

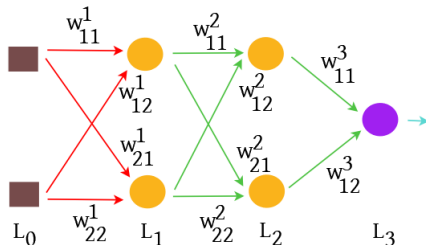
# Deep Learning - Theory and Practice

*IE 643*  
*Lecture 6*

Aug 14, 2025.

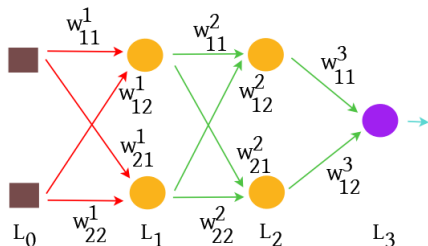
- 1 Recap
  - MLP-Data Perspective
- 2 Optimization Concepts
  - Gradient Descent
  - Stochastic Gradient Descent
  - Mini-batch SGD
- 3 Sample-wise Gradient Computation
  - MLP for prediction tasks

# Multi Layer Perceptron - Data Perspective



- **Input:** Training Data  $D = \{(x^s, y^s)\}_{s=1}^S$ .
- For each sample  $x^s$  the prediction  $\hat{y}^s = \text{MLP}(x^s)$ .
- **Error:**  $e^s = E(y^s, \hat{y}^s)$ .
- **Aim:** To minimize  $\sum_{s=1}^S e^s$ .

# Multi Layer Perceptron - Data Perspective

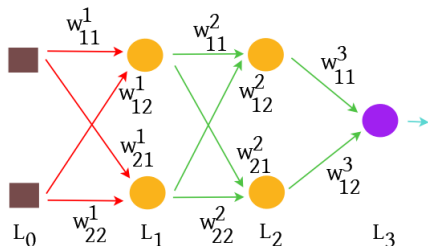


## Optimization perspective

- Given training data  $D = \{(x^s, y^s)\}_{s=1}^S$ ,

$$\min \sum_{s=1}^S e^s$$

# Multi Layer Perceptron - Data Perspective

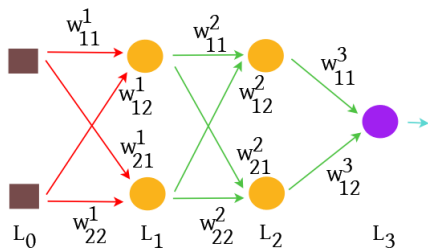


## Optimization perspective

- Given training data  $D = \{(x^s, y^s)\}_{s=1}^S$ ,

$$\min \sum_{s=1}^S e^s = \sum_{s=1}^S E(y^s, \hat{y}^s)$$

# Multi Layer Perceptron - Data Perspective

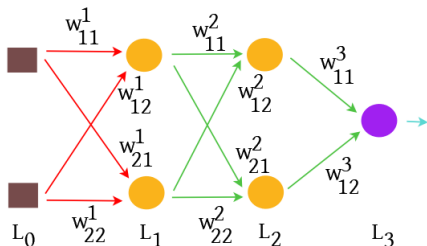


## Optimization perspective

- Given training data  $D = \{(x^s, y^s)\}_{s=1}^S$ ,

$$\min \sum_{s=1}^S e^s = \sum_{s=1}^S E(y^s, \hat{y}^s) = \sum_{s=1}^S E(y^s, \text{MLP}(x^s))$$

# Multi Layer Perceptron - Data Perspective



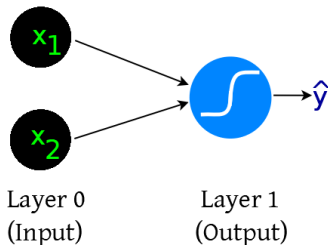
## Optimization perspective

- Given training data  $D = \{(x^s, y^s)\}_{s=1}^S$ ,

$$\min \sum_{s=1}^S e^s = \sum_{s=1}^S E(y^s, \hat{y}^s) = \sum_{s=1}^S E(y^s, \text{MLP}(x^s))$$

- Note:** The minimization is over the weights of the MLP  $W^1, \dots, W^L$ , where  $L$  denotes number of layers in MLP.

