loss

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$L = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$\frac{\partial L}{\partial \theta_i} = x_i(\hat{y} - y)$$

$$\frac{\partial^2 L}{\partial \theta_i^2} = \frac{\partial}{\partial \theta_i} [x_i(\hat{y} - y)] = x_i^2(\hat{y} - \hat{y}^2) \geq 0$$

$$x_i^2 \geq 0 \quad \hat{y} - \hat{y}^2 \in \left[0, \frac{1}{4}\right]$$

Derivative

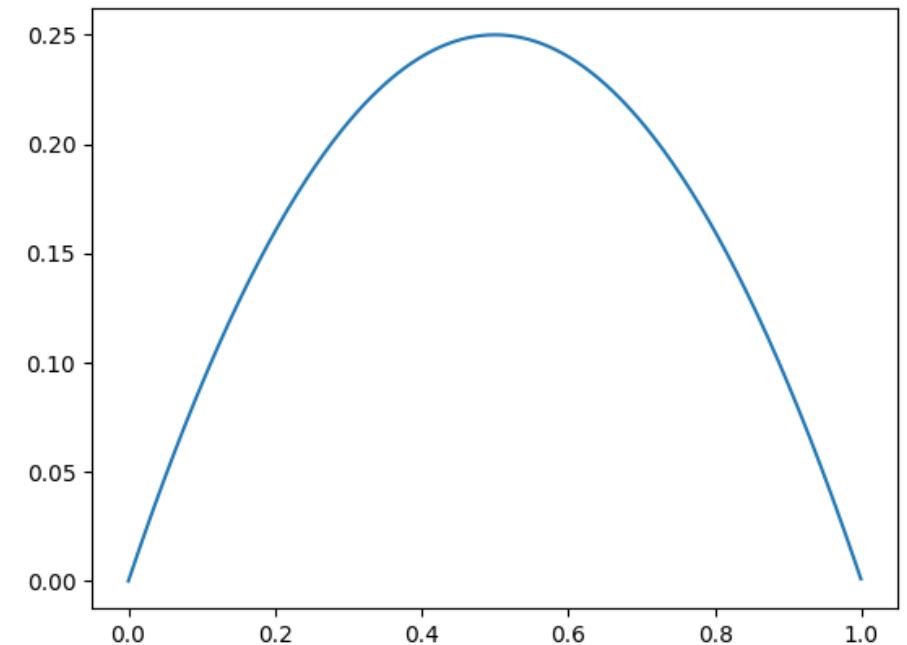
$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \theta_i}$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})}$$

$$\frac{\partial \hat{y}}{\partial z} = \hat{y}(1-\hat{y})$$

$$\frac{\partial z}{\partial \theta_i} = x_i$$

$$\frac{\partial L}{\partial \theta_i} = x_i(\hat{y} - y)$$



# Logistic Regression-MSE

## ❖ Construct loss

### Model and Loss

$$z = \theta^T x = x^T \theta$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

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

### Derivative

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \theta_i} \quad \frac{\partial \hat{y}}{\partial z} = \hat{y}(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) \quad \frac{\partial z}{\partial \theta_i} = x_i$$

$$\frac{\partial L}{\partial \theta_i} = 2x_i(\hat{y} - y)\hat{y}(1 - \hat{y})$$

# Mean Squared Error

## Model and Loss

$$z = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

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

## Derivative

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \theta_i} \quad \frac{\partial \hat{y}}{\partial z} = \hat{y}(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) \quad \frac{\partial z}{\partial \theta_i} = x_i$$

$$\frac{\partial L}{\partial \theta_i} = 2x_i(\hat{y} - y)\hat{y}(1 - \hat{y})$$

$$\frac{\partial L}{\partial \theta_i} = 2x_i(\hat{y} - y)\hat{y}(1 - \hat{y}) = 2x_i(-\hat{y}^3 + \hat{y}^2 - y\hat{y} + y\hat{y}^2)$$

$$\frac{\partial^2 L}{\partial \theta_i^2} = \frac{\partial}{\partial \theta_i} [2x_i(-\hat{y}^3 + \hat{y}^2 - y\hat{y} + y\hat{y}^2)]$$

$$= 2x_i[-3\hat{y}^2 x_i \hat{y}(1 - \hat{y}) + 2x_i \hat{y} \hat{y}(1 - \hat{y}) - y x_i \hat{y}(1 - \hat{y}) + 2x_i y \hat{y} \hat{y}(1 - \hat{y})]$$

$$= 2x_i^2 \hat{y}(1 - \hat{y})[-3\hat{y}^2 + 2\hat{y} - y + 2y\hat{y}]$$

# Mean Squared Error

$$\frac{\partial^2 L}{\partial \theta_i^2} = 2x_i^2 \hat{y}(1 - \hat{y})[-3\hat{y}^2 + 2\hat{y} - y + 2y\hat{y}]$$

