

# Lec12 Deep Learning

## 12-2 Multilayer Perceptron

Yang Shu

School of Data Science and Engineering

yshu@dase.ecnu.edu.cn

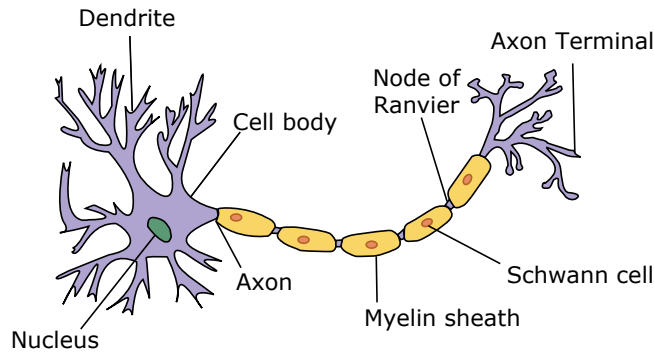
[Acknowledgement: Slides are adapted from Deep Learning Course, Mingsheng Long, THU]



# Outline

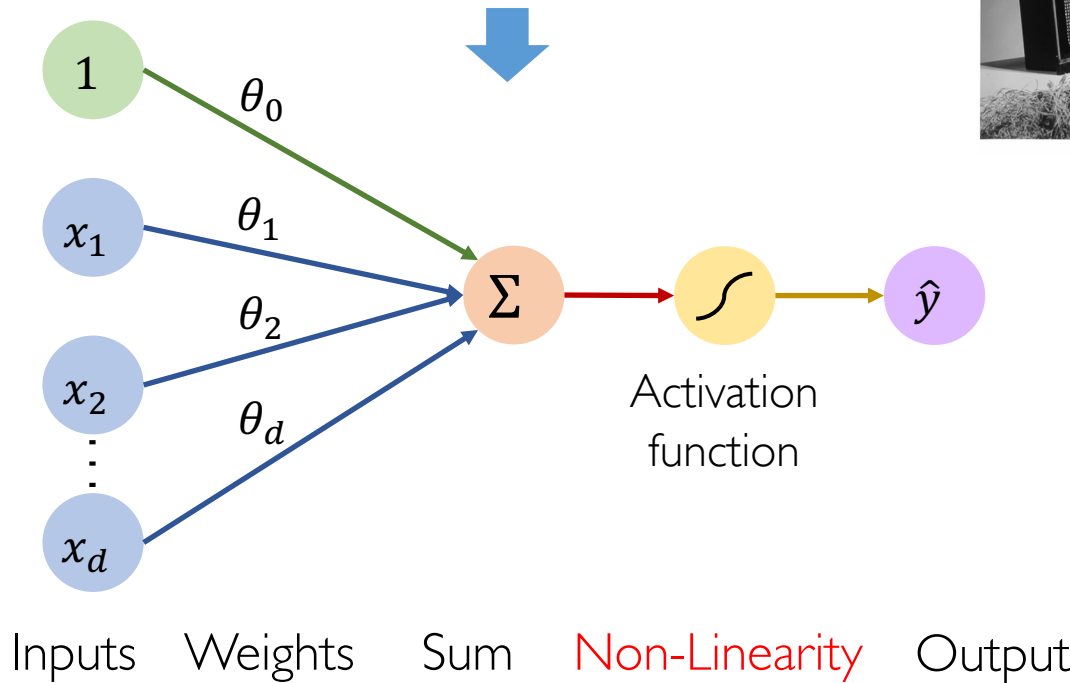
- **Feedforward Network**
  - **Perceptron**
  - Multilayer Perceptron
- Backpropagation
- Practical Training Strategies

# Neuron and Perceptron



Rosenblatt 1958  
An psychologist

Principles of Neurodynamics:  
Perceptrons and the Theory of  
Brain Mechanisms



Linear combination of inputs

$$\hat{y} = g \left( \theta_0 + \sum_{i=1}^d x_i \theta_i \right)$$

Output

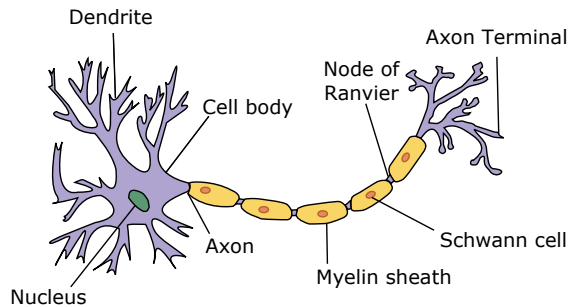
Non-Linearity Activation function

Bias

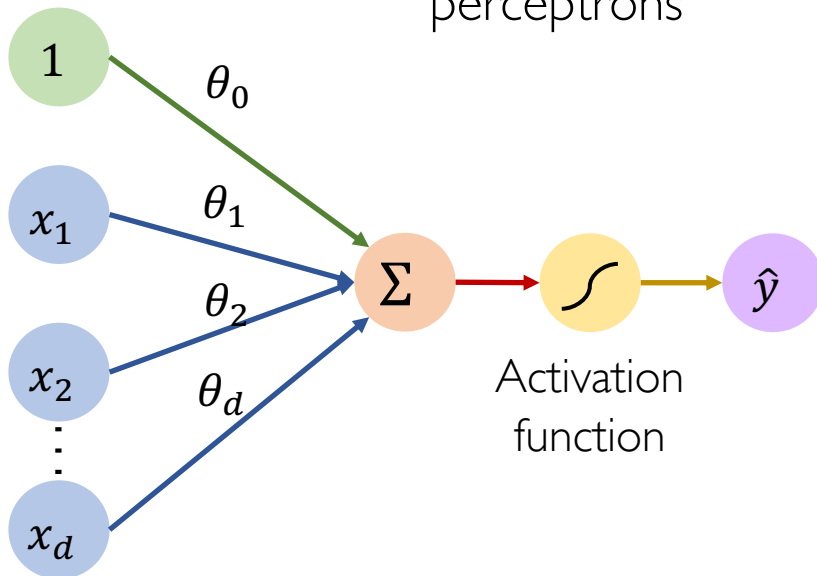
# Outline

- **Feedforward Network**
  - Perceptron
  - **Multilayer Perceptron**
- Backpropagation
- Practical Training Strategies

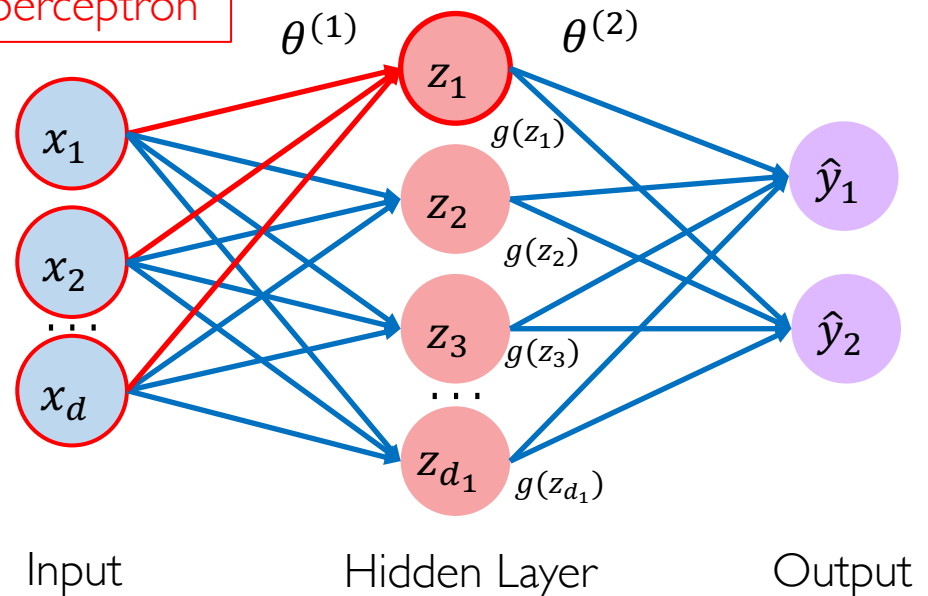
# Multilayer Perceptron (MLP)



