

2-1 Neural Network with One Hidden Layer I

(one training example)

Zhonglei Wang

WISE and SOE, XMU, 2025

Contents

1. Revisit logistic regression
2. Forward propagation
3. Backpropagation
4. (Batch) gradient descent algorithm

Intuition

1. Illustrate basic concepts using **only ONE** training example (\mathbf{x}, y)
2. A neural network can be viewed as a function of features \mathbf{x} with parameter $\boldsymbol{\theta}$
3. To obtain an estimator of the model parameter $\boldsymbol{\theta}$, we use a (batch) gradient descent algorithm
4. An essential step is

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \frac{\partial \mathcal{J}}{\partial \boldsymbol{\theta}} \left(\boldsymbol{\theta}^{(t)} \right)$$

- $\boldsymbol{\theta}^{(t)}$: current model parameter
- **What are the cost function and its partial derivatives?**

Intuition

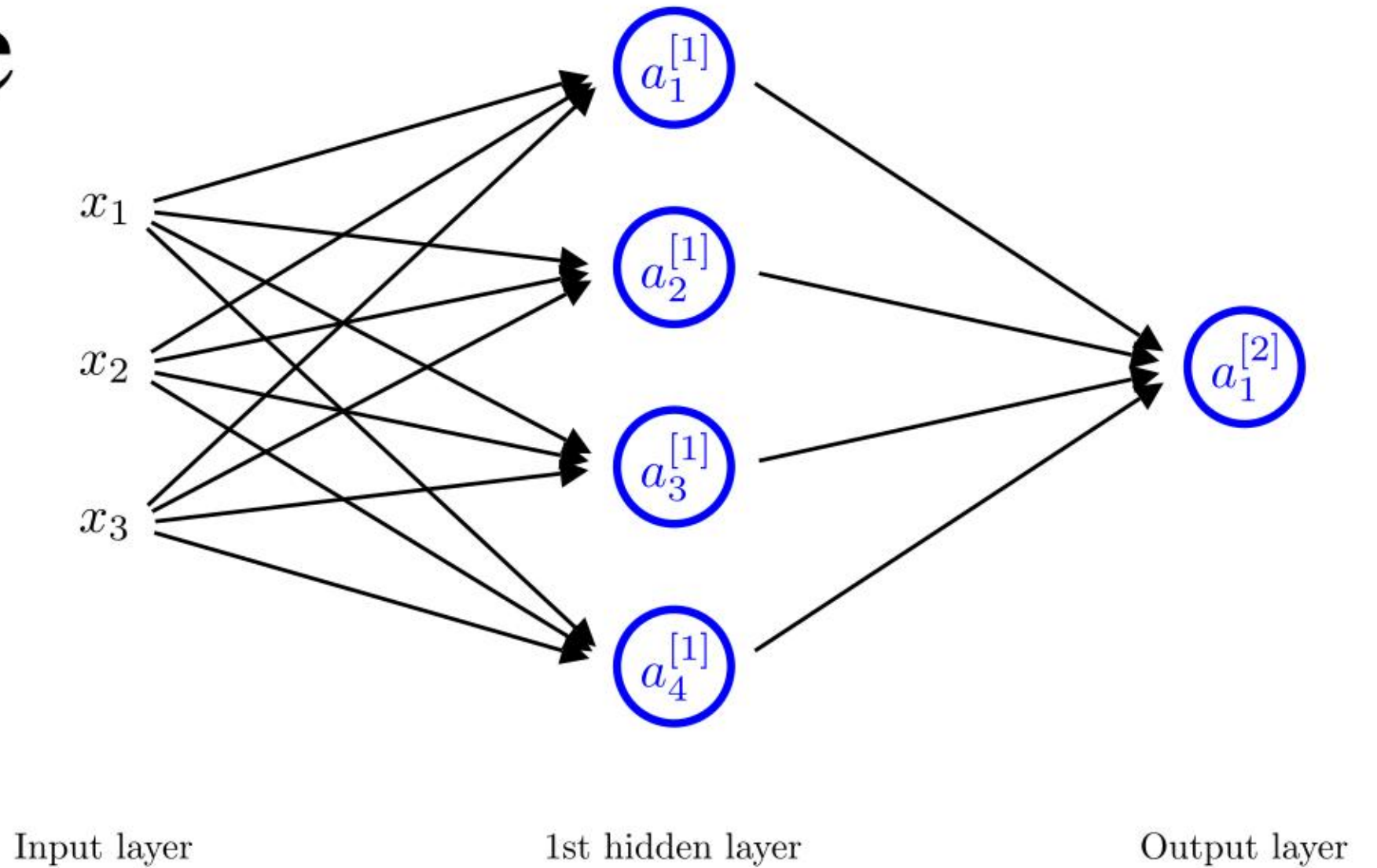
1. To obtain $\partial \mathcal{J}(\boldsymbol{\theta}^{(t)}) / \partial \boldsymbol{\theta}$, we introduce
 - **Forward propagation**: based on the current parameter, calculate the “activated” values as well as the cost function
 - **Backpropagation**: based on “those” values, obtain partial derivatives
2. We use a toy neural network to introduce forward propagation and backpropagation.