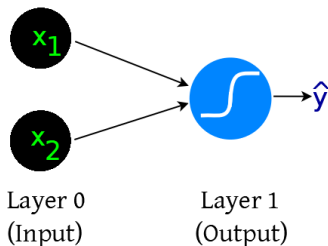
# MLP - Data Perspective: A Simple Example



$$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2) = \frac{1}{1 + \exp(-[w_{11}^1 x_1 + w_{12}^1 x_2])}$$



# MLP - Data Perspective: A Simple Example

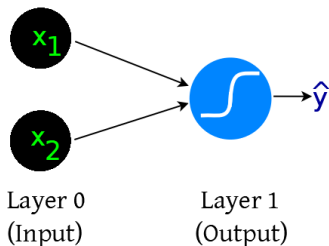


$$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2) = \frac{1}{1 + \exp(-[w_{11}^1 x_1 + w_{12}^1 x_2])}$$

**Property of 0-1 sigmoid  $\sigma : \mathbb{R} \rightarrow [0, 1]$**

- $\sigma$  is continuous
- $\sigma$  is monotonic
- $\sigma(z) \rightarrow \begin{cases} 0 & \text{if } z \rightarrow -\infty \\ 1 & \text{if } z \rightarrow +\infty \end{cases}$

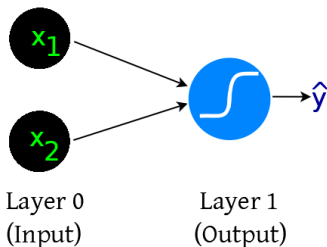
# MLP - Data Perspective: A Simple Example



- Let

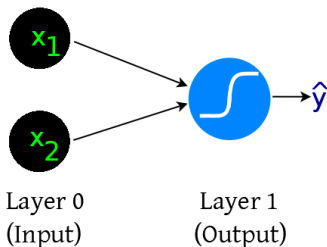
$$D = \{(x^1 = (-3, -3), y^1 = 1), \\ (x^2 = (-2, -2), y^2 = 1), \\ (x^3 = (4, 4), y^3 = 0), \\ (x^4 = (2, -5), y^4 = 0)\}.$$

# MLP - Data Perspective: A Simple Example



$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$

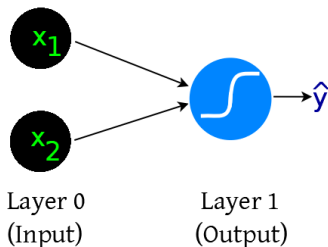
# MLP - Data Perspective: A Simple Example



$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$

- **Assume:**  $\text{Err}(y, \hat{y}) = (y - \hat{y})^2$ .

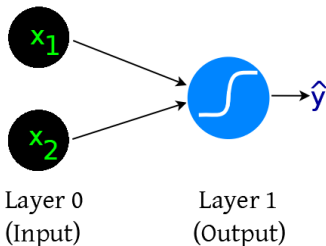
# MLP - Data Perspective: A Simple Example



$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$

- **Assume:**  $\text{Err}(y, \hat{y}) = (y - \hat{y})^2$ .
- Popularly called the **squared error**.

# MLP - Data Perspective: A Simple Example

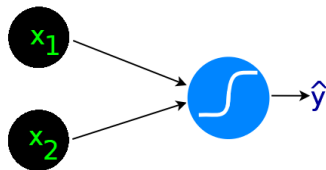


$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$

- Total error (or loss):

$$E = \sum_{i=1}^4 e^i = \sum_{i=1}^4 \text{Err}(y^i, \hat{y}^i)$$

# MLP - Data Perspective: A Simple Example



Layer 0  
(Input)

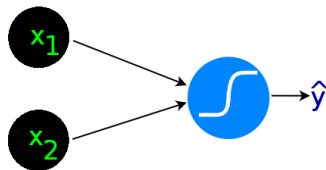
Layer 1  
(Output)

$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$

- Total error (or loss):

$$E = \sum_{i=1}^4 \left( y^i - \frac{1}{1 + \exp(-[w_{11}^1 x_1^i + w_{12}^1 x_2^i])} \right)^2$$

# MLP - Data Perspective: A Simple Example



Layer 0  
(Input)

Layer 1  
(Output)

$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$

- Aim: To minimize the total error (or loss), which is

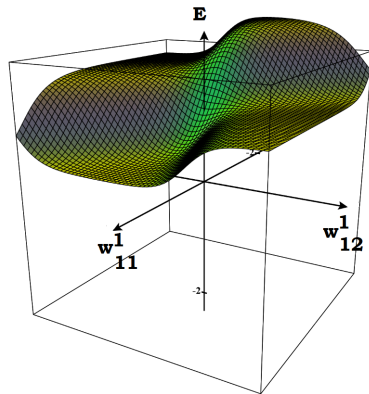
$$\min_{w_{11}^1, w_{12}^1} E = \sum_{i=1}^4 \left( y^i - \frac{1}{1 + \exp(-[w_{11}^1 x_1^i + w_{12}^1 x_2^i])} \right)^2$$



