

Lecture 3 Deep Feedforward Networks

Deep Learning (深度学习)

Three Steps for Deep Learning



Deep Learning is as simple as linear model.....

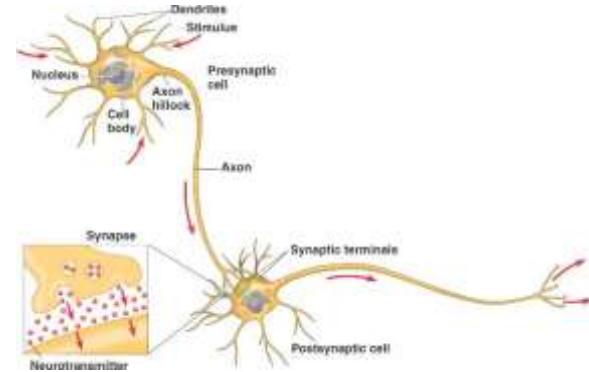
Overview

- Model Architectures
 - Artificial neurons
 - Activation function and saturation
 - Feedforward neural nets
- How to train a neural net
 - Loss Function Design
 - Optimization
 - Gradient Descent and Stochastic Gradient Descent
 - Back-propagation
 - Advanced Training tips

The Perceptron

- Invented in 1958 by Frank Rosenblatt
- Inspired by neurobiology

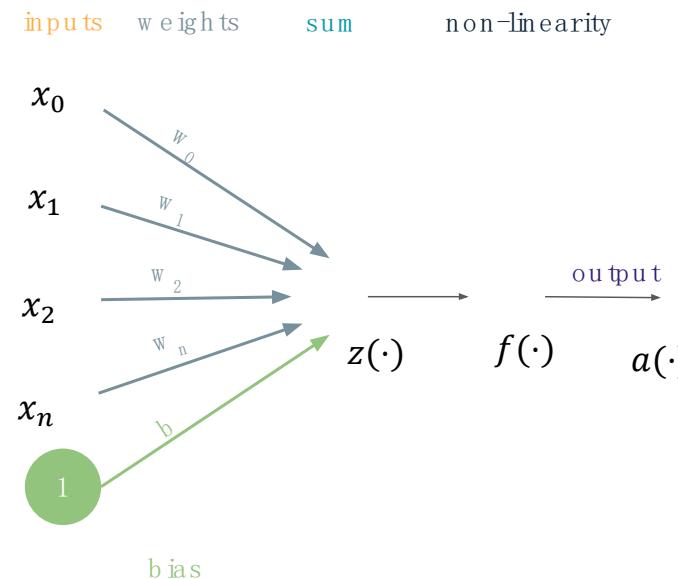
Artificial Neuron



- Each neuron is a very simple function
 - Pre-activation: $z(x) = \sum_i w_i x_i + b = w^T x + b$
 - Output activation: $a(x) = f(z(x)) = f(w^T x + b)$
- $f(\cdot)$: nonlinear activation function

w : weight

b : bias term

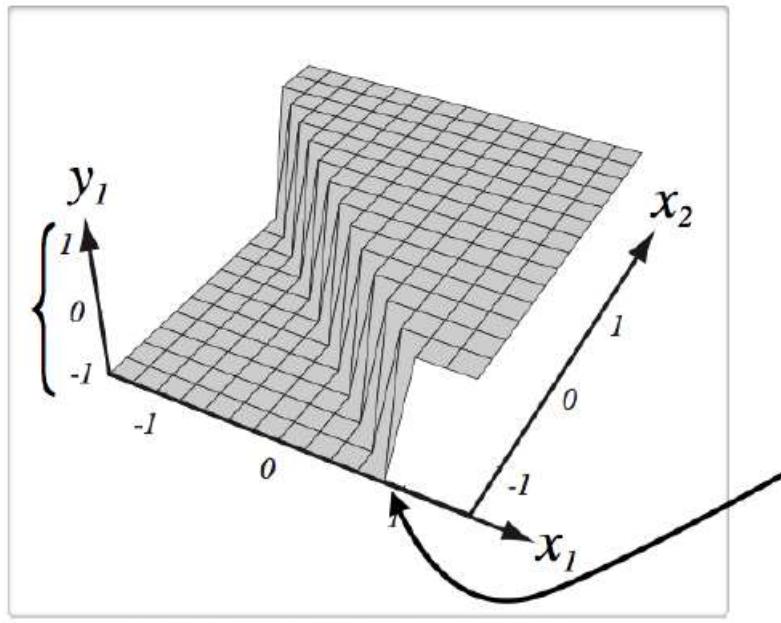


Artificial Neuron

- Output activation

$$a(x) = f(z(x)) = f(w^T x + b)$$

Range is determined by $g(\cdot)$



(from Pascal Vincent's slides)

Bias only changes the position of the riff

Activation Function

- Linear activation function
- Sigmoid activation function
- Hyperbolic tangent activation function
- Rectified linear (ReLU) activation function
- Softmax activation function



Non-linear activation function,
frequently used in deep neural networks.

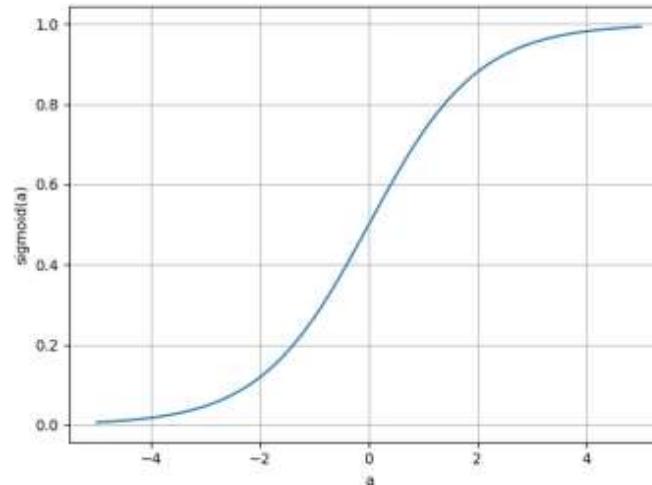
Linear Activation Function

- No input squashing
- No nonlinear transformation
$$f(z) = z$$
- Why non-linearity?
 - **Without non-linearity**, deep neural networks work the same as linear transform
 - $W_1(W_2 \cdot x) = (W_1 W_2)x = Wx$
 - **With non-linearity**, networks with more layers can approximate more complex function

Sigmoid Activation Function

- Squashes the neuron's output to (0,1)
- Bounded
- Strictly increasing

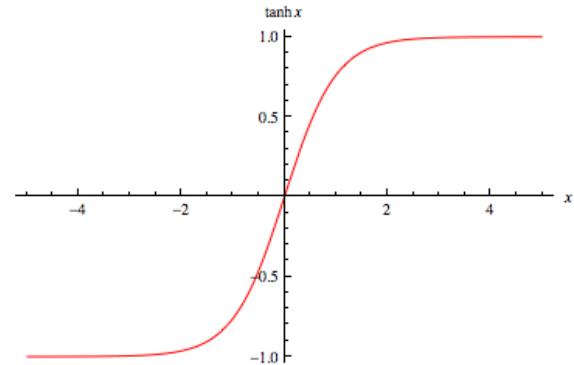
$$f(z) = \sigma(z) \stackrel{\text{def}}{=} \frac{1}{1 + \exp(-z)}$$



Hyperbolic Tangent (“tanh”) Function

- Squashes the neuron’s output to (-1,1)
- Can be positive or negative
- Bounded
- Strictly increasing

$$f(z) = \tanh(z) \stackrel{\text{def}}{=} \frac{\exp(2z) - 1}{\exp(2z) + 1}$$



Rectified Linear Activation Function (ReLU)

- Tends to produce units with sparse activities

- No upper-bound

- Increasing

$$f(z) = \text{reclin}(z) \stackrel{\text{def}}{=} \max(z, 0)$$

- ReLU variants:

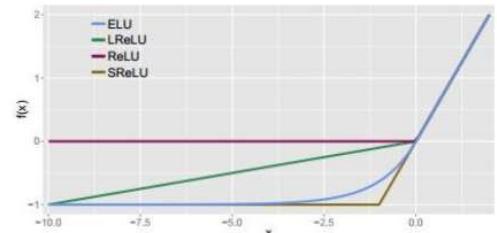
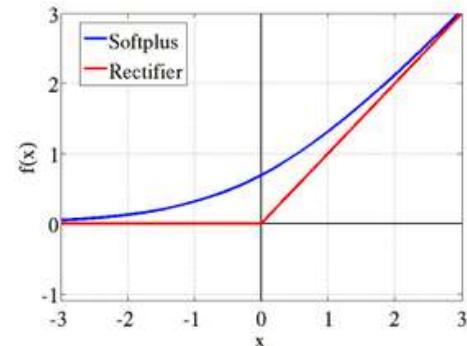
- Shift ReLU: $\max(-1, z)$

- Leaky ReLU: $\max(0.1z, z)$

- Parameter ReLU: $\max(\mu z, z)$

- Exponential Linear Units: $\max(z, \mu(\exp(z) - 1))$

- Maxout: $\max(w_1^T z + b_1, w_2^T z + b_2)$



Softmax Activation Function

- In multi-class classification (C classes), we need to
 - generate multiple output: $\mathbf{z} \in \mathbb{R}^C$ (in vector form)
 - estimate the conditional probability of each class:
$$p(y = i | \mathbf{z}) = \frac{\exp z_i}{\sum_c \exp z_c}$$
- Strictly positive
- Sum up to one