$$x_i^2 \geq 0$$

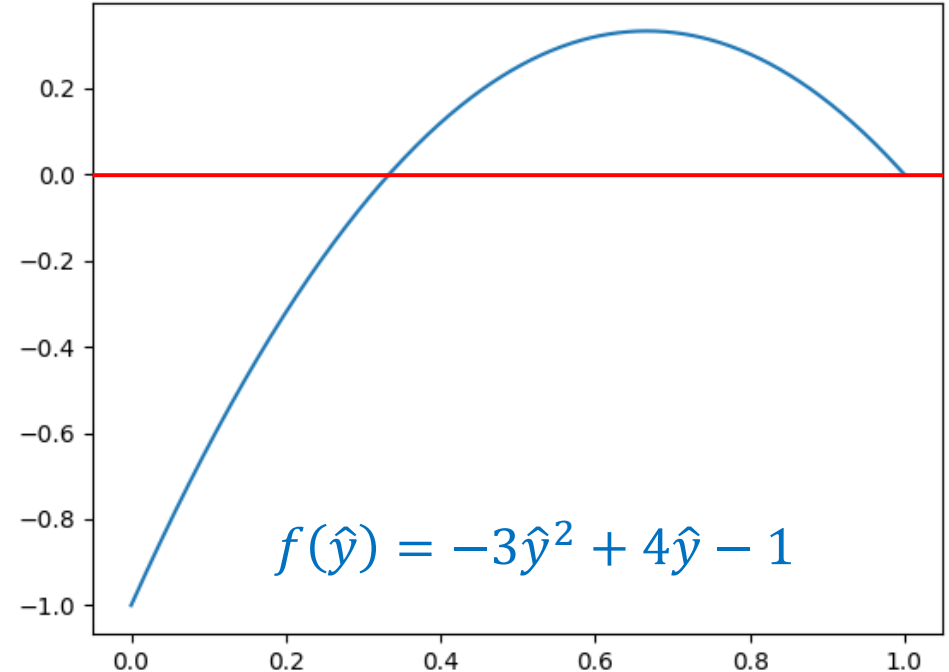
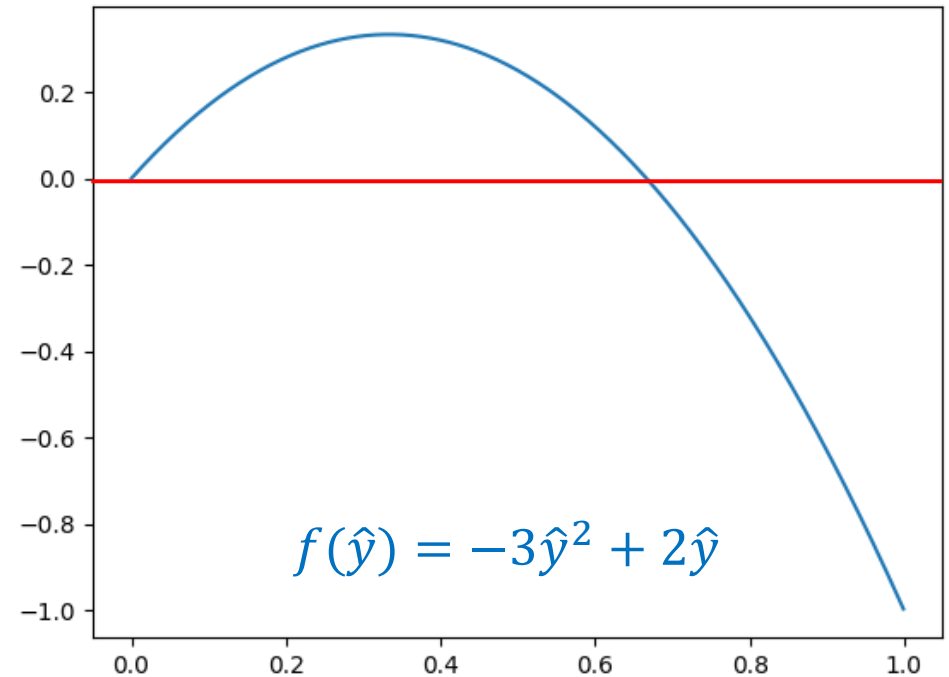
$$\hat{y}(1 - \hat{y}) \in \left[0, \frac{1}{4}\right]$$