# MLP - Data Perspective: A Simple Example

## Visualizing the loss surface:

$x_1$	$x_2$	$y$	$\hat{y} = \sigma(w_{11}^1 x_1 + w_{12}^1 x_2)$
-3	-3	1	$\sigma(-3w_{11}^1 - 3w_{12}^1)$
-2	-2	1	$\sigma(-2w_{11}^1 - 2w_{12}^1)$
4	4	0	$\sigma(4w_{11}^1 + 4w_{12}^1)$
2	-5	0	$\sigma(2w_{11}^1 - 5w_{12}^1)$



$$E = \sum_{i=1}^4 \left( y^i - \frac{1}{1 + \exp(-[w_{11}^1 x_1^i + w_{12}^1 x_2^i])} \right)^2$$

# Optimization Concepts

# General Optimization Problem

$$\min_{x \in \mathcal{C}} f(x)$$

# General Optimization Problem

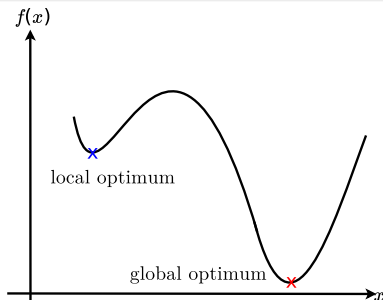
$$\min_{x \in \mathcal{C}} f(x)$$

- $f$  is called **objective function** and  $\mathcal{C}$  is called **feasible set**.
- Let  $f^* = \min_{x \in \mathcal{C}} f(x)$  denote the **optimal objective function value**.
- **Optimal Solution Set**  $S^* = \{x \in \mathcal{C} : f(x) = f^*\}$ .
- Let us denote by  $x^*$  an optimal solution in  $S^*$ .

# General Optimization Problem

$$\min_{x \in \mathcal{C}} f(x)$$

(OP)



# General Optimization Problem

$$\min_{x \in \mathcal{C}} f(x) \quad (\text{OP})$$

## Local Optimal Solution

A solution  $z$  to (OP) is called local optimal solution if  $f(z) \leq f(\hat{z})$ ,  $\forall \hat{z} \in \mathcal{N}(z, \epsilon)$  for some  $\epsilon > 0$ .

**Note:**  $\mathcal{N}(z, \epsilon)$  denotes suitable  $\epsilon$ -neighborhood of  $z$ .

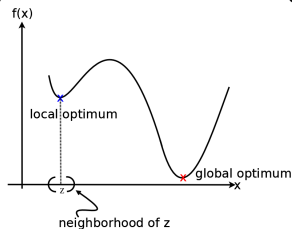
# General Optimization Problem

$$\min_{x \in \mathcal{C}} f(x) \quad (\text{OP})$$

## Local Optimal Solution

A solution  $z$  to (OP) is called local optimal solution if  $f(z) \leq f(\hat{z})$ ,  $\forall \hat{z} \in \mathcal{N}(z, \epsilon)$  for some  $\epsilon > 0$ .

**Note:**  $\mathcal{N}(z, \epsilon)$  denotes suitable  $\epsilon$ -neighborhood of  $z$ .



# General Optimization Problem

$$\min_{x \in \mathcal{C}} f(x) \quad (\text{OP})$$

## Local Optimal Solution

A solution  $z$  to (OP) is called local optimal solution if  $f(z) \leq f(\hat{z})$ ,  $\forall \hat{z} \in \mathcal{N}(z, \epsilon)$  for some  $\epsilon > 0$ .

**Note:**  $\mathcal{N}(z, \epsilon)$  denotes suitable  $\epsilon$ -neighborhood of  $z$ .

## $\epsilon$ - Neighborhood of $z \in \mathcal{C}$

$$\mathcal{N}(z, \epsilon) = \{u \in \mathcal{C} : \text{dist}(z, u) \leq \epsilon\}.$$



# General Optimization Problem

$$\min_{x \in \mathcal{C}} f(x) \quad (\text{OP})$$

## Local Optimal Solution

A solution  $z$  to (OP) is called local optimal solution if  $f(z) \leq f(\hat{z})$ ,  $\forall \hat{z} \in \mathcal{N}(z, \epsilon)$  for some  $\epsilon > 0$ .

## Global Optimal Solution

A solution  $z$  to (OP) is called global optimal solution if  $f(z) \leq f(\hat{z})$ ,  $\forall \hat{z} \in \mathcal{C}$ .