

2-7 Softmax Regression

Zhonglei Wang

WISE and SOE, XMU, 2025

Contents

1. Model

2. Forward propagation

3. Backpropagation

4. Example

Introduction

1. Binary classification models

- $y \in \{0, 1\}$
- Model is built for $P(y = 1 \mid \mathbf{x})$

2. In practice, we may have K classes: $y \in \{0, 1, \dots, K - 1\}$

- Consider one-hot representation
- $\mathbf{y} = (0, \dots, 1, \dots, 0)^T$
- If $y = k$, only the $(k + 1)$ th element of \mathbf{y} is 1

Model

1. For a feature \mathbf{x}

- We consider a neural network with K outputs
- Each output is a score for the corresponding class
- The scores may not sum up to 1

2. That is, the only difference is the number of neurons in the last layer

- For **binary** classification problems, we only have **one neuron** for the output layer
- For **general** classification problems, we have $d^{[L]} = K$ **neurons** for the output layer

Model

1. Recall that

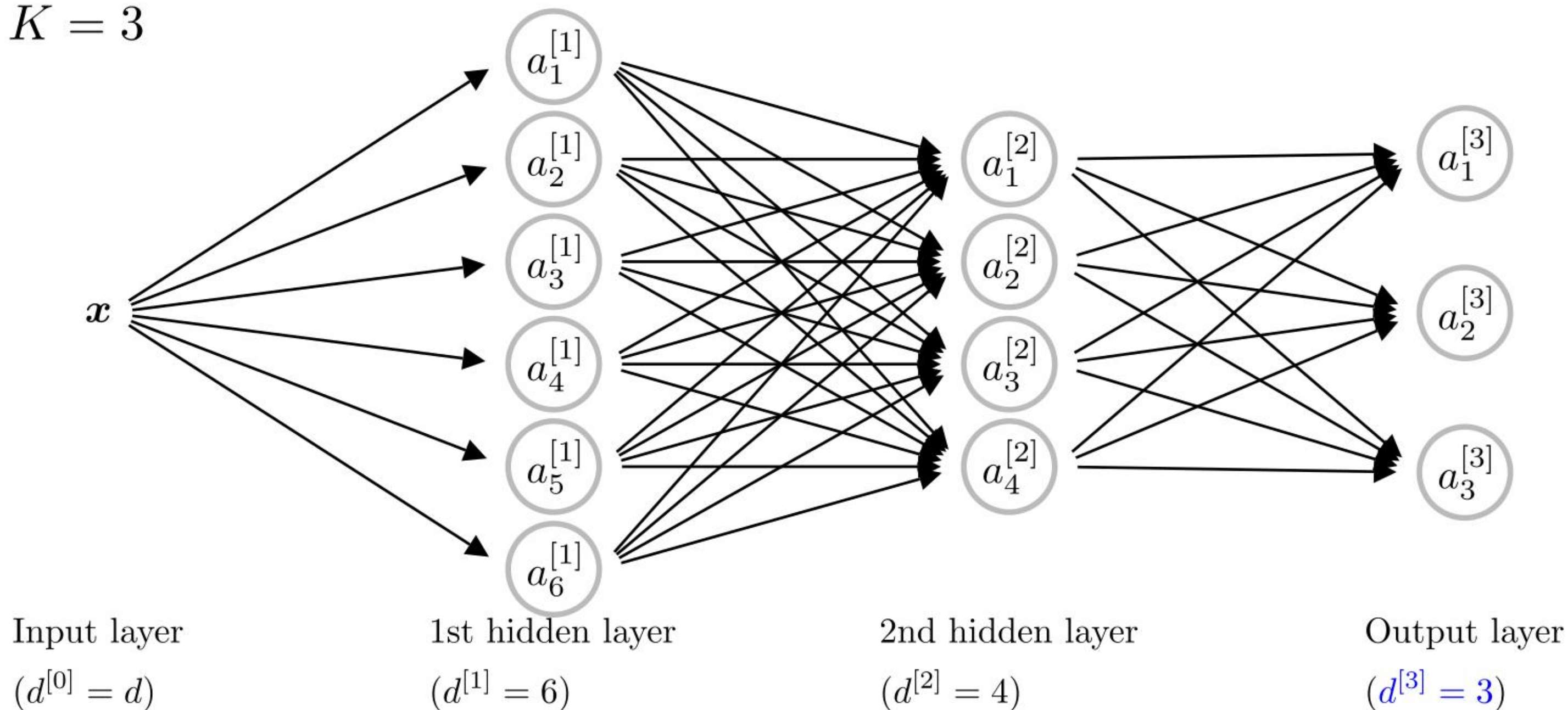
- L : number of layers in the neural network
- $d^{[l]}$: number of neurons in the l th layer ($l = 0, \dots, L$)
- $\mathbf{a}^{[l]} = (a_1^{[l]}, \dots, a_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]} \times 1}$
- $\mathbf{W}^{[l]} = (\mathbf{w}_1^{[l]}, \dots, \mathbf{w}_{d^{[l]}}^{[l]})^T \in \mathbb{R}^{d^{[l]}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