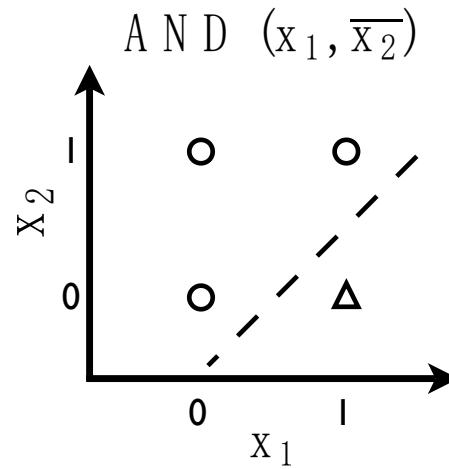
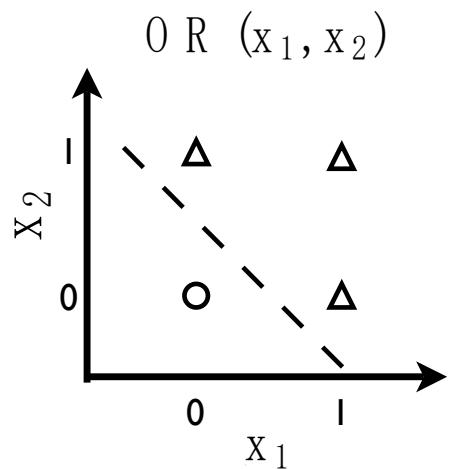
$$\mathbf{f}(\mathbf{z}) = \text{softmax}(\mathbf{z}) = \left[\frac{\exp z_1}{\sum_c \exp z_c} \cdots \frac{\exp z_C}{\sum_c \exp z_c} \right]^T$$

Characteristics of Activation Functions

- 连续可导（允许少数点上不可导）的非线性函数
 - 参数学习
- 激活函数及其导函数要尽可能的简单
 - 有利于网络计算效率
- 激活函数的导函数的值域要在一个合适的区间内
 - 太大或太小将影响训练的效率和稳定性
- 单调递增

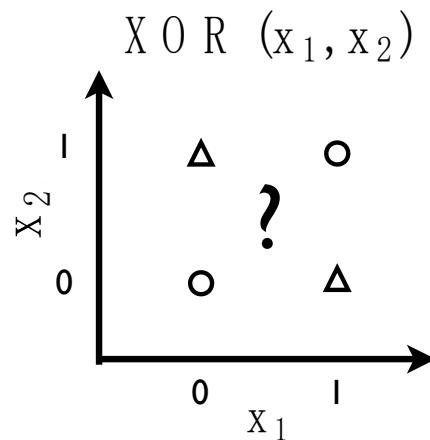
Capacities of a Single Neuron

- Solve linearly separable problems



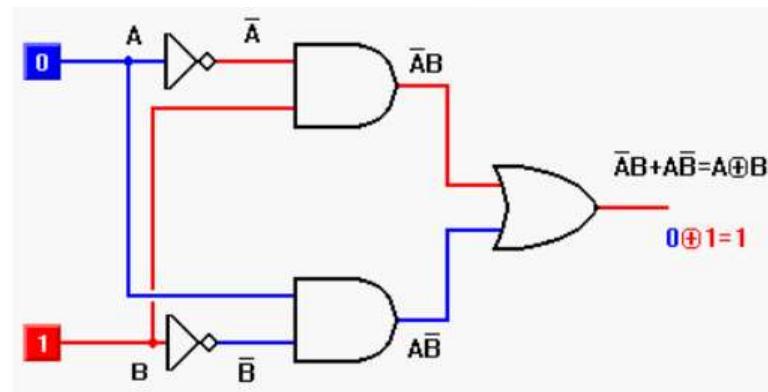
Capacities of a Single Neuron

- Can't solve linearly inseparable problems



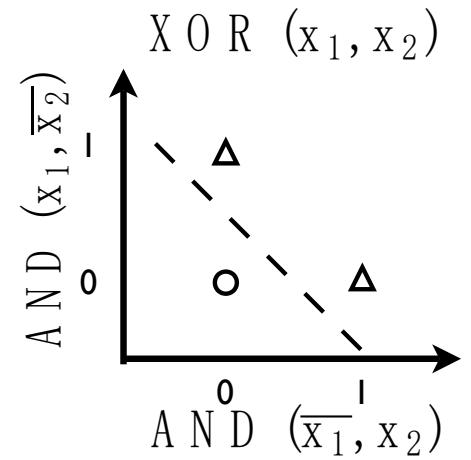
How to implement XOR?

Input		Output
A	B	
0	0	0
0	1	1
1	0	1
1	1	0

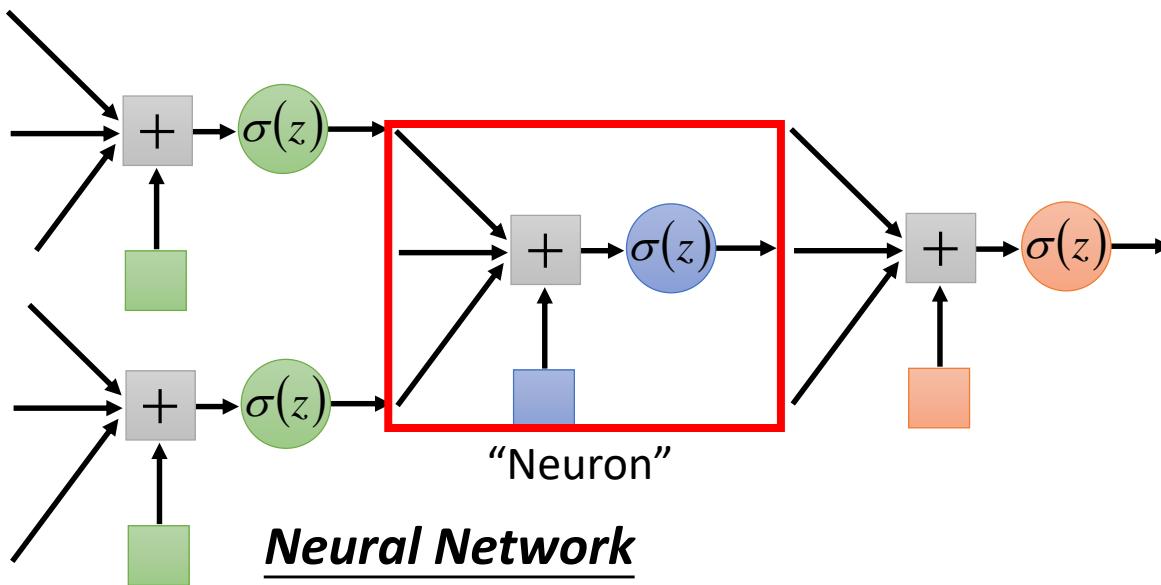


$$A \text{ XOR } B = AB' + A'B$$

Multiple operations can produce more complicate output



Neural Network

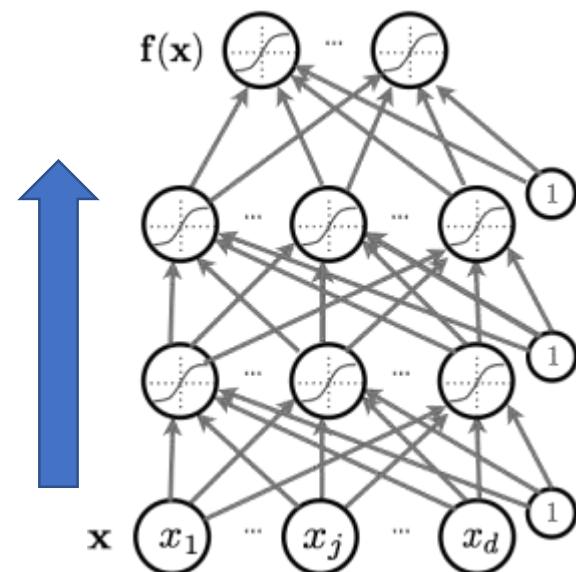


Different connection leads to different network structures

Network parameter θ : all the weights and biases in the “neurons”

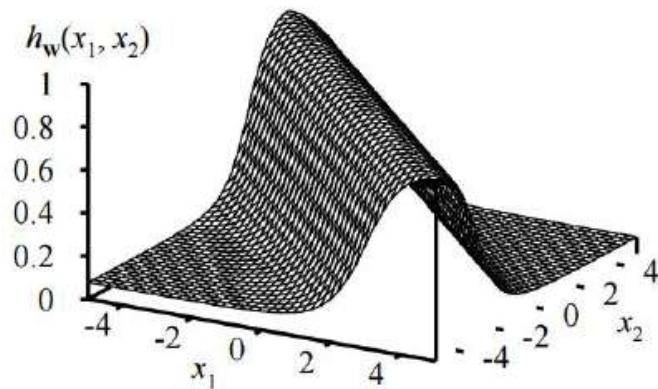
Multilayer Feedforward Neural Networks

- Each neuron in one layer has directed connections to the neurons of the subsequent layer
- Information propagates from input x to output $f(x)$, through many hidden layers