MLP is layering of perceptrons

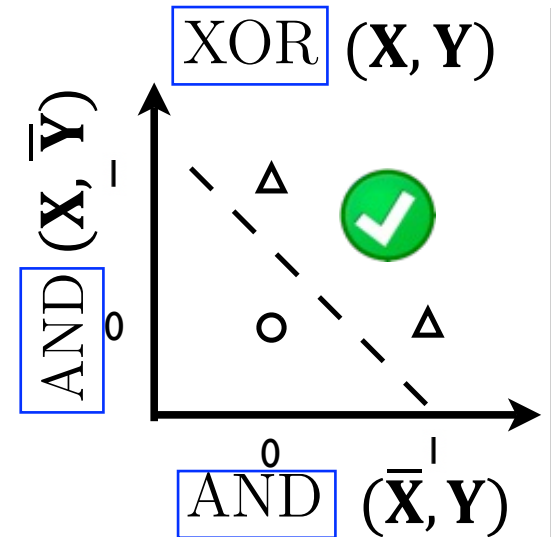
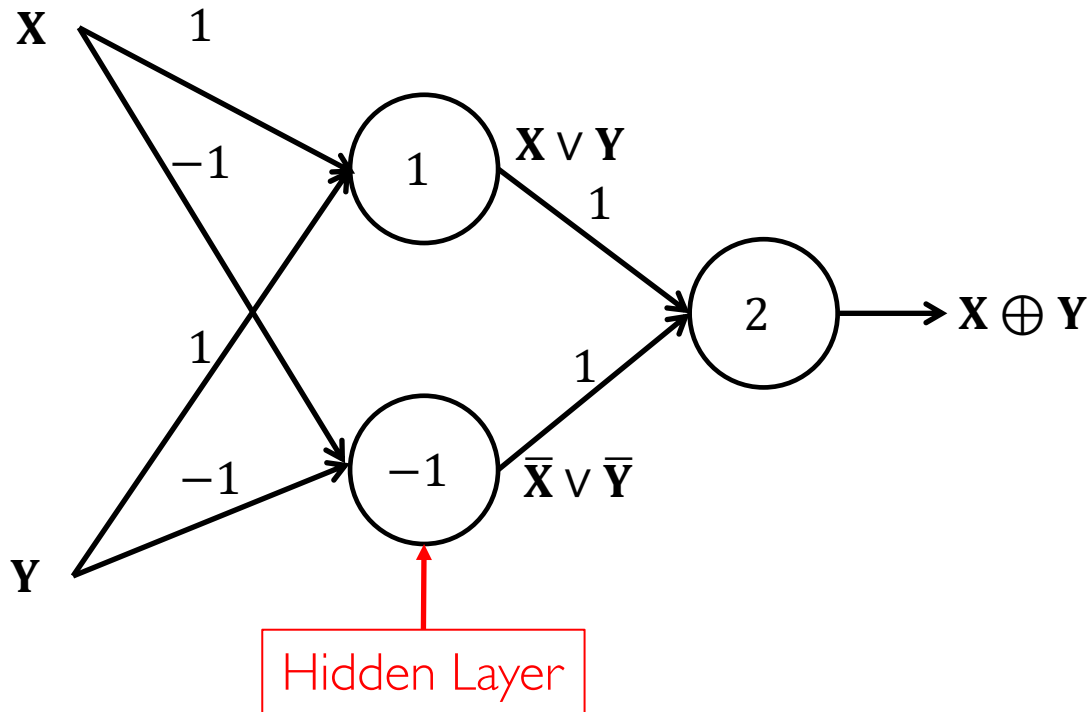


A single perceptron



# MLP for XOR

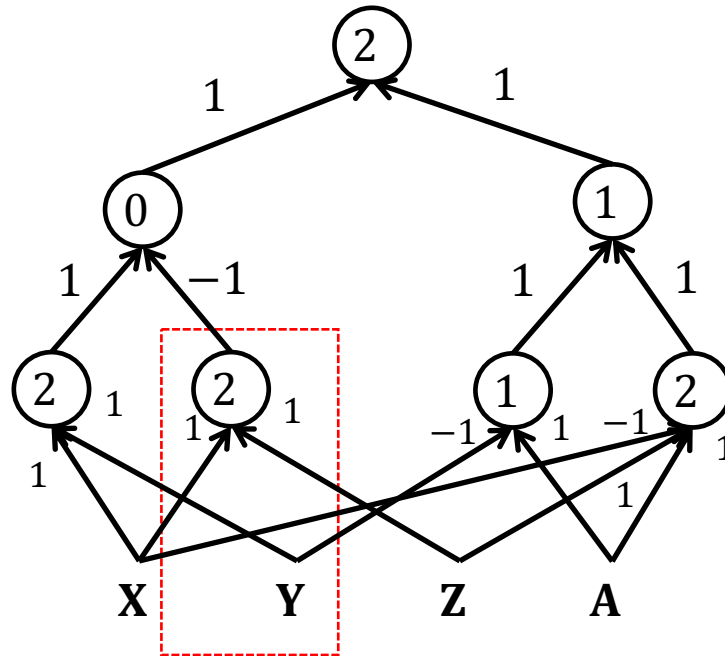
$$X \oplus Y = (X \vee Y) \& (\bar{X} \vee \bar{Y})$$



Marvin Minsky and Seymour Papert. Perceptrons. An Introduction to Computational Geometry. 1969.