Revisit logistic regression

1. The model is too **simple** in practice
2. Why not use more “circles”?

Neural network example

1. Assume $\mathbf{x} = (x_1, x_2, x_3)^T$



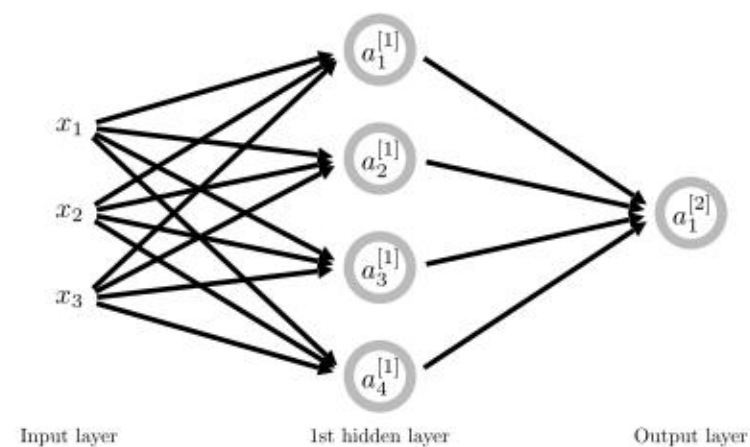
2. Each **circle** represents two operations

- **Linear** transformation with two model parameters, including a bias and a weight
- **Activation** (nonlinear transformation) without model parameters

Neural network example

1. Remarks

- Different model parameters are used for different “circles” to extract different information
- In other words, we “construct” a set of “new” features in the hidden layer
- Compared with logistic regression model, this neural network uses one hidden layer to extract more information
- For the output layer, the hidden layer can be regarded as its “input” layer



Calculation

Parameters

$$z_1^{[1]} = \mathbf{b}_1^{[1]} + \mathbf{x}^T \mathbf{w}_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$\mathbf{b}_1^{[1]} \quad \mathbf{w}_1^{[1]}$$

$$z_2^{[1]} = \mathbf{b}_2^{[1]} + \mathbf{x}^T \mathbf{w}_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$\mathbf{b}_2^{[1]} \quad \mathbf{w}_2^{[1]}$$

$$z_3^{[1]} = \mathbf{b}_3^{[1]} + \mathbf{x}^T \mathbf{w}_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$\mathbf{b}_3^{[1]} \quad \mathbf{w}_3^{[1]}$$

$$z_4^{[1]} = \mathbf{b}_4^{[1]} + \mathbf{x}^T \mathbf{w}_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$\mathbf{b}_4^{[1]} \quad \mathbf{w}_4^{[1]}$$

$$z_1^{[2]} = \mathbf{b}_1^{[2]} + (\mathbf{a}^{[1]})^T \mathbf{w}_1^{[2]}, \quad a_1^{[2]} = \sigma(z_1^{[2]})$$