$$y = 0$$

$$f(\hat{y}) = -3\hat{y}^2 + 2\hat{y}$$

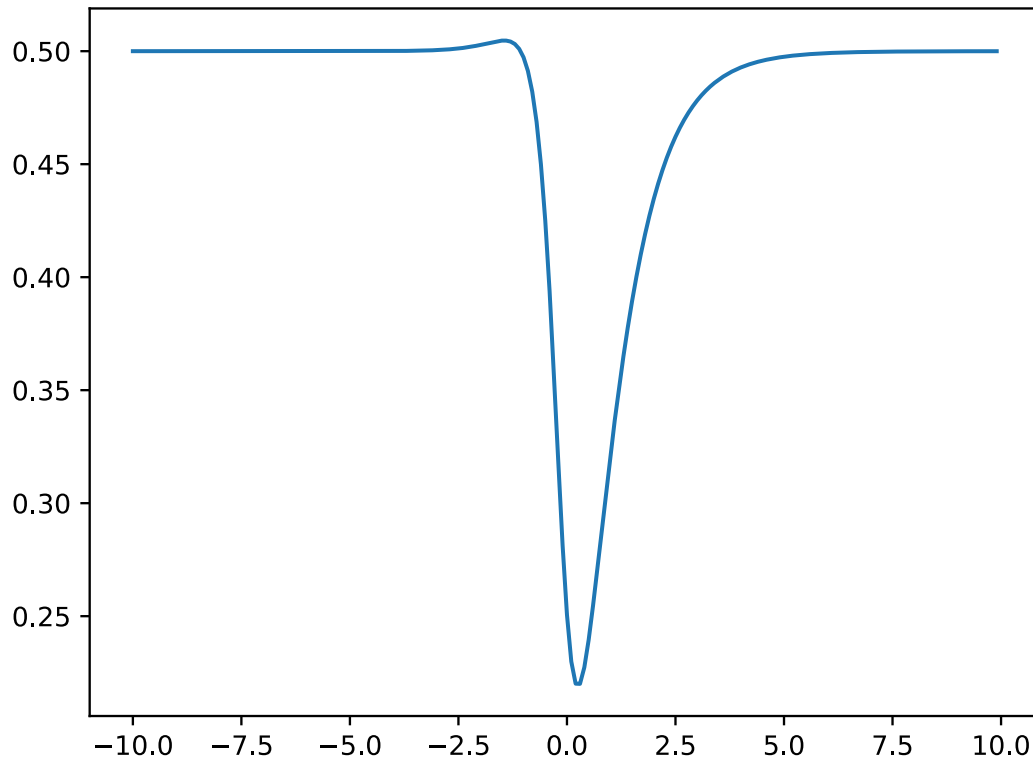
$$y = 1$$

$$f(\hat{y}) = -3\hat{y}^2 + 4\hat{y} - 1$$