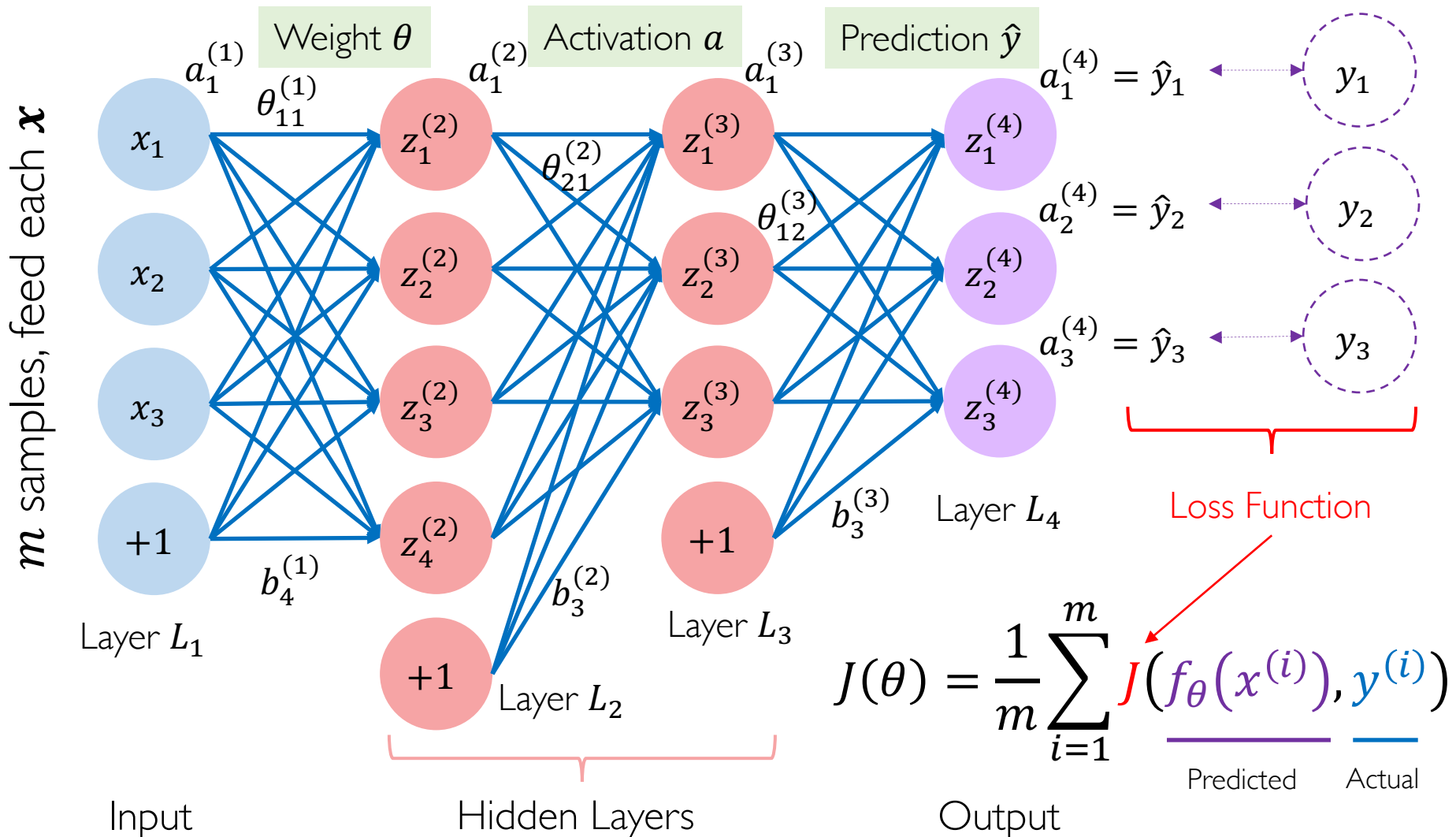
# MLP: Boolean Functions

$$((A \& \bar{X} \& Z) | (A \& \bar{Y})) \& ((X \& Y) | \overline{(X \& Z)})$$



- MLP is universal to **represent** arbitrarily complex Boolean functions
- The connections in MLP can be **sparse** - a phenomenon in the **Brain**

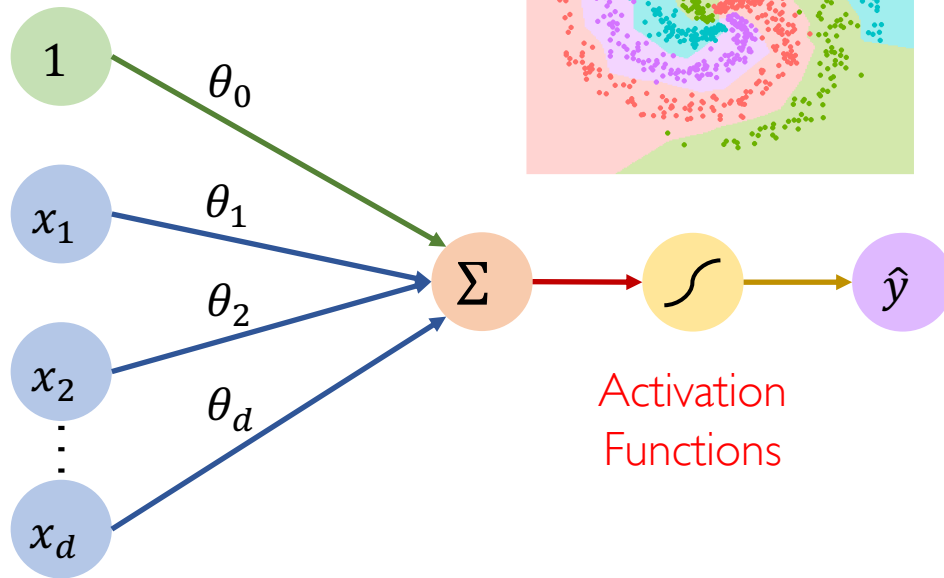
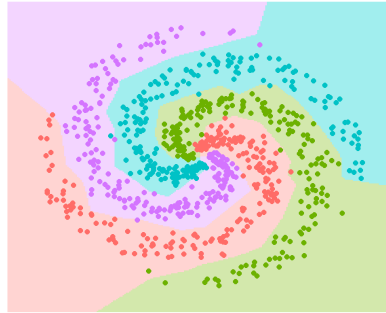
# Multilayer Perceptron (MLP)





# Activation Functions

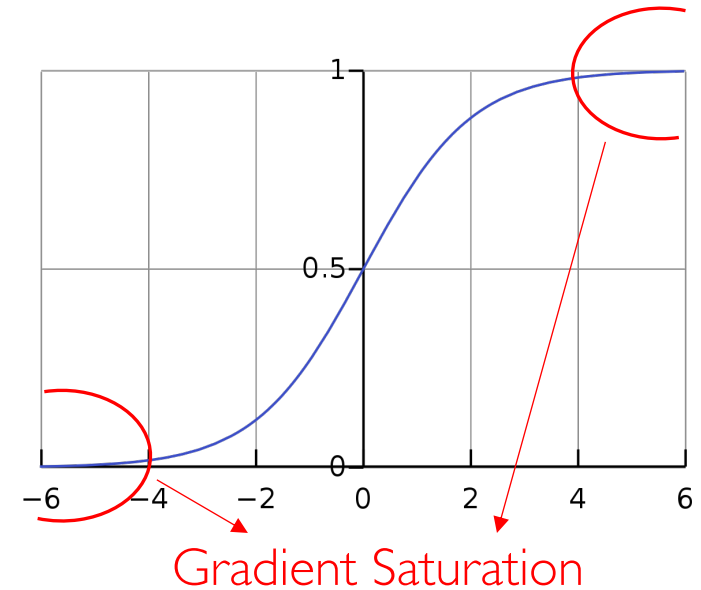
Non-linearity through activation functions  $\hat{y} = g(\theta_0 + \mathbf{x}^T \boldsymbol{\theta})$



Inputs    Weights    Sum    Non-Linearity    Output

sigmoid function

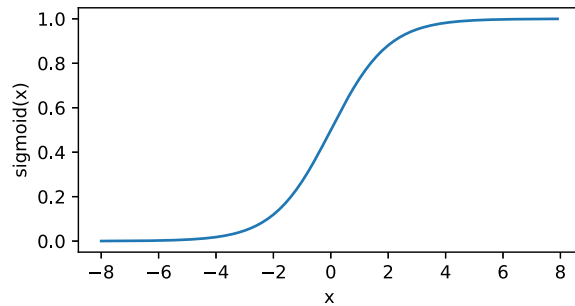
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



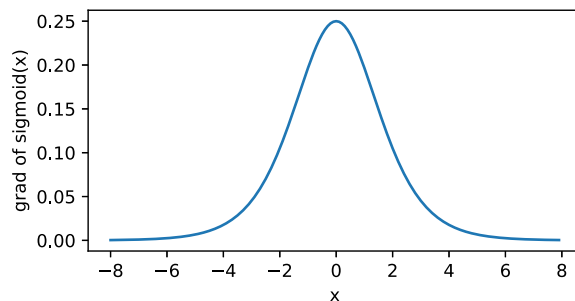
# Activation Functions

Sigmoid

$$g(z) = \frac{1}{1 + e^{-z}}$$



$$g'(z) = g(z)(1 - g(z))$$

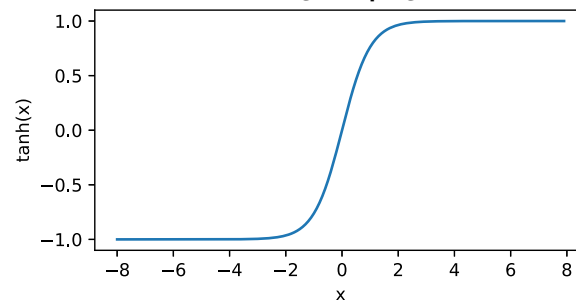


Range in (0, 1), probability

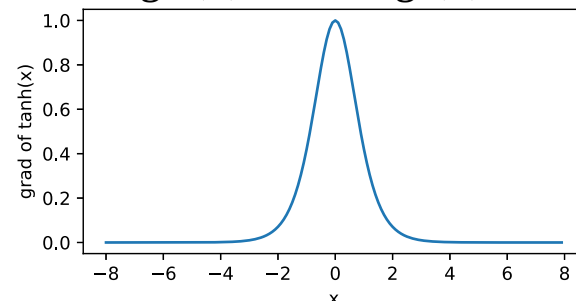
Relative smaller gradient

Hyperbolic Tangent (tanh)