$$\mathbf{b}_1^{[2]} \quad \mathbf{w}_1^{[2]}$$

Vectorization

1. Denote

- L : number of layers in the neural network
- $d^{[l]}$: number of neurons in the l th layer ($l = 0, \dots, L$)
- $\mathbf{a}^{[l]} = (a_1^{[l]}, \dots, a_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times 1}$
- $\mathbf{W}^{[l]} = (\mathbf{w}_1^{[l]}, \dots, \mathbf{w}_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times d^{[l-1]}}$
- $\mathbf{b}^{[l]} = (b_1^{[l]}, \dots, b_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times 1}$

2. Model parameters

- $\{(\mathbf{b}^{[l]}, \mathbf{W}^{[l]}) : l = 1, \dots, L\}$

Vectorization

1. For the previous example, we have

- $d^{[0]} = 3, d^{[1]} = 4, d^{[2]} = 1$
- $\mathbf{W}^{[1]} \in \mathbb{R}^{4 \times 3}, \mathbf{W}^{[2]} \in \mathbb{R}^{1 \times 4}$
- $\mathbf{b}^{[1]} \in \mathbb{R}^{4 \times 1}, \mathbf{b}^{[2]} \in \mathbb{R}^{1 \times 1}$

Forward propagation

1. Forward propagation obtains “activated” values using the CURRENT parameters
2. Assume the **current model parameters** are $\mathbf{b}^{[1]}, \mathbf{W}^{[1]}, b^{[2]}, \mathbf{W}^{[2]}$
3. The calculation can be simplified as

$$\mathbf{z}^{[1]} = \mathbf{b}^{[1]} + \mathbf{W}^{[1]} \mathbf{x}$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{b}^{[2]} + \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma(\mathbf{z}^{[2]})$$

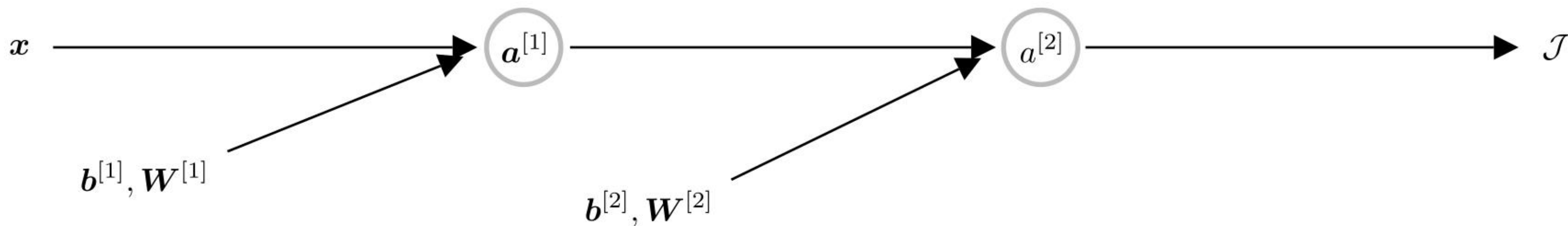
Forward propagation

1. Consider a binary response: $y \in \{0, 1\}$
2. $a^{[2]}$ is the estimated value for $P(y = 1 \mid \mathbf{x})$
3. Thus, we consider a loss function

$$\mathcal{J} = \mathcal{L} = - \left\{ y \log a^{[2]} + (1 - y) \log (1 - a^{[2]}) \right\}$$