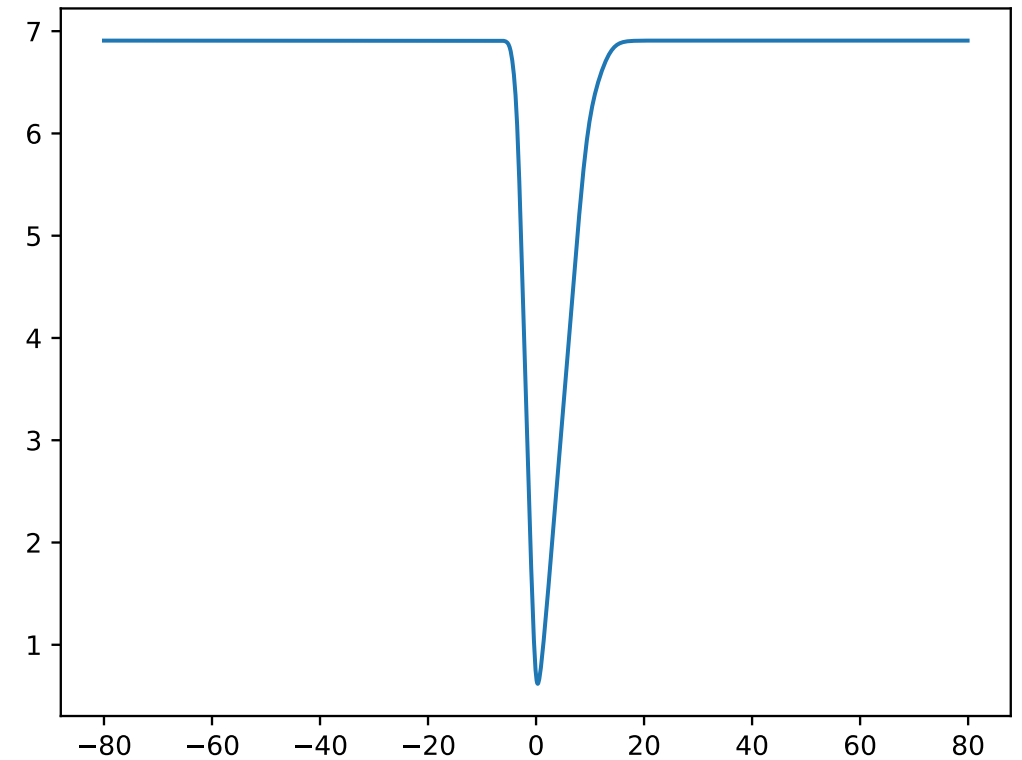


# MSE and BCE

## ❖ Visualization



Mean Squared Error



Binary Cross-Entropy

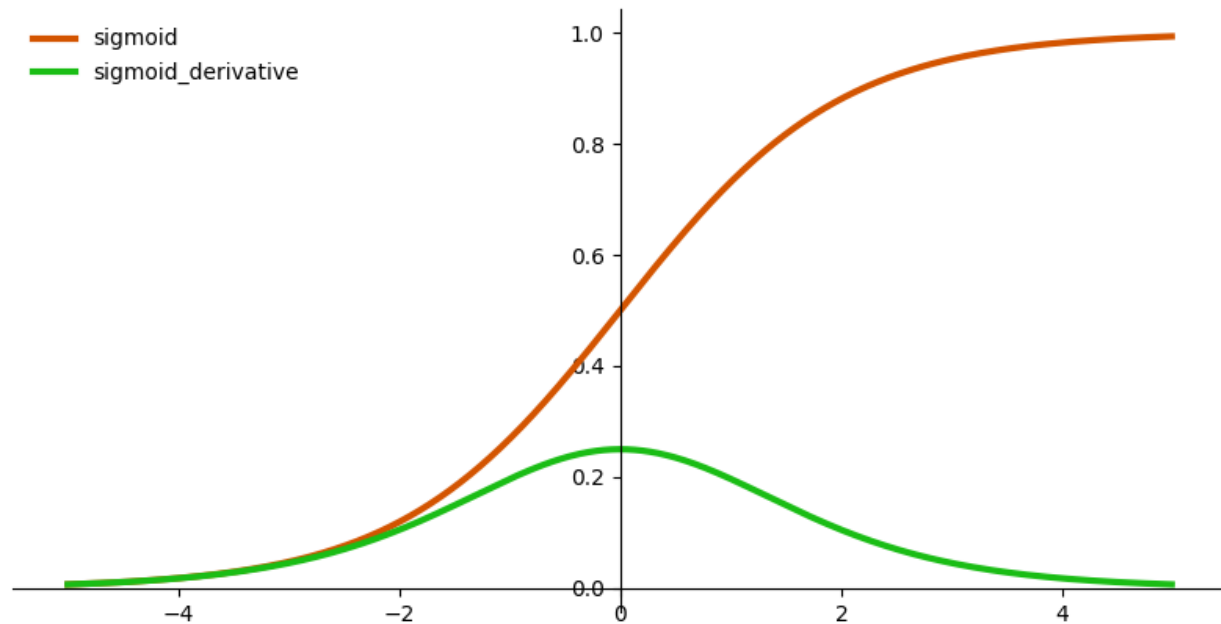