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



$$g'(z) = 1 - g(z)^2$$

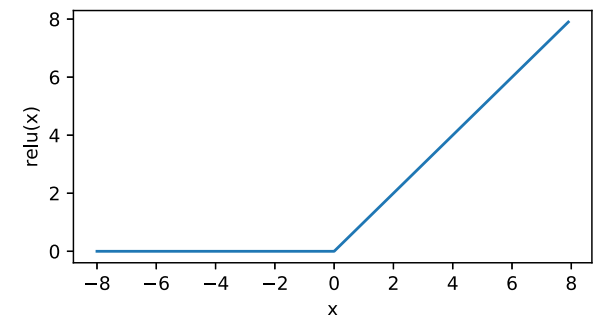


Zero-centered

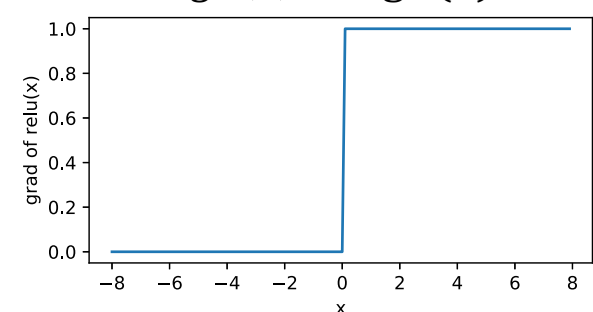
Relatively larger saturation region

Rectified Linear Unit (ReLU)

$$g(z) = \max(0, z)$$



$$g'(z) = \text{sgn}(z)$$



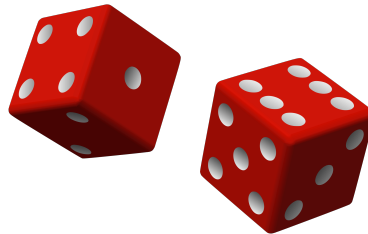
Time-efficient and faster  
convergence, but neurons  
are prone to death.

# Activation Functions

- Softmax function

$$g(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Categorical  
distribution



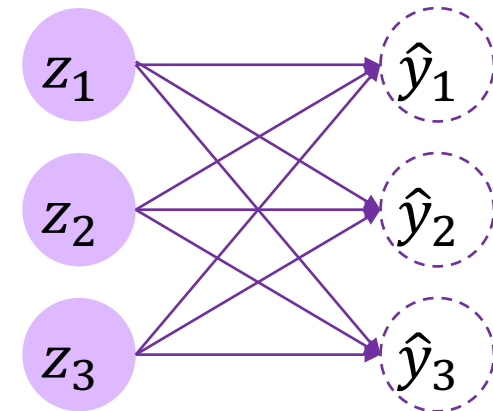
rooster  
cat  
dog  
donkey

- Winner takes all
- The biggest one dominates the outputs

- Exponent arithmetic

- Cause numerical overflow
- For numerical stability, implementation is

$$g(\mathbf{z})_i = \frac{e^{z_i - z_u}}{\sum_{j=1}^k e^{z_j - z_u}}, u = \operatorname{argmax}(z_j)$$



$$z_j \in (-\infty, +\infty) \quad \hat{y}_j \in [0,1], \sum_{j=1}^k \hat{y}_j = 1$$

# Cost Function

- Softmax function in the output layer

$$\hat{\mathbf{y}} = \mathbf{a}^{(n_l)} = f_{\theta}(\mathbf{x}^{(i)}) = \begin{bmatrix} p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ p(y^{(i)} = 2 | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ \vdots \\ p(y^{(i)} = k | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \end{bmatrix} = \frac{1}{\sum_{j=1}^k \exp(z_j^{(n_l)})} \begin{bmatrix} \exp(z_1^{(n_l)}) \\ \exp(z_2^{(n_l)}) \\ \vdots \\ \exp(z_k^{(n_l)}) \end{bmatrix}$$

- Cross-entropy  $J(q, p) = -\sum_{j=1}^k q_j \log p_j$

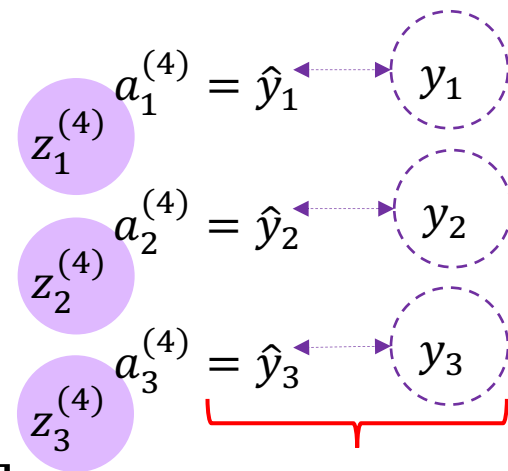
- Loss function:  $J(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^k y_j \log \hat{y}_j$

-  $\mathbf{y}$  follows one-hot coding

»  $y_j = 1$  if  $y^{(i)} = j$  and  $y_j = 0$  otherwise

- Cost function

$$\min_{\theta} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[ \sum_{j=1}^k \mathbf{1}\{y^{(i)} = j\} \log \frac{\exp(z_j^{(n_l)})}{\sum_{j'=1}^k \exp(z_{j'}^{(n_l)})} \right]$$



Loss Function

# Outline

- Feedforward Network
  - Perceptron
  - Multilayer Perceptron
- **Backpropagation**
- Practical Training Strategies

# Multilayer Perceptron: Notations

- Recursive transformations in propagation

$$- a_i^{(1)} = x_i \quad \text{\#units}$$

$$- z_i^{(l)} = \sum_{j=1}^{s_{l-1}} \theta_{ij}^{(l-1)} a_j^{(l-1)} + b_i^{(l-1)}$$