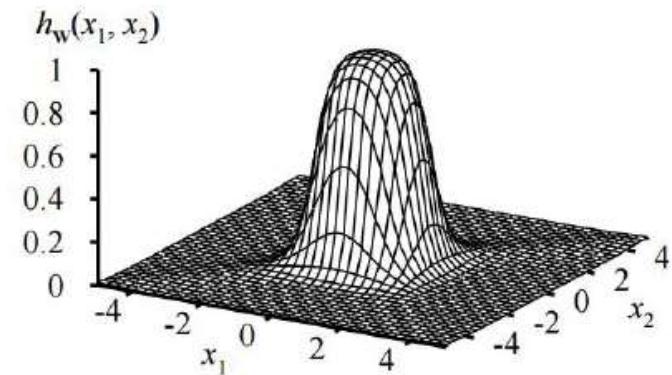


Expressions of Multi-Layer Neural Network

Continuous function w/ 2 layers



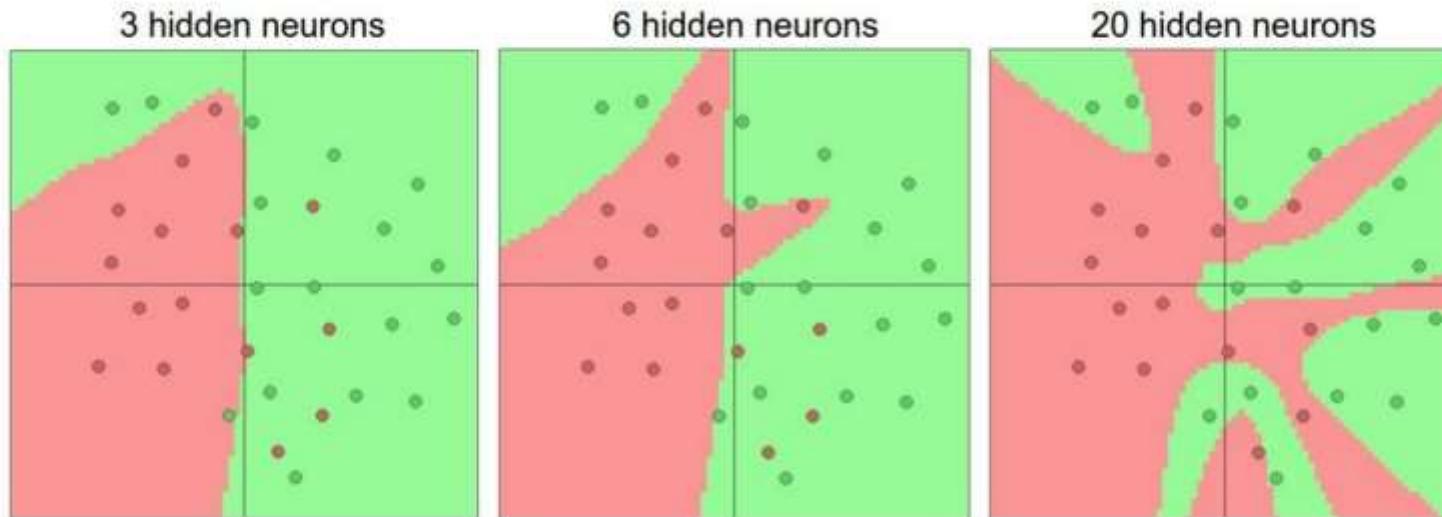
Continuous function w/ 3 layers



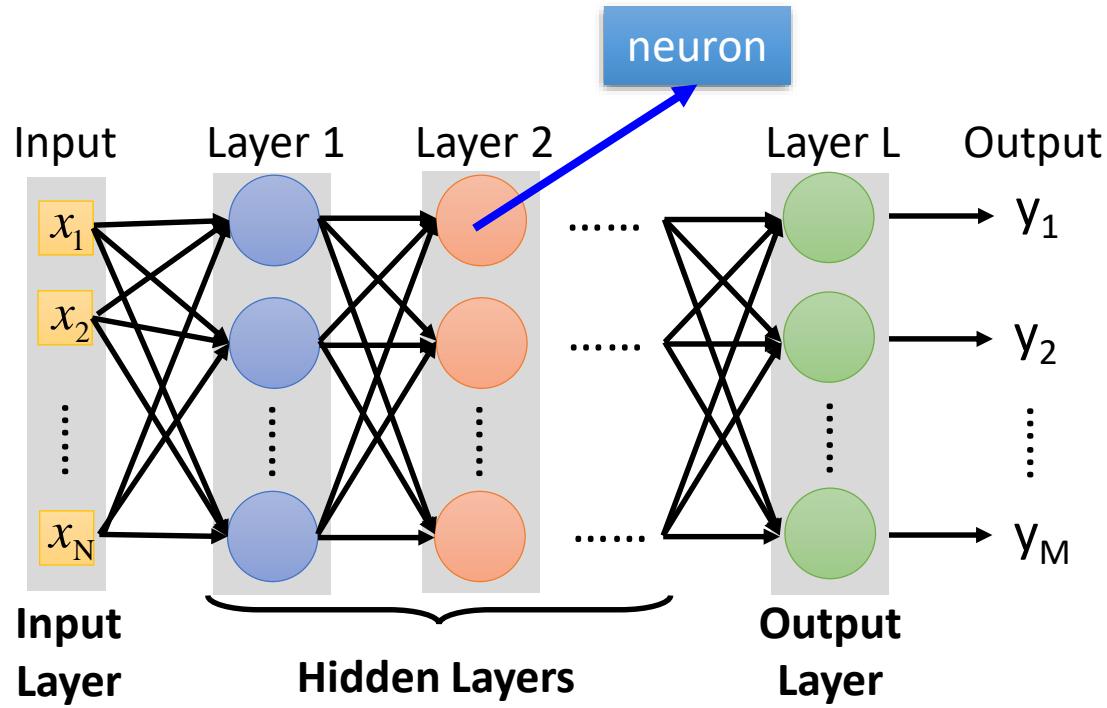
Multiple layers enhance the model expression -> the model can approximate more complex functions