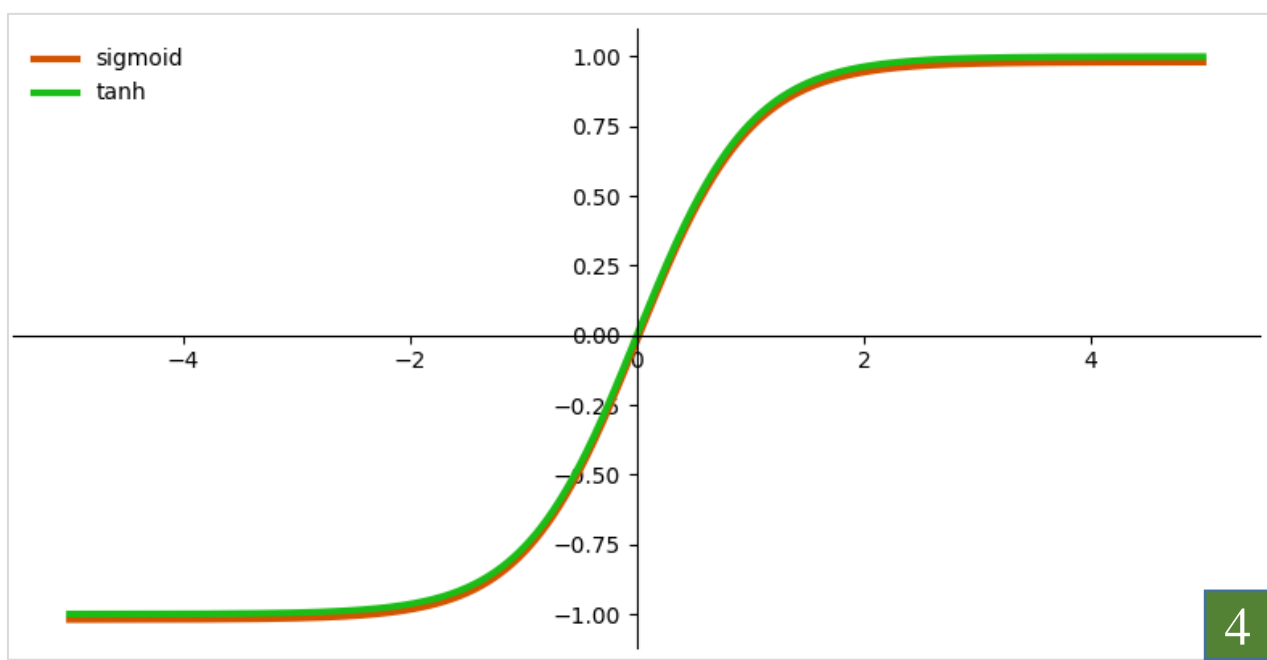
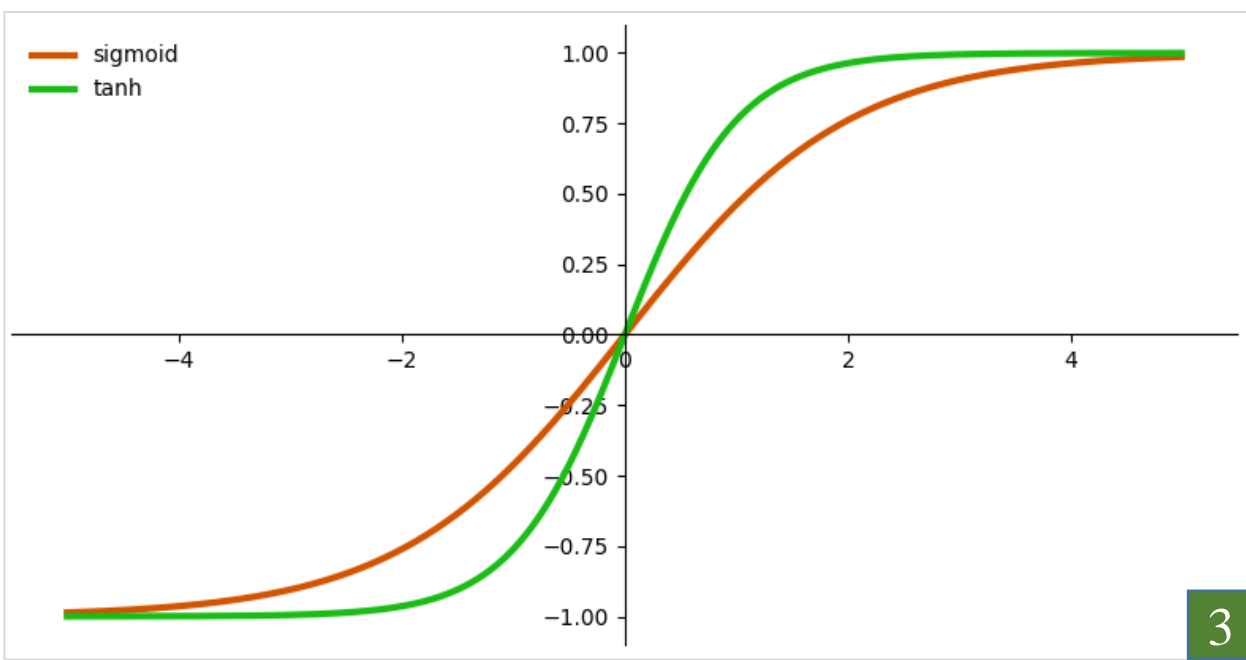
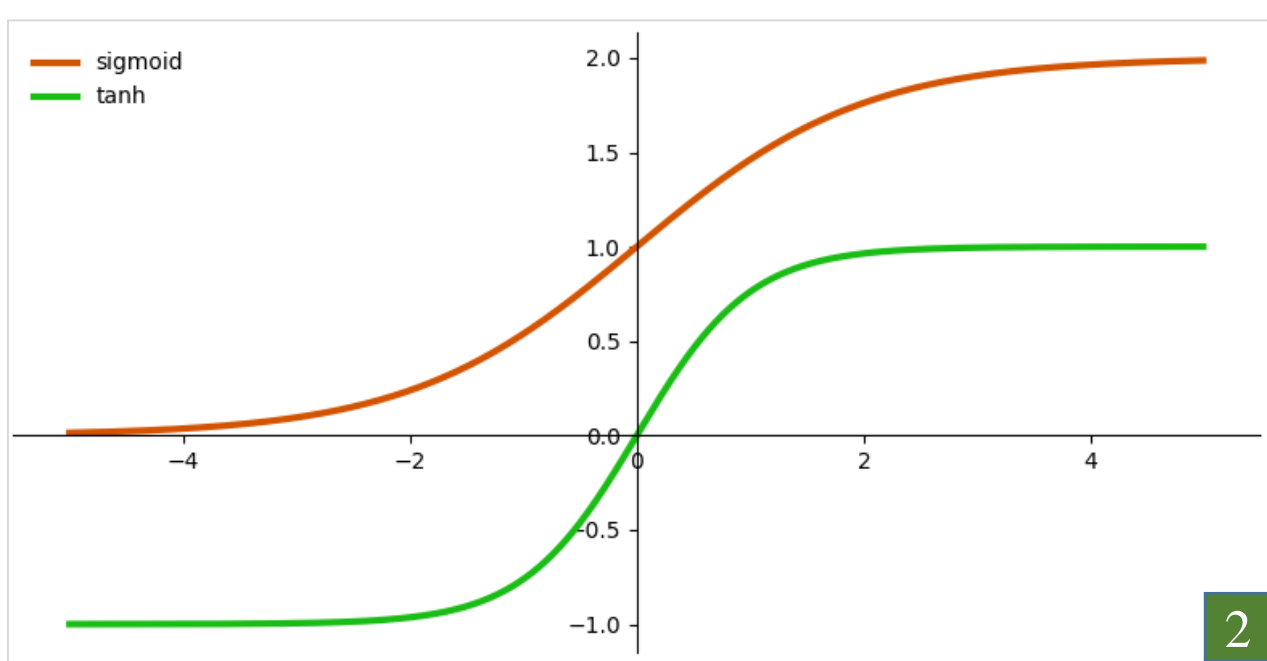
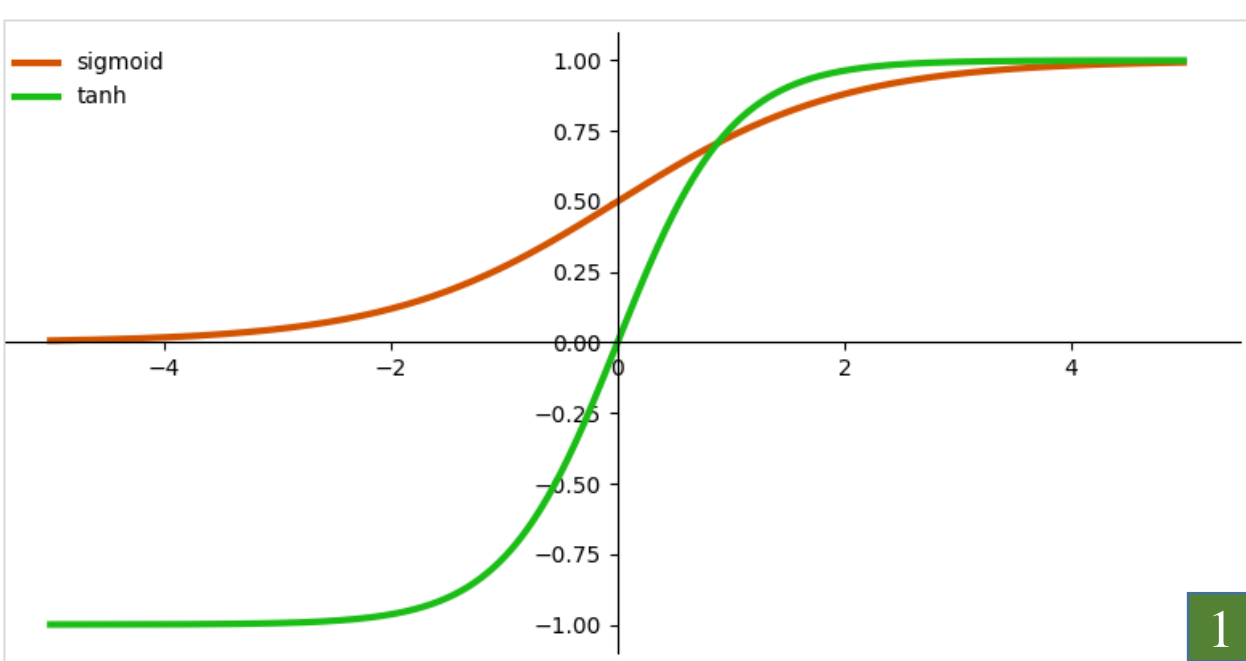