2. For **binary classification** problems, we have $d^{[L]} = 1$

3. For **general softmax regression** problems, we have $d^{[L]} = K$

Model

1. For $K = 3$



Forward propagation

1. Let $\mathbf{A}^{[0]} = \mathbf{X}$
2. For $l = 1, \dots, L$,

$$\mathbf{Z}^{[l]} = \left(\mathbf{b}^{[l]}\right)^{\mathrm{T}} + \mathbf{A}^{[l-1]} \left(\mathbf{W}^{[l]}\right)^{\mathrm{T}}$$
$$\mathbf{A}^{[l]} = \sigma^{[l]} \left(\mathbf{Z}^{[l]}\right)$$

- $\sigma^{[l]}(z)$: activation function for the l th layer
- Broadcasting is used for activation functions

3. For the last layer,

$$\mathbf{W}^{[L]} \in \mathbb{R}^{K \times d_L - 1}, \quad \mathbf{b} \in \mathbb{R}^{K \times 1}$$

Forward propagation

1. The cost function for softmax regression is

$$\mathcal{J}(\boldsymbol{\theta}) = -n^{-1} \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log a_{ik}^{[L]}$$

- $\boldsymbol{\theta}$: model parameters
- $\mathbf{y}_i = (y_{i1}, \dots, y_{iK})^T$: one-hot representation for the i th example
- $a_{ik}^{[L]}$: estimated probability for the k th class of the i th example

2. The dimension of the $\mathbf{A}^{[L]}$, containing estimated probabilities in the last layer

- $\mathbf{A}^{[L]} \in \mathbb{R}^{n \times 1}$ for **binary classification** problems
- $\mathbf{A}^{[L]} \in \mathbb{R}^{n \times K}$ for **softmax regression** problems

Backpropagation

1. $\mathbf{dA}^{[L]}$ can be obtained from the cost function

2. Assume $\mathbf{dA}^{[l]}$ is available ($l = L, \dots, 2$)

$$\mathbf{dZ}^{[l]} = \mathbf{dA}^{[l]} \circ \sigma^{[l]'}(\mathbf{Z}^{[l]})$$

$$\mathbf{dW}^{[l]} = \left(\mathbf{dZ}^{[l]}\right)^T \mathbf{dA}^{[l-1]}$$

$$\mathbf{db}^{[l]} = \left(\mathbf{dZ}^{[l]}\right)^T \mathbf{1}$$

$$\mathbf{dA}^{[l-1]} = \mathbf{dZ}^{[l]} \mathbf{W}^{[l]}$$

3. It remains to obtain $\mathbf{dA}^{[L]}$ for softmax regression

Backpropogation

1. The cost function is

$$\mathcal{J}(\boldsymbol{\theta}) = -n^{-1} \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log a_{ik}^{[L]}$$

2. Thus, the ik th component of $d\mathbf{A}^{[L]}$ is

$$\frac{\partial \mathcal{J}}{\partial a_{ik}^{[L]}} = -\frac{y_{ik}}{a_{ik}^{[L]}}$$

3. Done!

Example

Test image



Model (FNN with 2 hidden layers)

Estimated result