Setting the Number of Neurons and Layers

- More neurons = more capacity
- More layers = more capacity

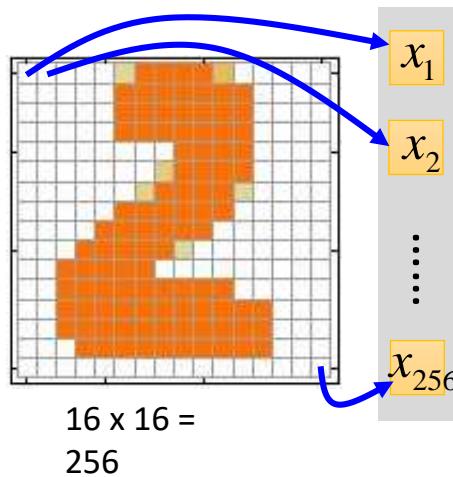


Fully Connect Feedforward Network



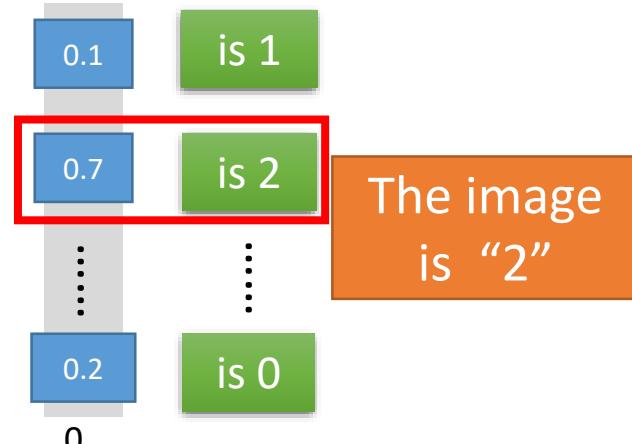
Example Application

Input



Ink → 1
No ink → 0

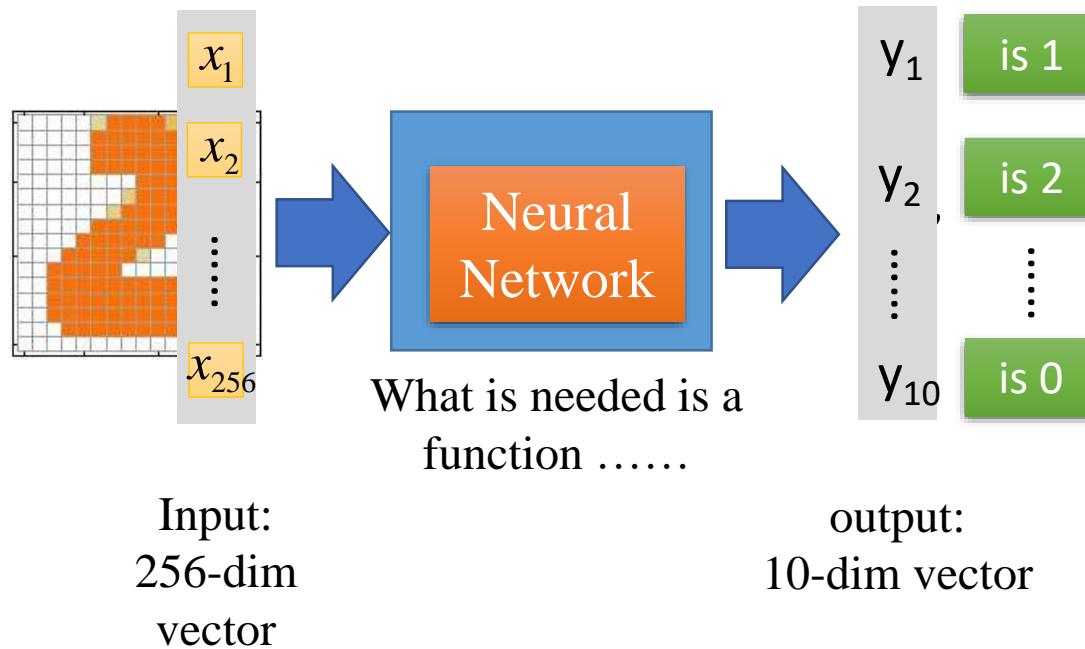
Output



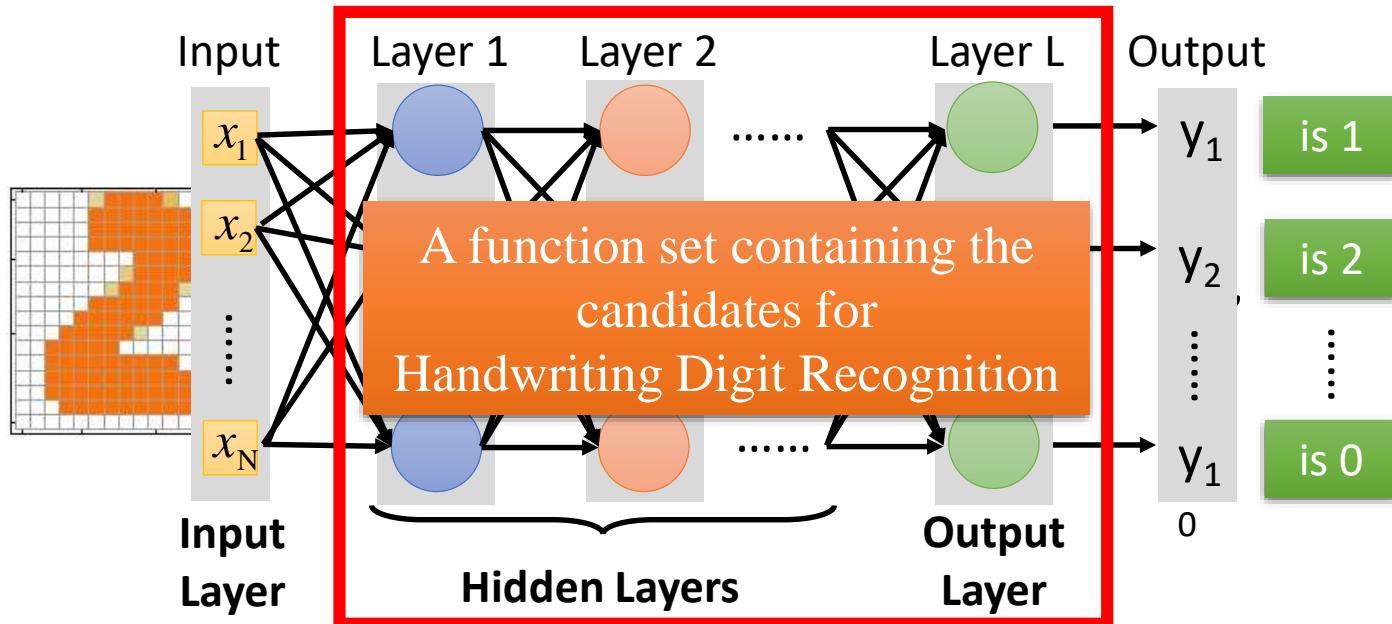
Each dimension represents the confidence of a digit.

Example Application

- Handwriting Digit Recognition

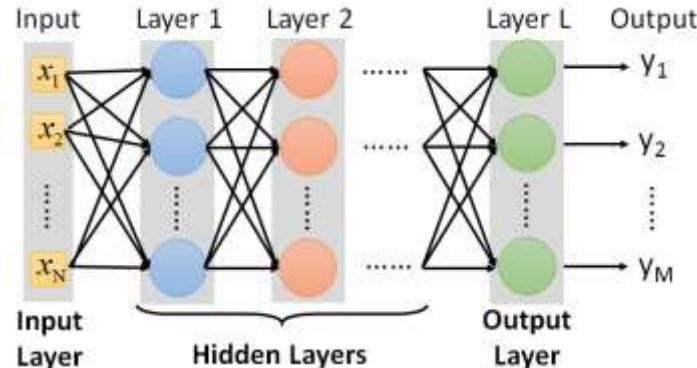


Example Application



You need to decide the network structure to let a good function in your function set.

FAQ



- Q: How many layers? How many neurons for each layer?

Trial and Error

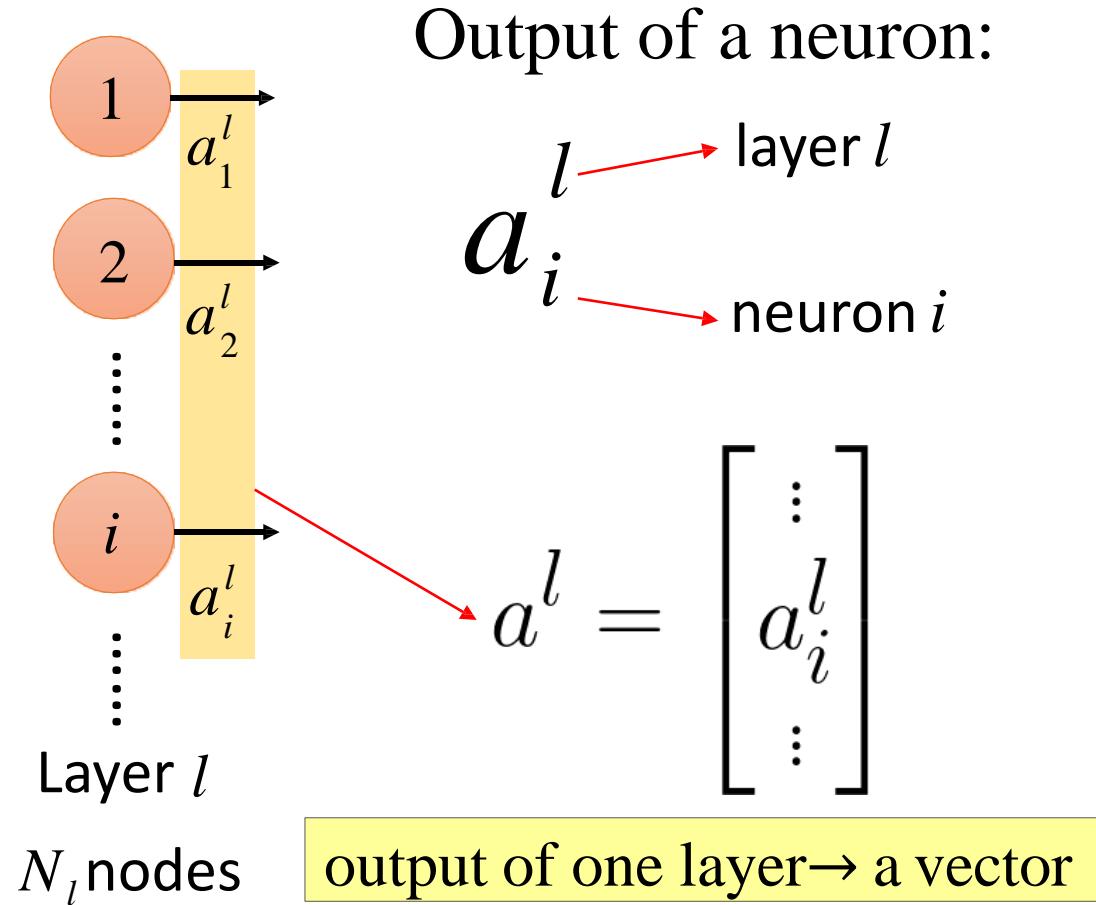
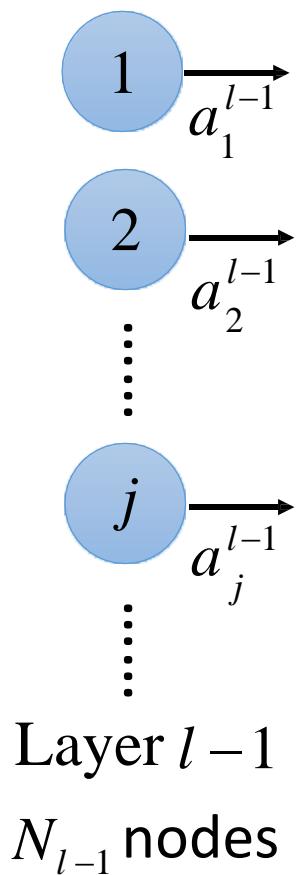
+

Intuition

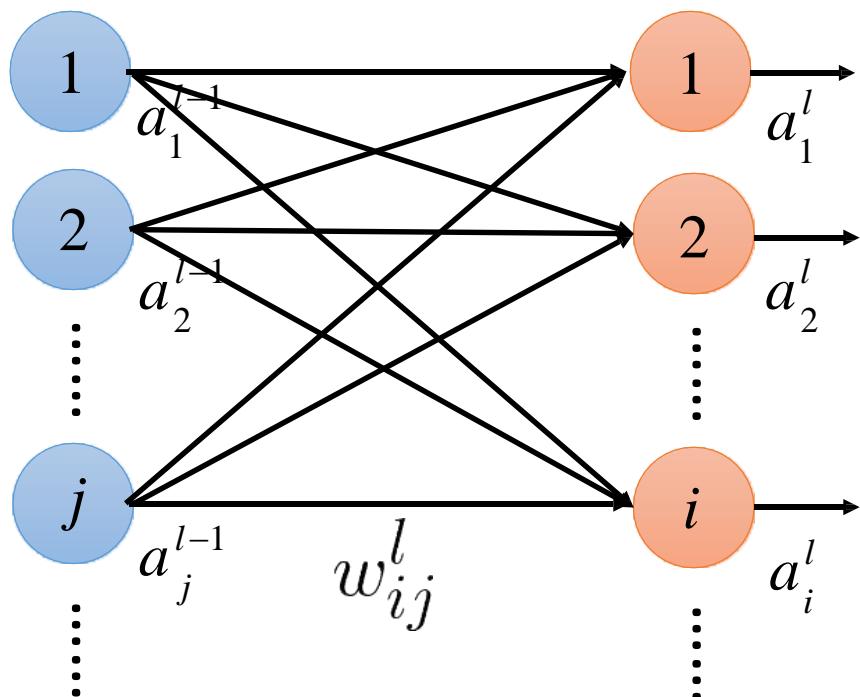
- Q: Can the structure be automatically determined?
 - E.g. Evolutionary Artificial Neural Networks
- Q: Can we design the network structure?