# Sigmoid and Tanh Functions

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



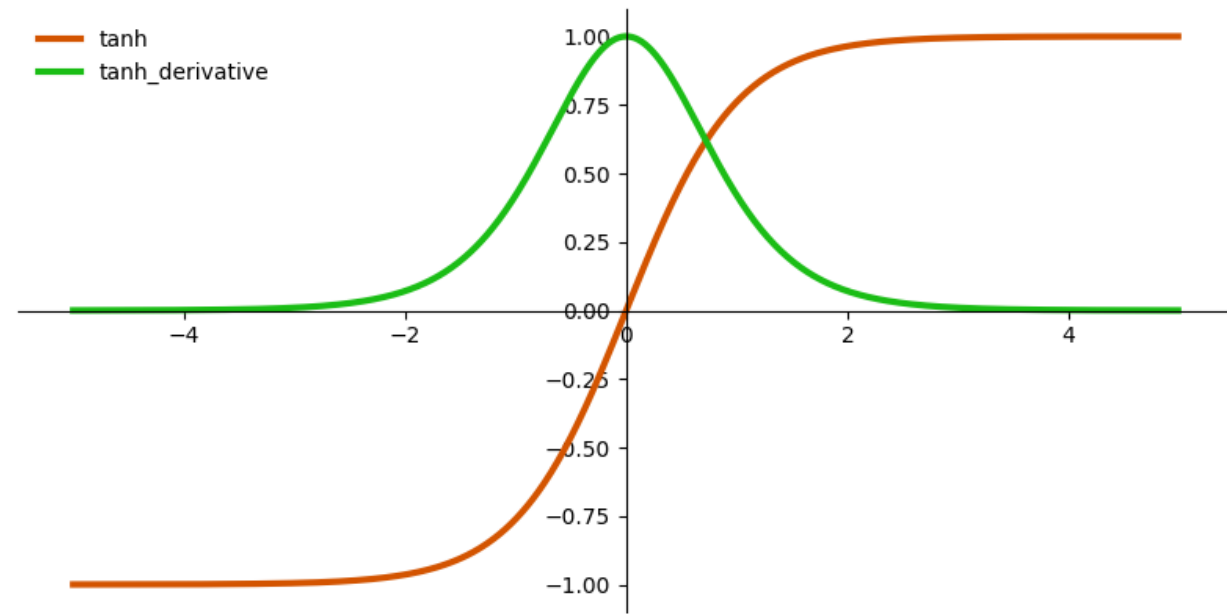
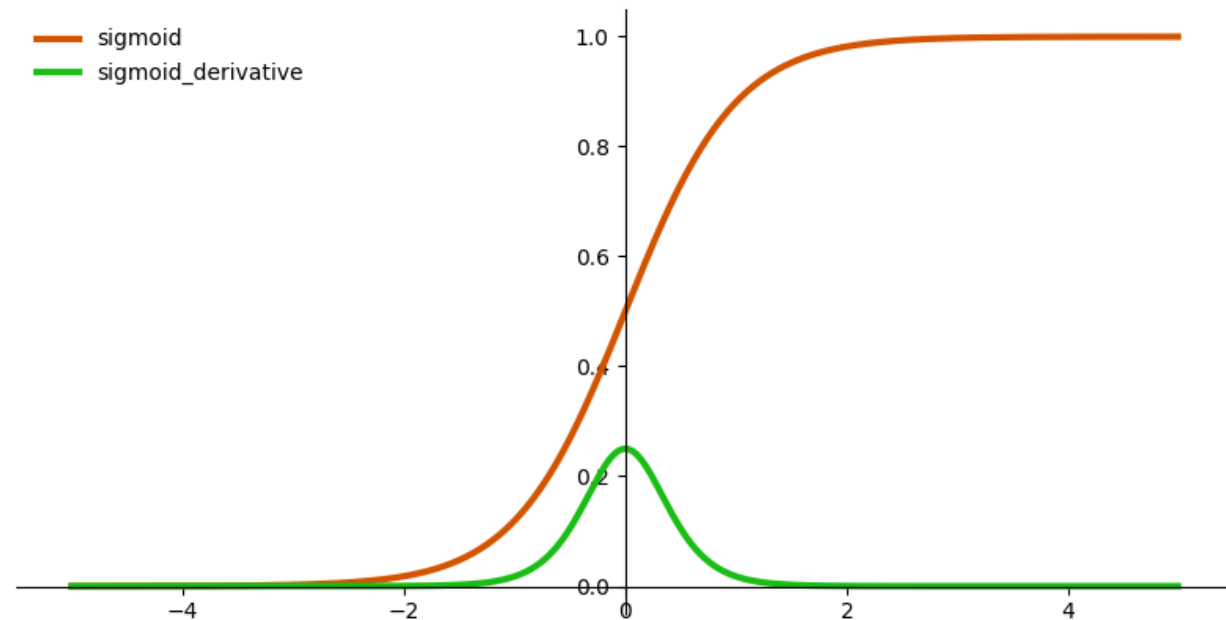


# Sigmoid and Tanh Functions

$$\text{sigmoid}(2x) = \frac{1}{1 + e^{-2x}}$$

$$\tanh(x) = 2 \times \frac{1}{1 + e^{-2x}} - 1$$

$$\tanh(x) = 2 \times \text{sigmoid}(2x) - 1$$



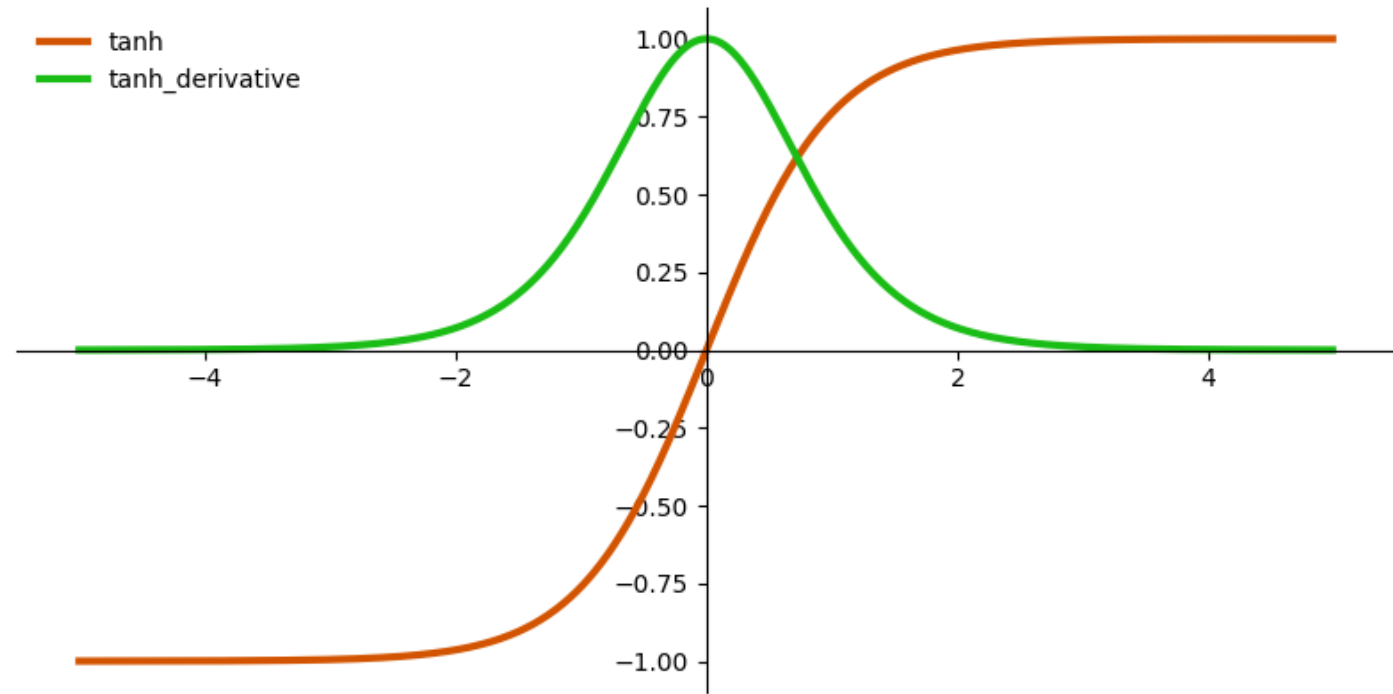
# Outline

- Vectorization
- Optimiztion for 1+ samples
- Logistic Regression – Mini-batch
- Logistic Regression – Batch
- BCE and MSE Loss Functions
- Sigmoid and Tanh Function (Optional)