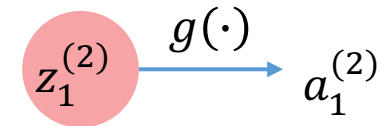
$$- a_i^{(l)} = g(z_i^{(l)})$$

$$- \hat{y}_i = a_i^{(n_l)}$$

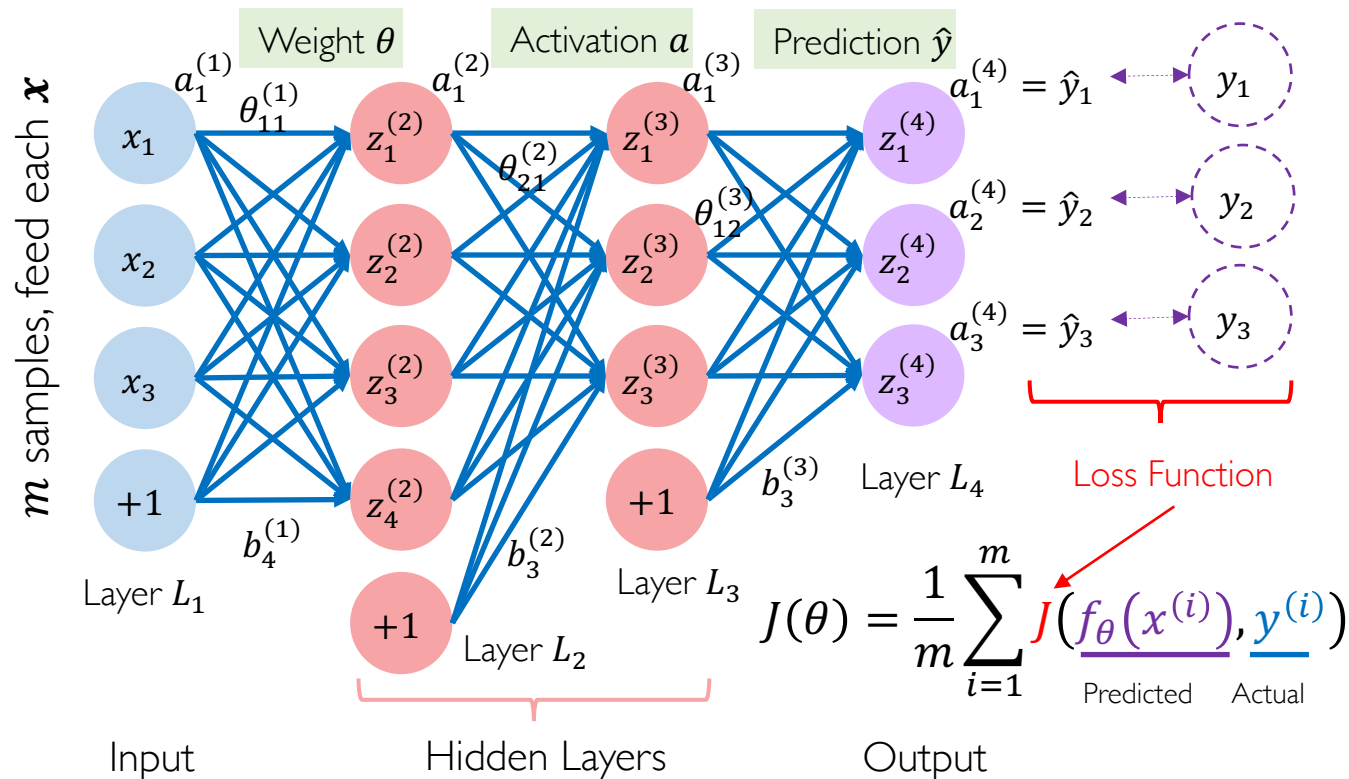
Superscript:  
layer index

\#layers

Subscript:  
unit index



Activation function (element-wise)



# Gradient-Based Training

$$\arg \min_{\theta} O(\mathcal{D}; \theta) = \sum_{i=1}^m L(y_i, f(x_i); \theta) + \Omega(\theta)$$

Risk  
Minimization

Data

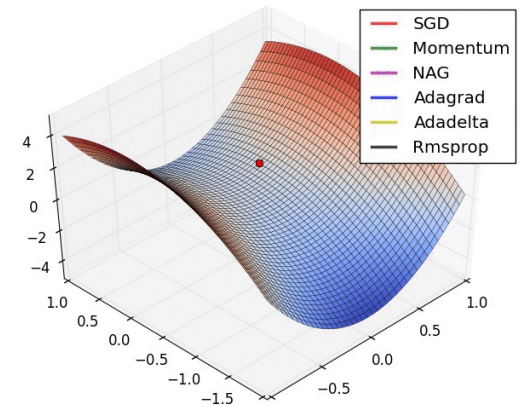
Hypothesis

Parameters

- Iterative Algorithm (Convergence Guarantee)

```
for ( $t = 1$  to  $T$ ) {  
  1. ForwardPropagation()  
  2. BackwardPropagation()  
  3.  $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla_{\theta} O(\mathcal{D}; \theta^{(t)})$   
}
```

Parameter Updates



# Gradient Descent (GD)

$$J(\theta, b) = \left[ \frac{1}{m} \sum_{i=1}^m J(\theta, b; x^{(i)}, y^{(i)}) \right]$$

- Gradient descent

$$- \theta_{ij}^{(l)} = \theta_{ij}^{(l)} - \eta \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b)$$

$$- b_i^{(l)} = b_i^{(l)} - \eta \frac{\partial}{\partial b_i^{(l)}} J(\theta, b)$$

$$- \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b; x^{(i)}, y^{(i)})$$

$$- \frac{\partial}{\partial b_i^{(l)}} J(\theta, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b_i^{(l)}} J(\theta, b; x^{(i)}, y^{(i)})$$

Derivatives of last layer

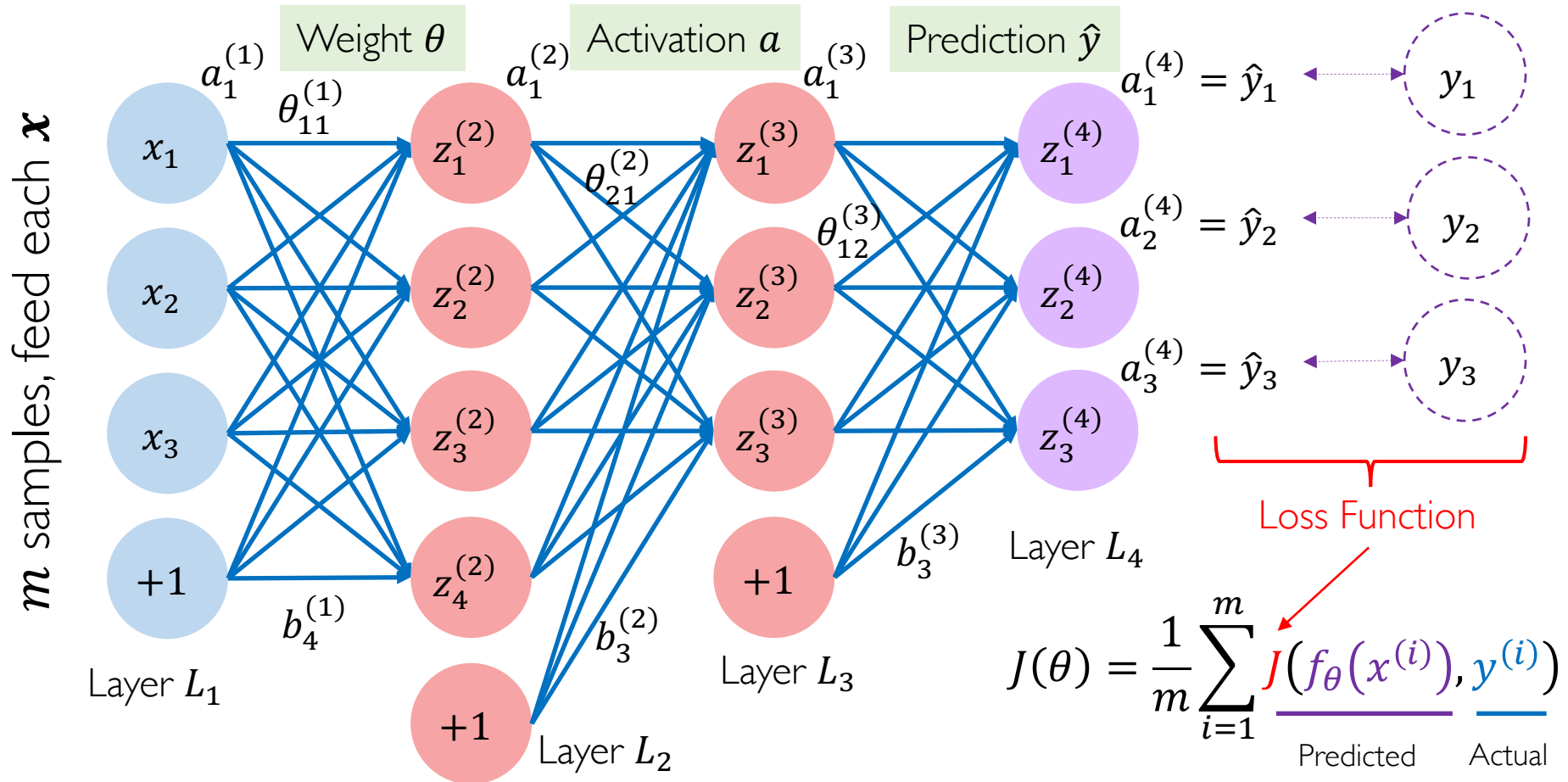
$$J(\theta, b; x^{(i)}, y^{(i)}) = - \sum_{j=1}^k \mathbf{1}\{y^{(i)} = j\} \log \frac{\exp(z_j^{(n_l)})}{\sum_{j'=1}^k \exp(z_{j'}^{(n_l)})}$$
$$\frac{\partial J(\theta, b)}{\partial z_j^{(n_l)}} = - \left( \mathbf{1}\{y^{(i)} = j\} - p(y^{(i)} = j | x^{(i)}; \theta) \right)$$

How to compute the derivatives  
for parameters in hidden layers?