Convolutional Neural Network (CNN)

Notation Definition



Notation Definition



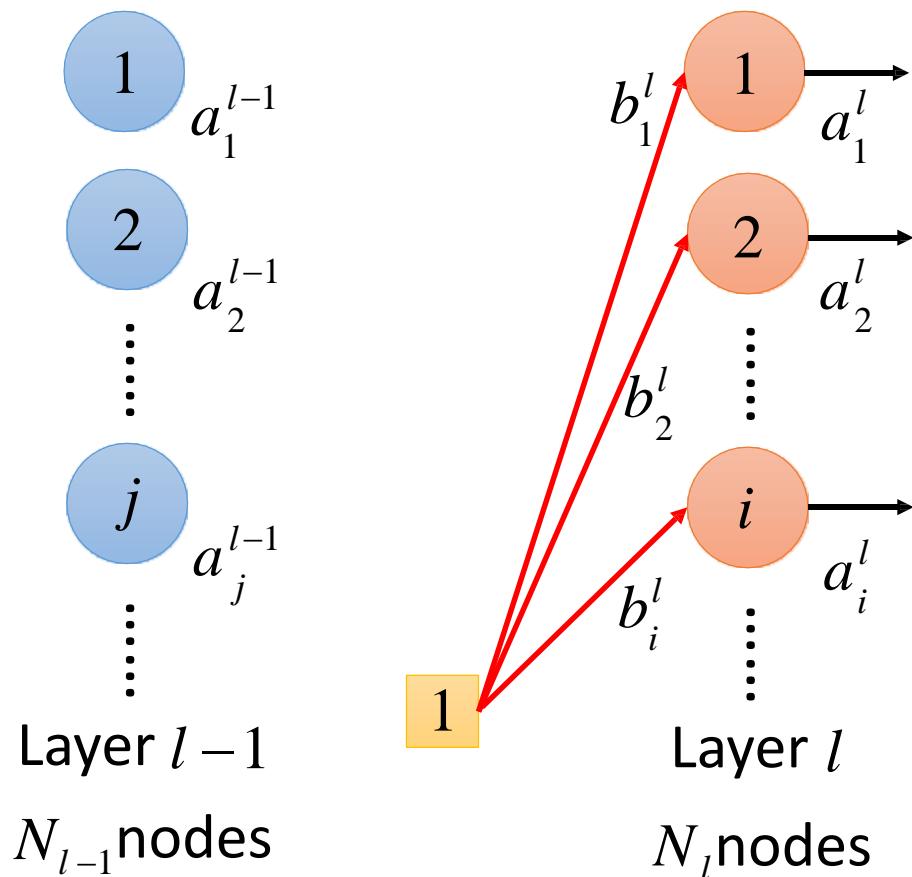
Layer $l-1$
 N_{l-1} nodes

Layer l
 N_l nodes

$$w_{ij}^l \quad \begin{array}{l} \text{layer } l-1 \\ \text{to layer } l \\ \text{from neuron } j \text{ (layer } l-1\text{)} \\ \text{to neuron } i \text{ (layer } l\text{)} \end{array}$$
$$N_{l-1} \quad W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & & \ddots \end{bmatrix} N_l$$