# Tanh function

$$\begin{aligned}\tanh(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} \\ &= 1 - \frac{2}{e^{2x} + 1}\end{aligned}$$

$$\begin{aligned}\tanh(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} \\ &= -\frac{e^{-2x} - 1}{e^{-2x} + 1} = \frac{2}{e^{-2x} + 1} - 1\end{aligned}$$



# Tanh function

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} = 1 - \frac{2}{e^{2x} + 1}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} = -\frac{e^{-2x} - 1}{e^{-2x} + 1} = \frac{2}{e^{-2x} + 1} - 1$$

$$\begin{aligned} \tanh'(x) &= \left( \frac{e^x - e^{-x}}{e^x + e^{-x}} \right)' = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} \\ &= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} \\ &= 1 - \left( \frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2 = 1 - \tanh^2(x) \end{aligned}$$

# Tanh function

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} = 1 - \frac{2}{e^{2x} + 1}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} = -\frac{e^{-2x} - 1}{e^{-2x} + 1} = \frac{2}{e^{-2x} + 1} - 1$$

$$\begin{aligned} \tanh'(x) &= \left( \frac{2}{e^{-2x} + 1} - 1 \right)' = \frac{4e^{-2x}}{(e^{-2x} + 1)^2} = 4 \left( \frac{e^{-2x} + 1 - 1}{(e^{-2x} + 1)^2} \right) \\ &= 4 \left( \frac{1}{e^{-2x} + 1} - \frac{1}{(e^{-2x} + 1)^2} \right) = - \left( \frac{4}{(e^{-2x} + 1)^2} - \frac{4}{e^{-2x} + 1} \right) \\ &= - \left( \frac{4}{(e^{-2x} + 1)^2} - \frac{4}{e^{-2x} + 1} + 1 - 1 \right) = 1 - \left( \frac{2}{e^{-2x} + 1} - 1 \right)^2 = 1 - \tanh^2(x) \end{aligned}$$

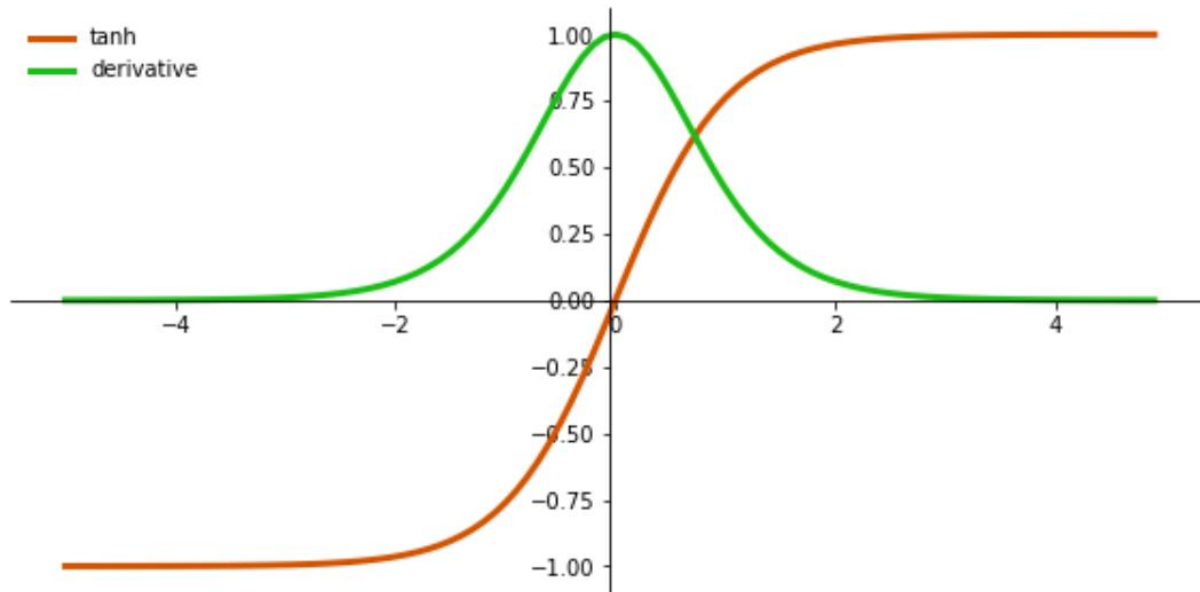


# Logistic Regression

## Tanh

### ❖ Construct loss

$$\tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$



### Model and Loss

$$z = \theta^T x = x^T \theta$$

$$\hat{y} = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\hat{y}_s = \frac{\hat{y} + 1}{2}$$

$$L = -y \log(\hat{y}_s) - (1 - y) \log(1 - \hat{y}_s)$$

### Derivative

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}_s} \frac{\partial \hat{y}_s}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \theta_i}$$

$$\frac{\partial L}{\partial \hat{y}_s} = -\frac{y}{\hat{y}_s} + \frac{1 - y}{1 - \hat{y}_s} = \frac{\hat{y}_s - y}{\hat{y}_s(1 - \hat{y}_s)}$$

$$\frac{\partial \hat{y}_s}{\partial \hat{y}} = \frac{1}{2}$$

$$\frac{\partial \hat{y}}{\partial z} = 1 - \hat{y}^2$$

$$\frac{\partial z}{\partial \theta_i} = x_i$$

$$\frac{\partial L}{\partial \theta_i} = x_i \frac{(\hat{y}_s - y)(1 - \hat{y}^2)}{2\hat{y}_s(1 - \hat{y}_s)}$$