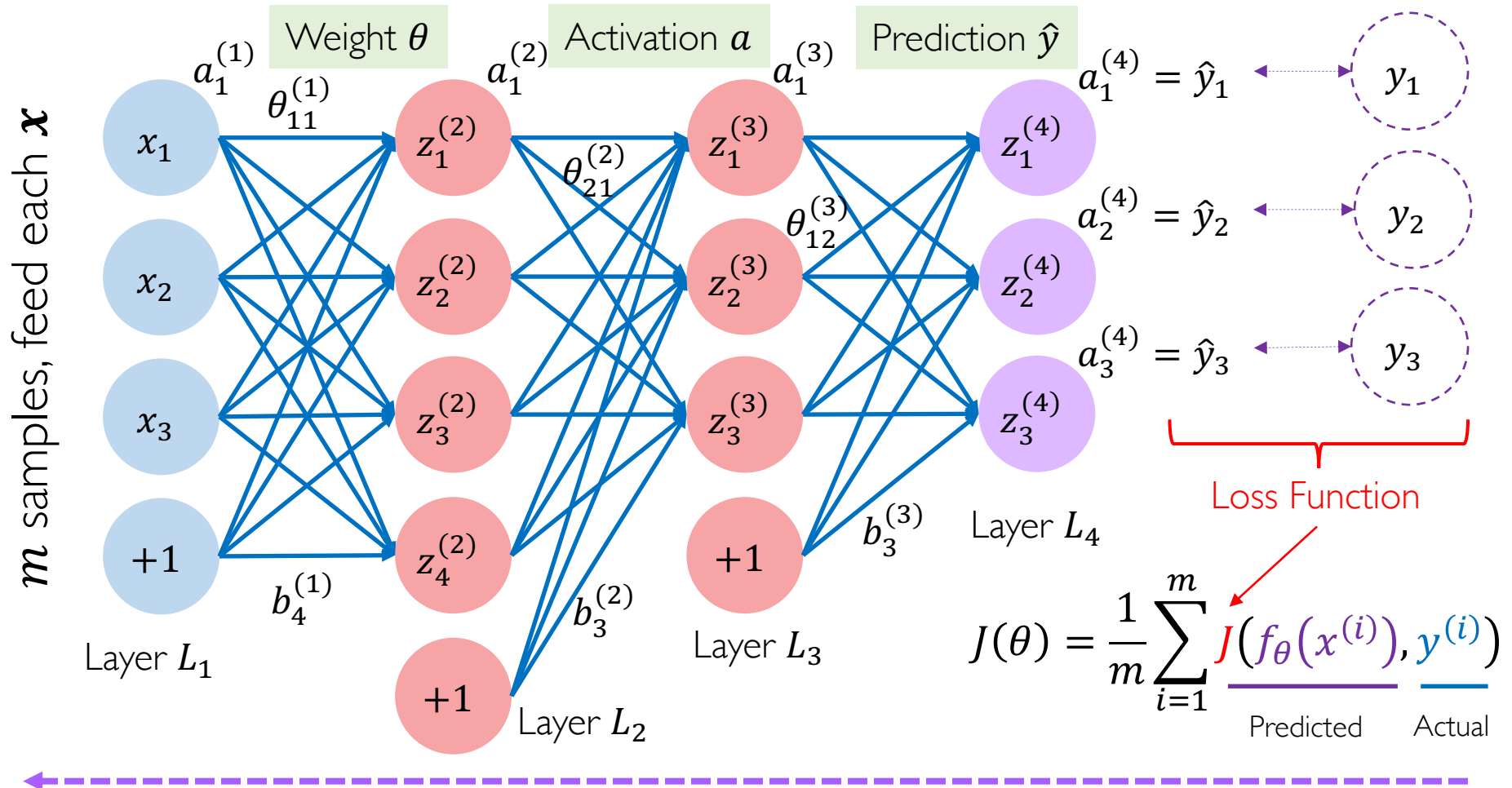


# Step I: Forward Propagation



- Forward propagate all  $(\mathbf{x}, \mathbf{y})$ 's to compute activations and objective  $J(\theta, \mathbf{b})$

# Step 2: Backward Propagation



- Backward propagate the gradients  $\nabla J(\theta, b)$  to update parameters in all layers

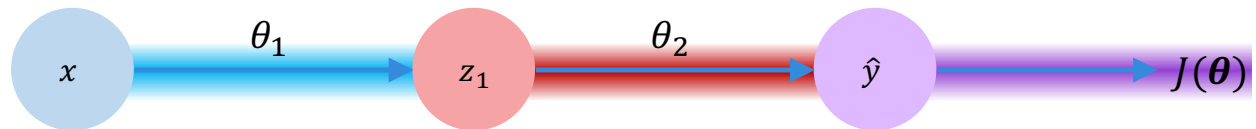
# Backpropagation (BP) Algorithm

- The backpropagation algorithm applies the chain rule of calculus to the parameters  $\theta$  of each layer

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

$$\nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$$

Jacobian matrix



$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{\partial J(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial \theta_1}$$

- Backpropagation is an efficient implementation of the chain rule
  - Uses **dynamic programming** (table filling)
  - Avoids re-computing repeated subexpressions (in dependency)
  - Speed vs memory tradeoff (sensitive to #samples  $m$ )

# Computing the Residual

- For each (node  $i$  in layer  $l$ ), compute the residual  $\delta_i^{(l)}$

$$\delta_i^{(l)} \triangleq \frac{\partial}{\partial z_i^{(l)}} J(\theta, b; x, y)$$

- $\delta_i^{(l)}$  measures how much each (node  $i$  in layer  $l$ ) is **responsible** for any errors  $J(\theta, b)$  of the network output  $z_i^{(n_l)}$
- Among all the nodes, output nodes are **directly** related to the loss, and we firstly compute  $\delta_i^{(n_l)}$  for each **output node**  $i$  in layer  $n_l$

$$\delta_i^{(n_l)} = \frac{\partial}{\partial z_i^{(n_l)}} J(\theta, b; x, y) = \frac{\partial}{\partial \hat{y}_i} J(\theta, b; x, y) g'(z_i^{(n_l)})$$

- How about hidden nodes (**not directly** related to the loss)?

# Computing the Residual

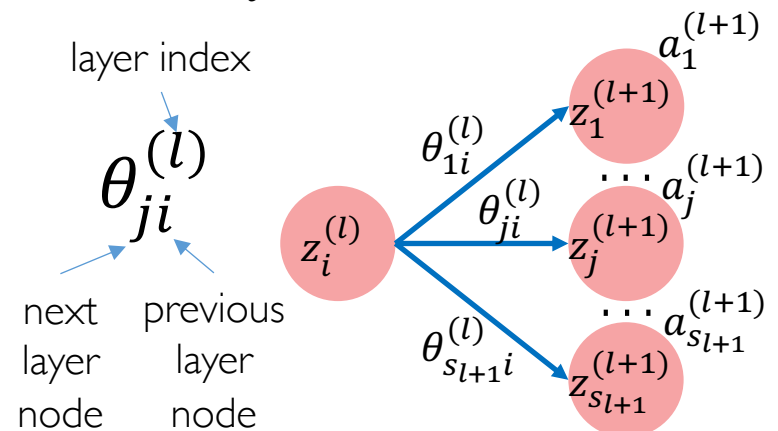
- For each (node  $i$  in layer  $l$ ), compute the residual  $\delta_i^{(l)}$

$$\delta_i^{(l)} \triangleq \frac{\partial}{\partial z_i^{(l)}} J(\theta, b; x, y)$$