weights between two layers
→ a matrix

Notation Definition

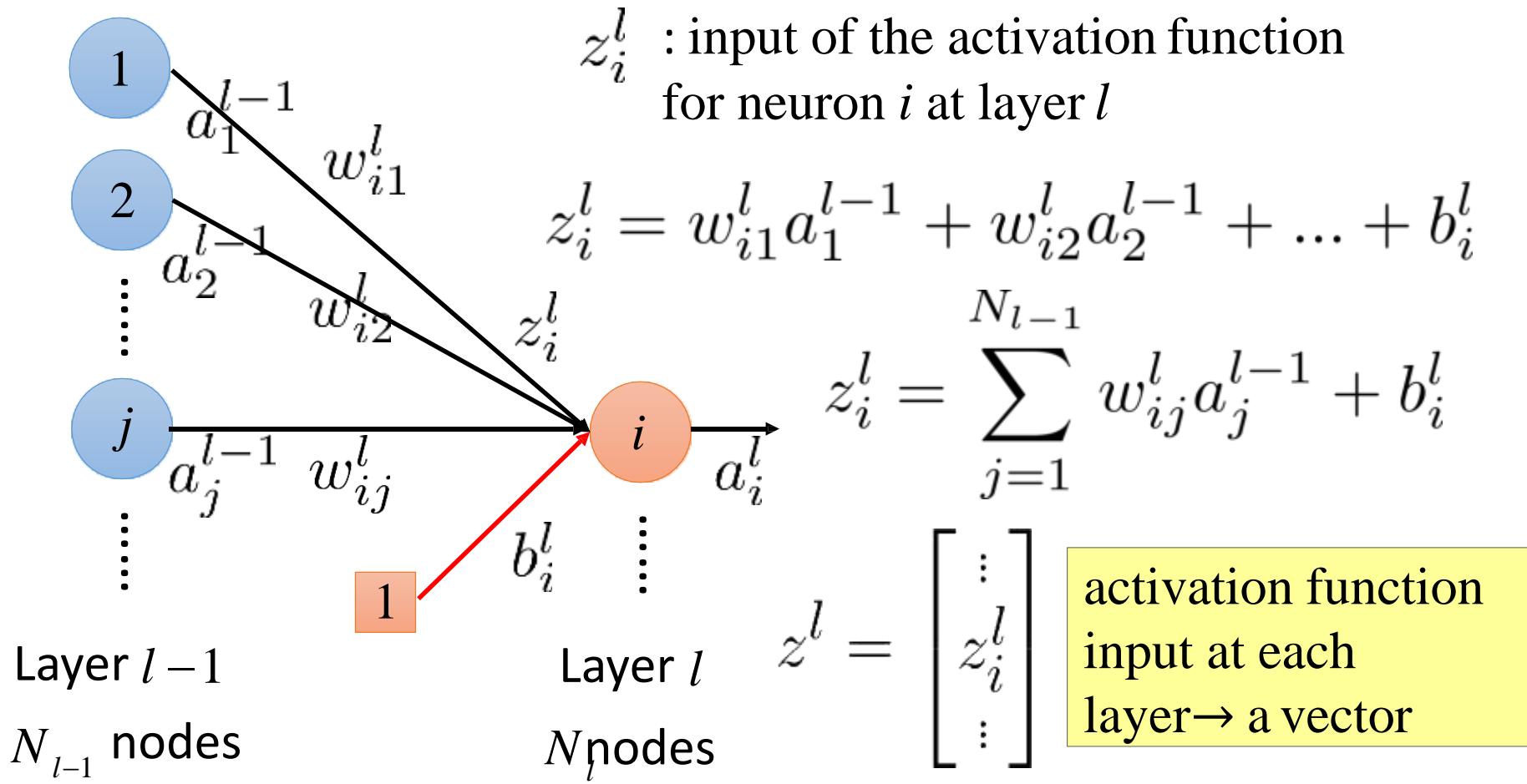


b_i^l : bias for neuron i at layer l

$$b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

bias of all neurons at each layer → a vector

Notation Definition



Notation Summary

a_i^l : output of a neuron

w_{ij}^l : a weight

a^l : output vector of a layer

W^l : a weight matrix

z_i^l : input of activation
function

b_i^l : a bias

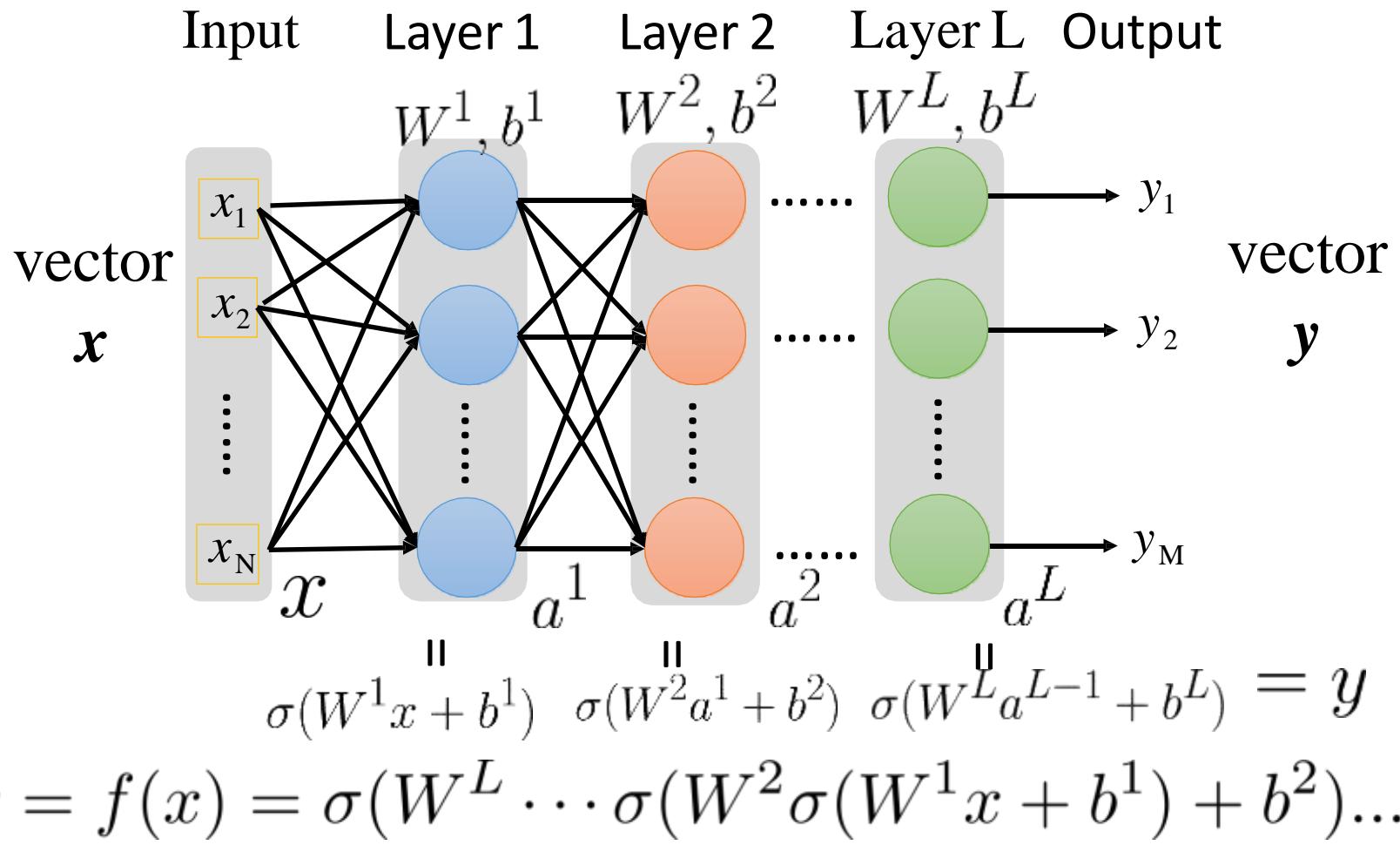
z^l : input vector of activation
function for a layer

b^l : a bias vector

Neural Network Formulation

$$f : R^N \rightarrow R^M$$

Fully connected feedforward network



Function = model parameters

Forward propagation

$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

Different parameters W and $b \rightarrow$ different functions

Formal definition:

$$f(x; \theta); \quad \theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

Pick a function $f =$ pick a set of model parameters θ

Training

- Empirical risk minimization $J(\theta)$

$$\arg \min_{\theta} \frac{1}{T} \sum_t l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) + \lambda \Omega(\theta)$$



Loss function Regularizer

- Learning is cast as optimization
 - Find a model parameter set that minimize $J(\theta)$
 - Loss function can sometimes be viewed as a surrogate for what we want to optimize

Loss Function

- In discriminative model (判别模型), model $y|x$.
- Learning the maximum likelihood, equivalently the cross entropy between training data and model distribution:

$$l(\theta) = -E_{x,y \sim \text{Data}} \log p(y|x)$$

- The specific form of loss function depends on the model distribution $p(\cdot)$

Loss Functions

- Loss function evaluates the performance of our model, it is chosen according to the output units
 - Normal: $\hat{y} = \mathbf{w}^T \mathbf{x} + b$
 - Bernoulli: $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b)$
 - Multinomial: $\hat{\mathbf{y}} = \text{softmax}(\mathbf{W}^T \mathbf{x} + \mathbf{b})$
- Consider regularization $\Omega(\theta)$
- Equivalent form of Loss function: $J = l(y, \hat{y}) + \lambda \Omega(\theta)$

Frequently Used Loss Functions

- Square loss

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

- Hinge loss

$$l(y, \hat{y}) = \max(0, 1 - \hat{y}y)$$

- Logistic loss

$$l(y, \hat{y}) = \log(1 + \exp(-\hat{y}y))$$

- Cross entropy loss

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

How to Train Multilayer Neural Nets?

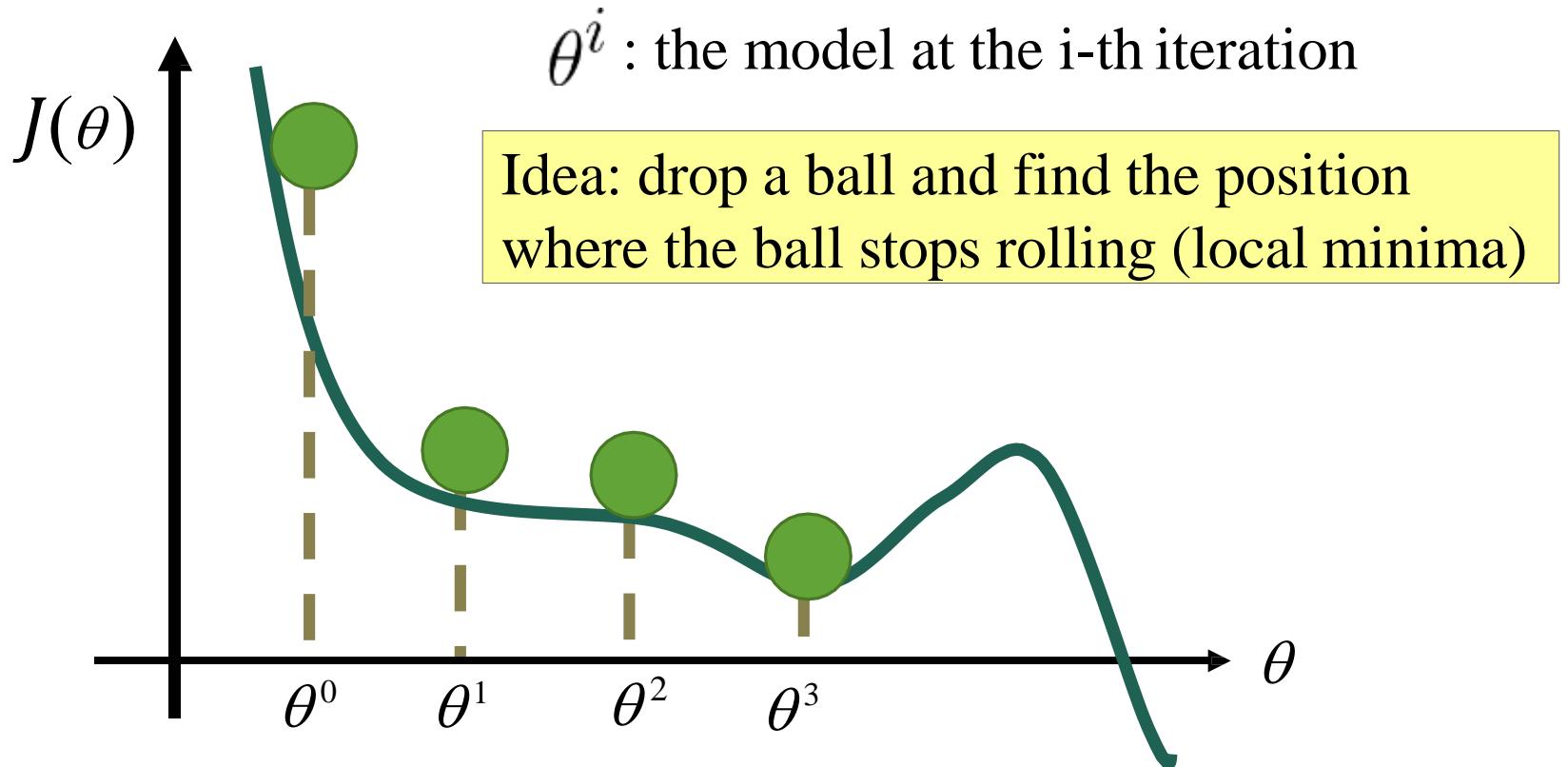
- Learning is reduced to optimization.
 - Given a loss function $J(\theta)$ and several parameter sets
 - Find a model parameter set that minimize $J(\theta)$

Overview

- Model Architectures
 - Artificial neurons
 - Activation function and saturation
 - Feedforward neural nets
- How to train a neural net
 - Loss Function Design
 - Optimization
 - Gradient Descent and Stochastic Gradient Descent
 - Backward propagation