# Logistic Regression

## Tanh

### ❖ Construct loss

$$\tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

$$z = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta}$$

Model and Loss

$$\hat{y} = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\hat{y}_s = \frac{\hat{y} + 1}{2}$$

$$L = -y \log(\hat{y}_s) - (1 - y) \log(1 - \hat{y}_s)$$

Derivative

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}_s} \frac{\partial \hat{y}_s}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \theta_i}$$

$$\frac{\partial L}{\partial \hat{y}_s} = -\frac{y}{\hat{y}_s} + \frac{1 - y}{1 - \hat{y}_s} = \frac{\hat{y}_s - y}{\hat{y}_s(1 - \hat{y}_s)}$$

$$\frac{\partial \hat{y}_s}{\partial \hat{y}} = \frac{1}{2}$$

$$\frac{\partial \hat{y}}{\partial z} = 1 - \hat{y}^2$$

$$\frac{\partial z}{\partial \theta_i} = x_i$$

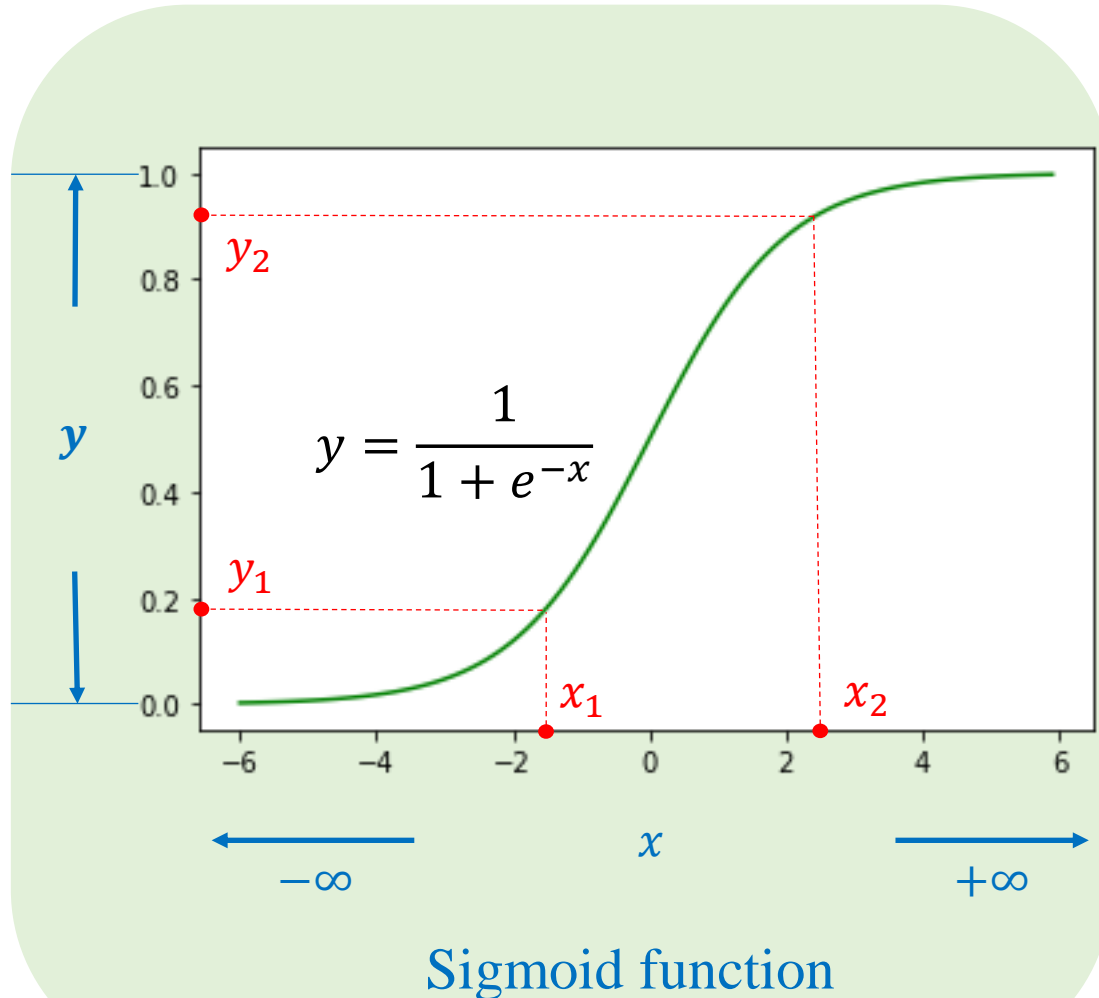
$$\frac{\partial L}{\partial \theta_i} = x_i \frac{(\hat{y}_s - y)(1 - \hat{y}^2)}{2\hat{y}_s(1 - \hat{y}_s)}$$

$$\frac{\partial L}{\partial \theta_i} = x_i \frac{(\frac{\hat{y} + 1}{2} - y)(1 - \hat{y}^2)}{2 \frac{\hat{y} + 1}{2} (1 - \frac{\hat{y} + 1}{2})}$$

$$\frac{\partial L}{\partial \theta_i} = x_i \frac{(\hat{y} + 1 - 2y)(1 - \hat{y}^2)}{(\hat{y} + 1)(1 - \hat{y})}$$

$$\frac{\partial L}{\partial \theta_i} = x_i(\hat{y} + 1 - 2y)$$

# Summary



1) Pick all the samples from training data

2) Compute output  $\hat{y}$

$$\mathbf{z} = \mathbf{x}\boldsymbol{\theta}$$

$$\hat{y} = \sigma(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

3) Compute loss (binary cross-entropy)

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

4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \mathbf{x}^T (\hat{\mathbf{y}} - \mathbf{y})$$

5) Update parameters

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

$\eta$  is learning rate