- Note that  $z_j^{(l+1)} = \sum_{i=1}^{s_l} \theta_{ji}^{(l)} a_i^{(l)} + b_j^{(l)}$
- For each (hidden node  $i$  in layer  $l$ ), compute  $\delta_i^{(l)}$  by a **weighted average** of the residuals of the nodes taking  $a_i^{(l)}$  as input

$$\frac{\partial J(\theta, b; x, y)}{\partial z_i^{(l)}} = \sum_{j=1}^{s_{l+1}} \frac{\partial J(\theta, b; x, y)}{\partial z_j^{(l+1)}} \frac{\partial z_j^{(l+1)}}{\partial z_i^{(l)}}$$

$$\delta_i^{(l)} = \sum_{j=1}^{s_{l+1}} \delta_j^{(l+1)} \theta_{ji}^{(l)} g'(z_i^{(l)})$$



# Computing the Gradients

- For (each node  $i$  in layer  $l$ ), we have computed
  - Forward propagation:  $z_j^{(l+1)} = \sum_{i=1}^{s_l} \theta_{ji}^{(l)} a_i^{(l)} + b_j^{(l)}$ 
    - » Non-linearity  $a_j^{(l)} = g(z_j^{(l)})$
  - Backward propagation:  $\delta_i^{(l)} \triangleq \frac{\partial}{\partial z_i^{(l)}} J(\theta, b; x, y)$
- Applying the chain rule, the gradients of parameters  $\theta_{ij}^{(l)}$  and  $b_i^{(l)}$ :

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b; x, y) = \frac{\partial}{\partial z_i^{(l+1)}} J(\theta, b; x, y) \frac{\partial z_i^{(l+1)}}{\partial \theta_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l+1)}$$

$$\frac{\partial}{\partial b_i^{(l)}} J(\theta, b; x, y) = \frac{\partial}{\partial z_i^{(l+1)}} J(\theta, b; x, y) \frac{\partial z_i^{(l+1)}}{\partial b_i^{(l)}} = \delta_i^{(l+1)}$$

# Step 3: Parameter Updates

$$J(\theta, b) = \left[ \frac{1}{m} \sum_{i=1}^m J(\theta, b; x^{(i)}, y^{(i)}) \right]$$

- Gradient descent

$$- \theta_{ij}^{(l)} = \theta_{ij}^{(l)} - \eta \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b)$$

$$- b_i^{(l)} = b_i^{(l)} - \eta \frac{\partial}{\partial b_i^{(l)}} J(\theta, b)$$

$$- \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b; x^{(i)}, y^{(i)})$$

$$- \frac{\partial}{\partial b_i^{(l)}} J(\theta, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b_i^{(l)}} J(\theta, b; x^{(i)}, y^{(i)})$$

Derivatives of last layer

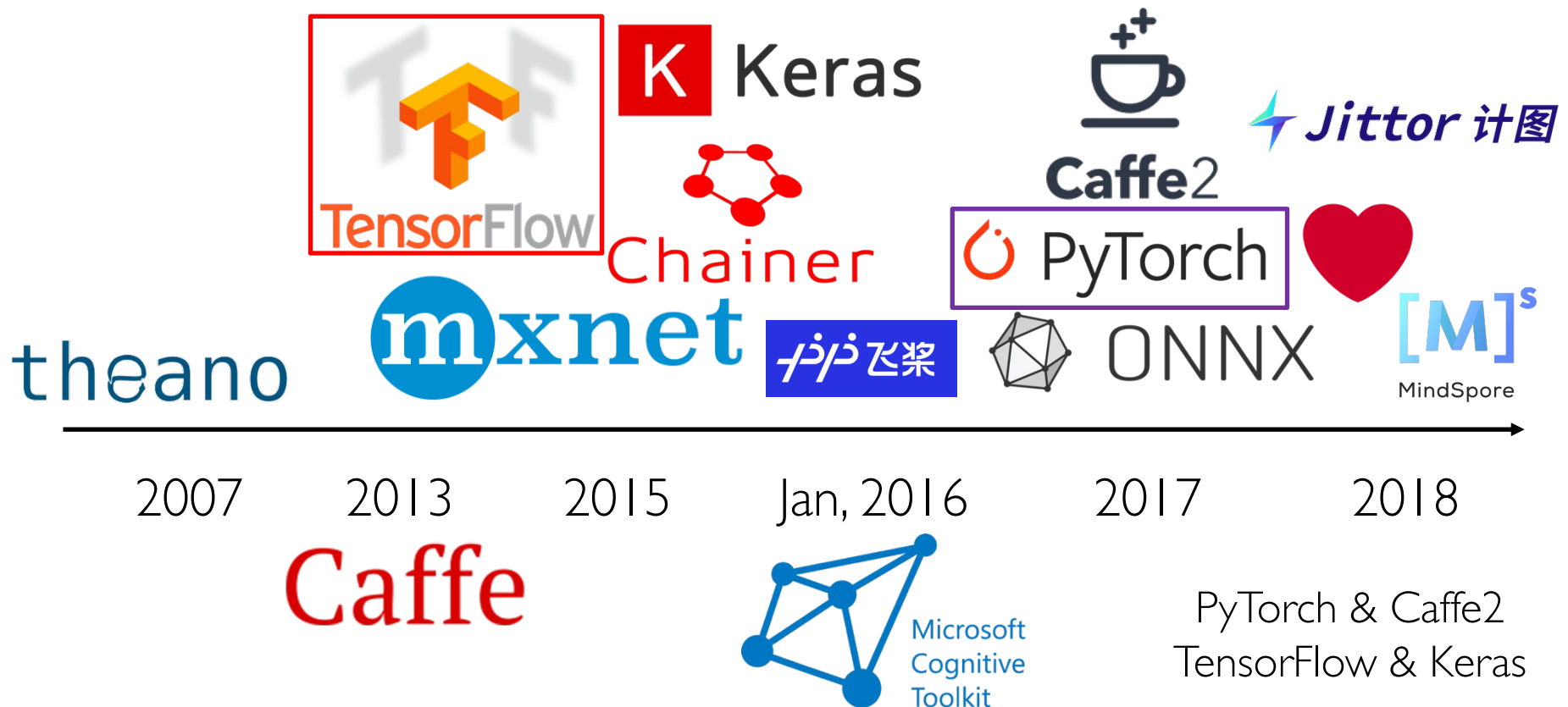
$$J(\theta, b; x^{(i)}, y^{(i)}) = - \sum_{j=1}^k \mathbf{1}\{y^{(i)} = j\} \log \frac{\exp(z_j^{(n_l)})}{\sum_{j'=1}^k \exp(z_{j'}^{(n_l)})}$$
$$\frac{\partial J(\theta, b)}{\partial z_j^{(n_l)}} = - \left( \mathbf{1}\{y^{(i)} = j\} - p(y^{(i)} = j | x^{(i)}; \theta) \right)$$



Computing all these  
by hand is painful

# Deep Learning Systems

- Symbolic vs. Imperative
- Computational Graph and Automatic Differentiation





# Outline

- Feedforward Network
  - Perceptron
  - Multilayer Perceptron
- Backpropagation
- **Practical Training Strategies**

# Stochastic Gradient Descent (SGD)

$$J(\theta, b) = \left[ \frac{1}{m} \sum_{i=1}^m J(\theta, b; x^{(i)}, y^{(i)}) \right]$$

- Stochastic Gradient Descent

- $\theta_{ij}^{(l)} = \theta_{ij}^{(l)} - \eta \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b)$  Stochasticity helps escaping from saddle points.

- $b_i^{(l)} = b_i^{(l)} - \eta \frac{\partial}{\partial b_i^{(l)}} J(\theta, b)$  Too expensive when  $m$  is very large!