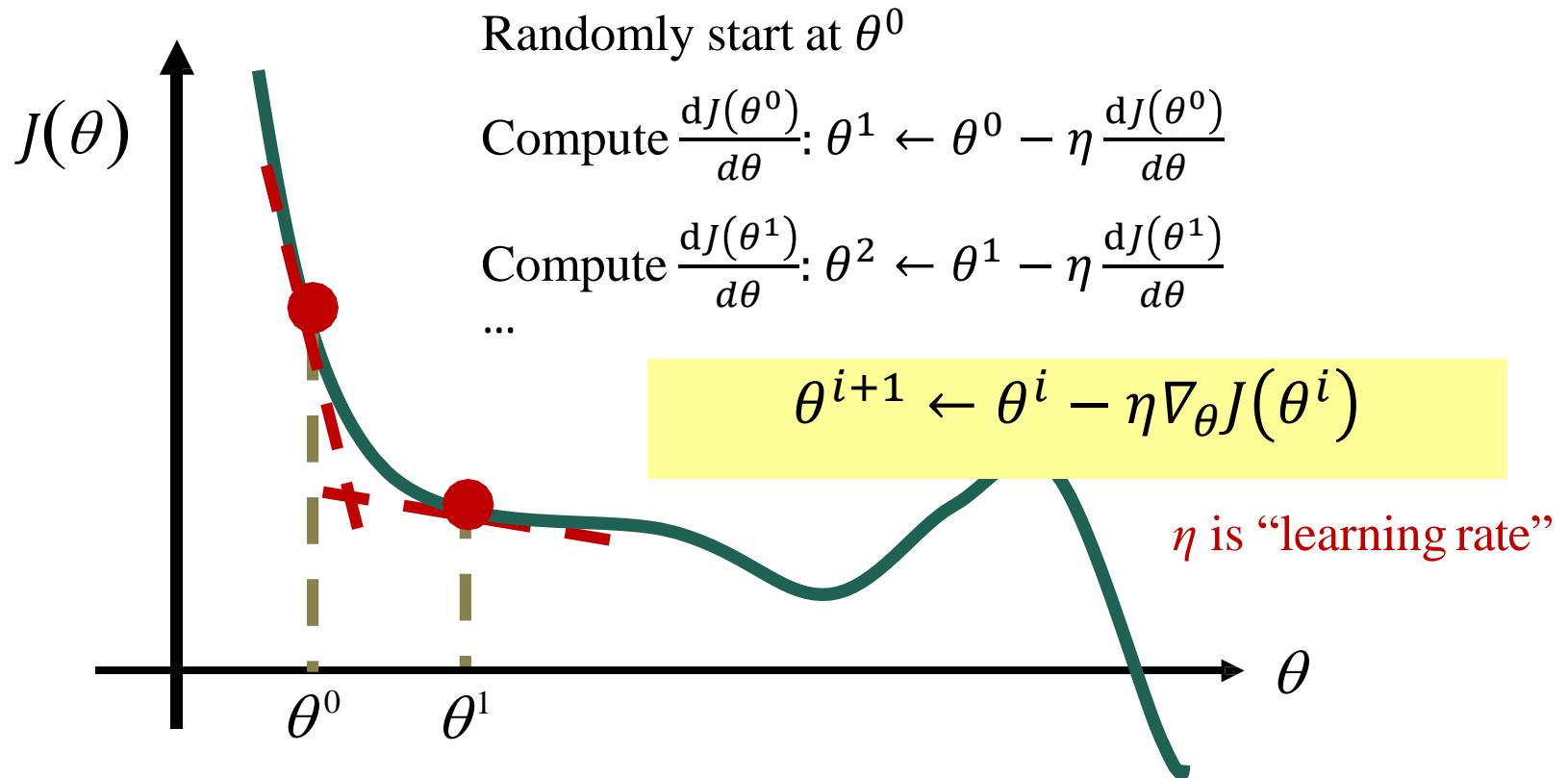
Gradient Descent for Optimization

Assume that θ has only one variable

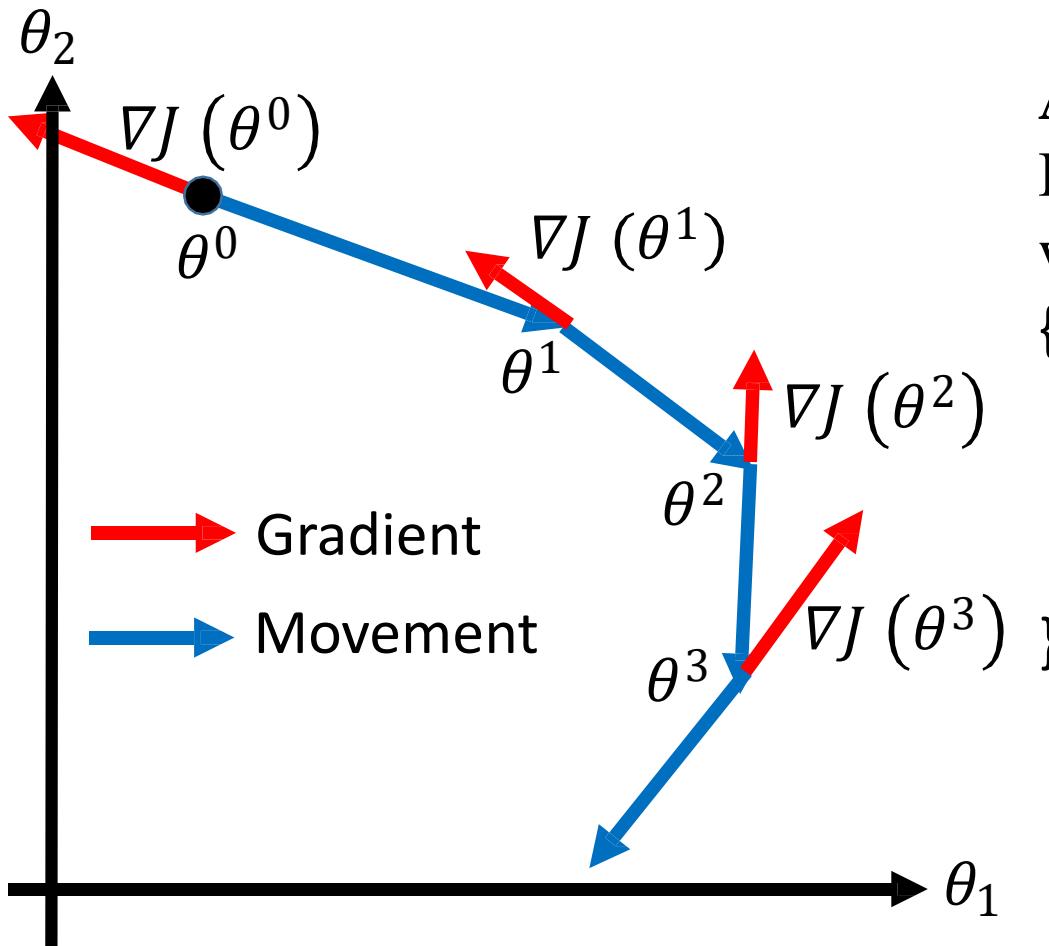


Gradient Descent for Optimization

Assume that θ has only one variable



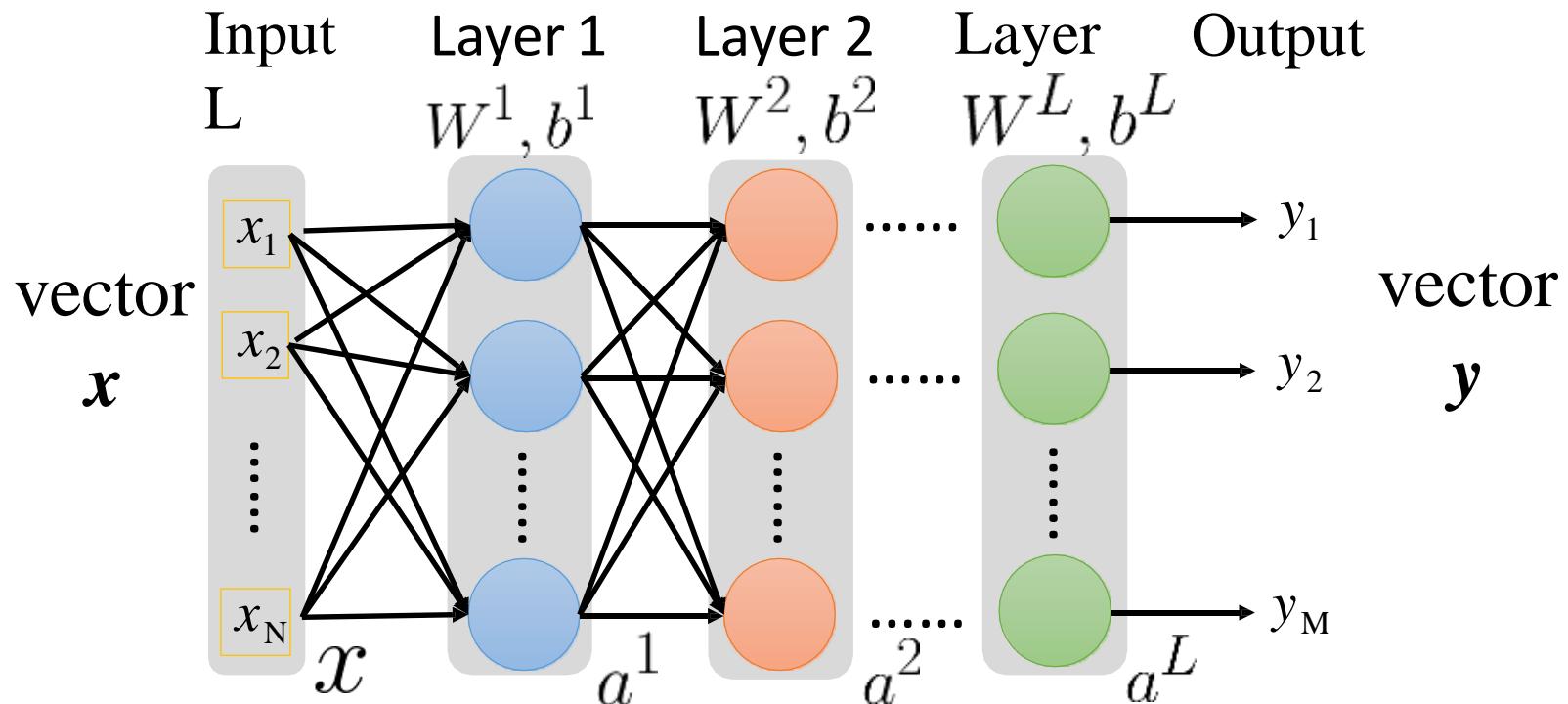
Gradient Descent for Optimization



Algorithm

```
Initialization: start at  $\theta^0$ 
while( $\theta^{(i+1)} \neq \theta^i$ )
{
    compute gradient at  $\theta^i$ 
    update parameters
     $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} J(\theta^i)$ 
}
```