Forward propagation

1. Visualize the calculation details



Forward propagation

$$\mathbf{z}^{[1]} = \mathbf{b}^{[1]} + \mathbf{W}^{[1]} \mathbf{x}$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{b}^{[2]} + \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma(\mathbf{z}^{[2]})$$

$$\mathcal{J} = - \left\{ y \log a^{[2]} + (1 - y) \log (1 - a^{[2]}) \right\}$$

Backpropagation

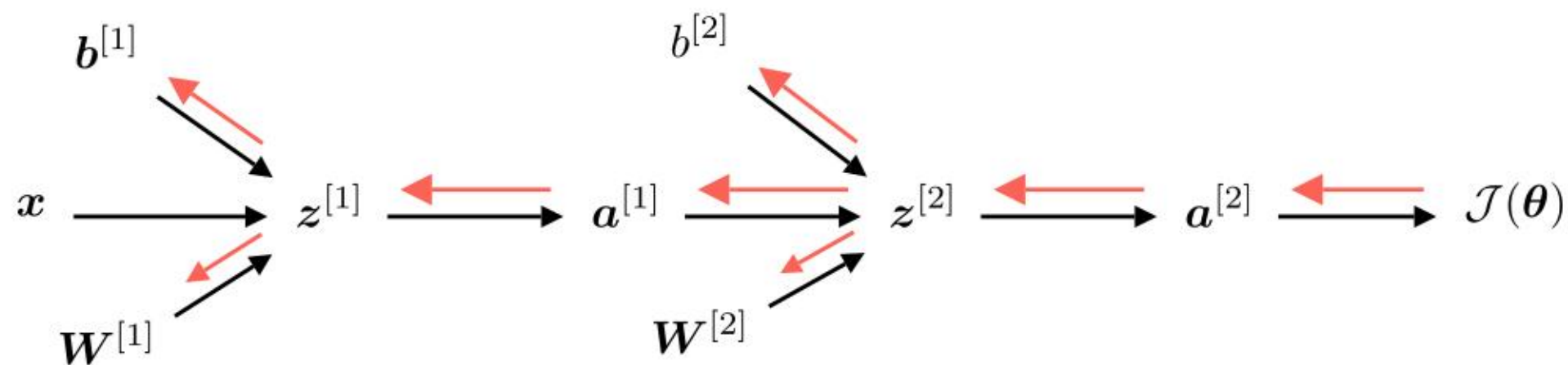
1. Based on the current model parameters, backpropagation is to obtain

$$\frac{\partial \mathcal{J}}{\partial \mathbf{b}^{[l]}}, \quad \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{[l]}} \quad (l = 1, \dots, L)$$

2. Core technique: **chain rule**

- Denote $f(\mathbf{x}) = g \circ h(\mathbf{x})$
- Then, we have

$$\frac{\partial f}{\partial \mathbf{x}} = \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial \mathbf{x}}$$



$$\frac{\partial \mathcal{J}(\theta)}{\partial b^{[2]}} = a^{[2]} - y \quad \frac{\partial \mathcal{J}(\theta)}{\partial W^{[2]}} = (a^{[2]} - y)(\mathbf{a}^{[1]})^T$$

$$\frac{\partial \mathcal{J}(\theta)}{\partial \mathbf{b}^{[1]}} = (a^{[2]} - y) \mathbf{D}(\mathbf{W}^{[2]})^T \quad \frac{\partial \mathcal{J}(\theta)}{\partial \mathbf{W}^{[1]}} = (a^{[2]} - y) \mathbf{D}(\mathbf{W}^{[2]})^T \mathbf{x}^T$$

$$\mathbf{D} = \text{diag}(\{a_j^{[1]}(1 - a_j^{[1]}) : j = 1, \dots, d^{[1]}\})$$

We only need to cache $\mathbf{a}^{[1]}$ and $a^{[2]}$

(Batch) gradient descent algorithm

Step1. Randomly initialize $\boldsymbol{\theta}^{(0)}$

Step2. Based on $\boldsymbol{\theta}^{(t)}$ obtain

$$\nabla \mathcal{J} \left(\boldsymbol{\theta}^{(t)} \right) = \frac{\partial \mathcal{J}}{\partial \boldsymbol{\theta}} \left(\boldsymbol{\theta}^{(t)} \right)$$

Step3. Update parameter

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \nabla \mathcal{J} \left(\boldsymbol{\theta}^{(t)} \right)$$

Step4. Go back to Step 2 until convergence