- $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b) = \left[ \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta, b; x^{(i)}, y^{(i)}) \right]$  For each iteration, compute gradients for a mini-batch (小批量) of data points

- $\frac{\partial}{\partial b_i^{(l)}} J(\theta, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b_i^{(l)}} J(\theta, b; x^{(i)}, y^{(i)})$  For each epoch, Shuffle (洗牌) the dataset



# Learning Rate Decay

- All SGD algorithms take **learning rate  $\eta$**  as a **hyperparameter**
- Though some algorithms can adjust learning rate adaptively, **a good choice of learning rate  $\eta$**  could result in better performance
- To make network converge **stably and quickly**, we could set learning rate that **decays over time**

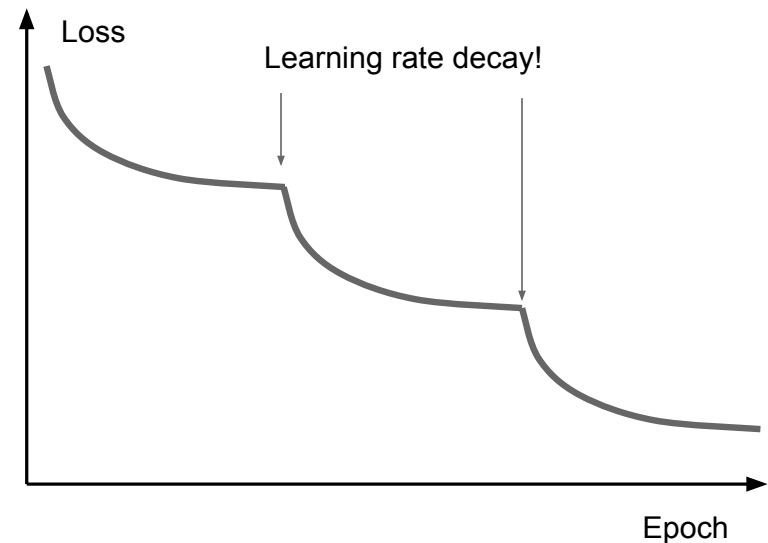
- Exponential decay strategy:

$$\eta = \eta_0 e^{-kt}$$

- 1/t decay strategy:

$$\eta = \eta_0 / (1 + kt)$$

- Step strategy: decay every  $T$  iterations (widely used in practice!)

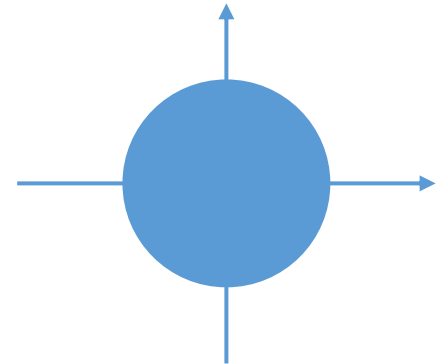


# Weight Decay

- L2 regularization:

$$\Omega(\theta) = \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ji}^{(l)})^2$$

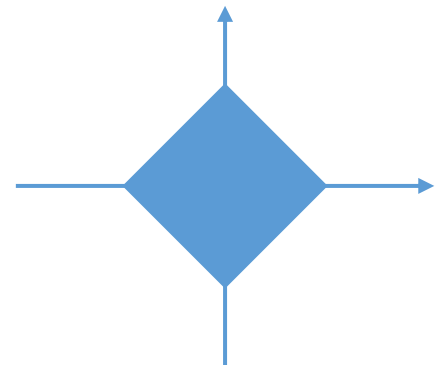
$$\frac{\partial}{\partial \theta^{(l)}} \Omega(\theta) = \lambda \theta^{(l)}$$



- L1 regularization:

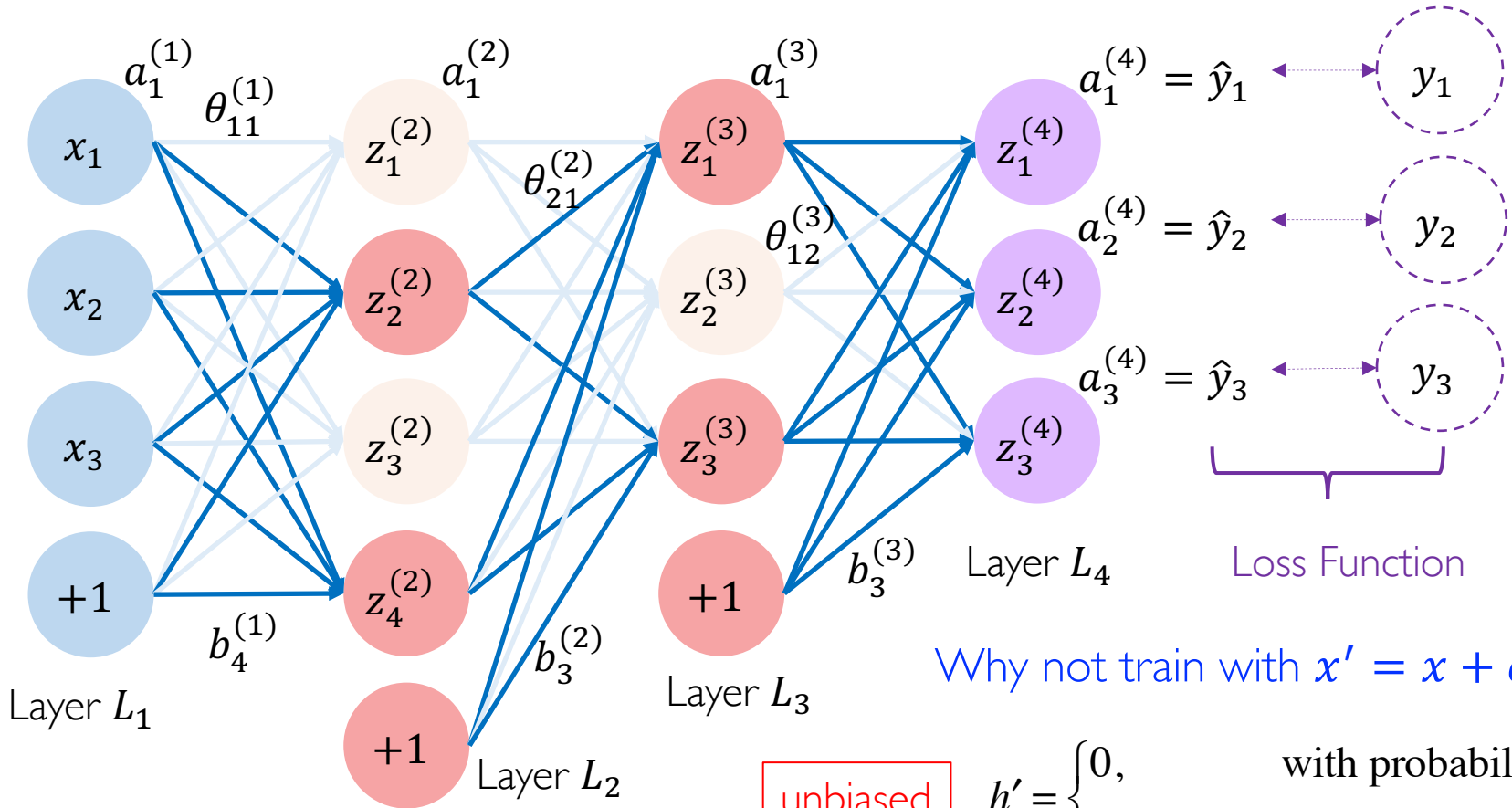
$$\Omega(\theta) = \lambda \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} |\theta_{ji}^{(l)}|$$

$$\frac{\partial}{\partial \theta^{(l)}} \Omega(\theta)_{ji} = \lambda (1_{\theta_{ji}^{(l)} > 0} - 1_{\theta_{ji}^{(l)} < 0})$$



# Dropout

- During training, randomly drop 50% of activations in each layer to **0**



Why not train with  $x' = x + \epsilon$ ?

$$\boxed{\text{unbiased}} \quad h' = \begin{cases} 0, & \text{with probability } p \\ h / (1 - p), & \text{otherwise} \end{cases}$$