Revisit Neural Network Formulation



$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

Gradient Descent for Neural Network

$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

$$W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \dots \\ w_{21}^l & w_{22}^l & \dots \\ \vdots & & \dots \end{bmatrix} \quad b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$\nabla J(\theta) = \begin{bmatrix} \vdots \\ \frac{\partial J(\theta)}{\partial w_{ij}^l} \\ \vdots \\ \frac{\partial J(\theta)}{\partial b_i^l} \\ \vdots \end{bmatrix}$$

Algorithm

Initialization: start at θ^0

while($\theta^{(i+1)} \neq \theta^i$)

{

 compute gradient at θ^i

 update parameters

}

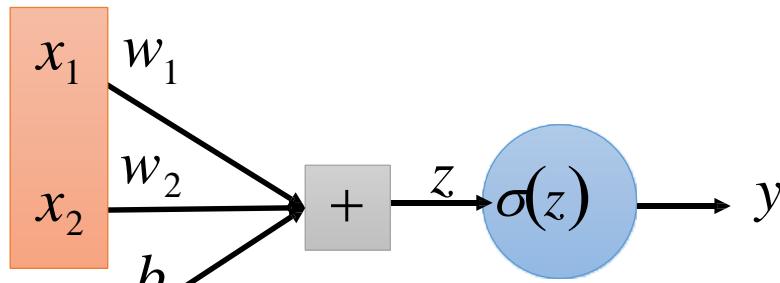
$$\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta J(\theta^i)$$

Gradient Descent for Optimization

Simple Case

$$y = f(x; \theta) = \sigma(Wx + b)$$

$$\theta = \{W, b\} = \{w_1, w_2, b\}$$



$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial w_1} \\ \frac{\partial J(\theta)}{\partial w_2} \\ \frac{\partial J(\theta)}{\partial b} \end{bmatrix}$$

Algorithm

Initialization: start at θ^0

while($\theta^{(i+1)} \neq \theta^i$)

{

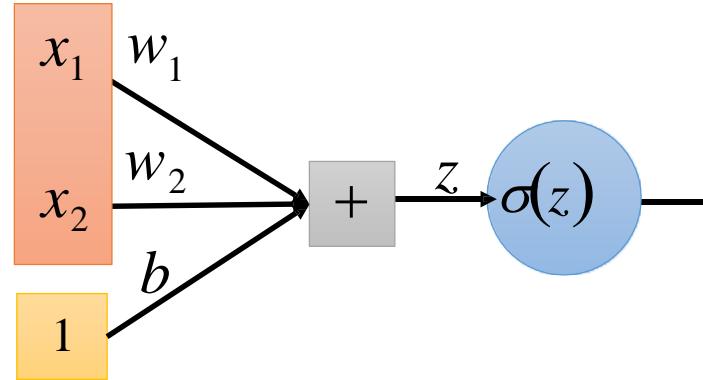
 compute gradient at θ^i
 update parameters

}

$$\theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} J(\theta^i)$$

$$\begin{bmatrix} w_1^{i+1} \\ w_2^{i+1} \\ b^{i+1} \end{bmatrix} = \begin{bmatrix} w_1^i \\ w_2^i \\ b^i \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial J(\theta^i)}{\partial w_1} \\ \frac{\partial J(\theta^i)}{\partial w_2} \\ \frac{\partial J(\theta^i)}{\partial b} \end{bmatrix}$$

To compute the Gradients



- If square loss

- $\hat{y} = \sigma(Wx + b) = \sigma(w_1x_1 + w_2x_2 + b)$
- $J(\theta) = (\sigma(Wx + b) - y)^2$
- $\frac{\partial J(\theta)}{\partial w_1} = 2(\sigma(Wx + b) - y)(1 - \sigma(Wx + b))\sigma(Wx + b)x_1$
- $\frac{\partial J(\theta)}{\partial w_2} = 2(\sigma(Wx + b) - y)(1 - \sigma(Wx + b))\sigma(Wx + b)x_2$
- $\frac{\partial J(\theta)}{\partial b} = 2(\sigma(Wx + b) - y)(1 - \sigma(Wx + b))\sigma(Wx + b)$

Algorithm

Initialization: start at θ^0

while($\theta^{(i+1)} \neq \theta^i$)

{compute gradient at θ^i

 update parameters

$\theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} J(\theta^i)$ }

Gradient Descent for Neural Network

$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

$$W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \dots \\ w_{21}^l & w_{22}^l & \dots \\ \vdots & & \dots \end{bmatrix} \quad b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$\nabla J(\theta) = \begin{bmatrix} \vdots \\ \frac{\partial J(\theta)}{\partial w_{ij}^l} \\ \vdots \\ \frac{\partial J(\theta)}{\partial b_i^l} \\ \vdots \end{bmatrix}$$

Algorithm

Initialization: start at θ^0

while($\theta^{(i+1)} \neq \theta^i$)

{

 compute gradient at θ^i

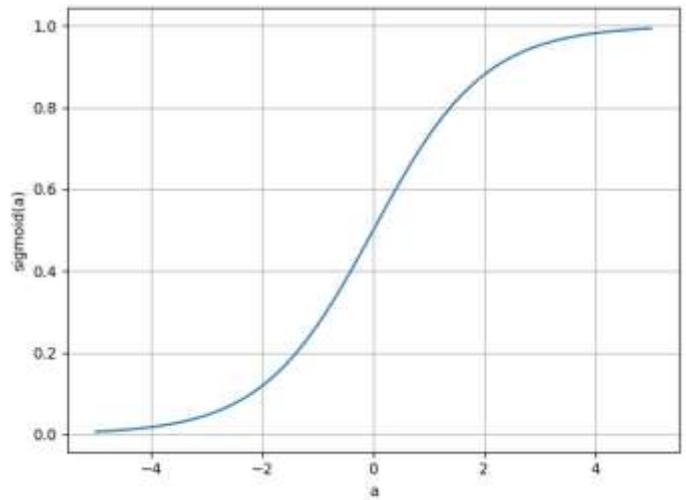
 update parameters

}

$$\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta J(\theta^i)$$

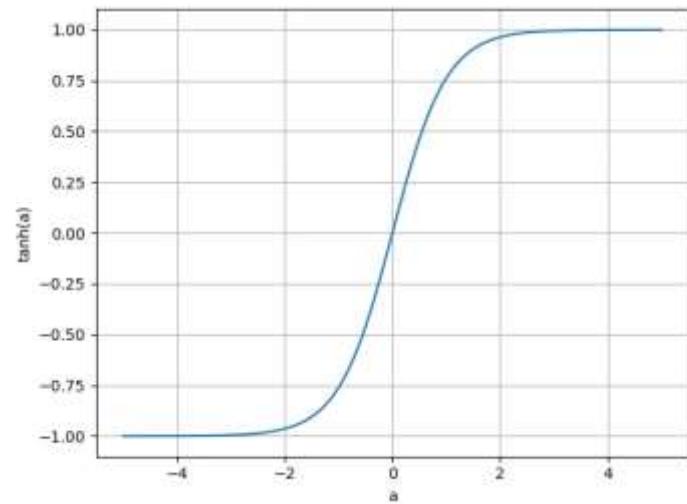
Gradient Computation: Sigmoid Unit

$$f(z) = \sigma(z), \quad f'(z) = \sigma(z)(1 - \sigma(z))$$



Gradient Computation:Tanh Unit

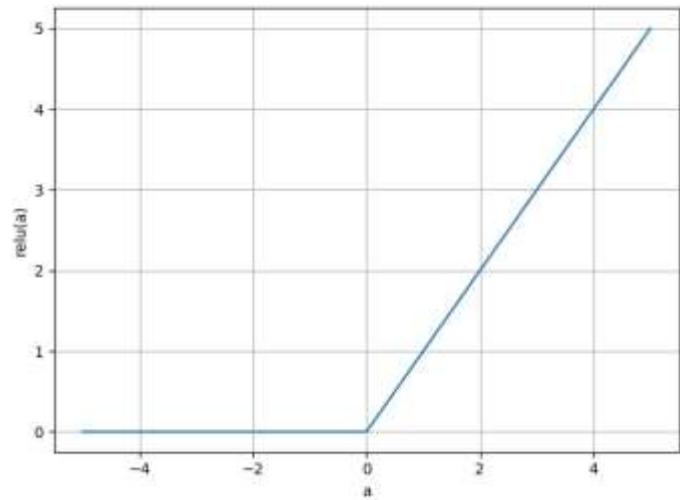
- $f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ and $\tanh'(z) = 1 - f(z)^2$
- $\tanh(z)$ approximates linear function when z is small
- Often it is preferable to sigmoid in feedforward neural nets (zero-centered)
- Problem: still kill gradient when saturated ☺



[LeCun et al., 1991]

Gradient Computation: ReLU

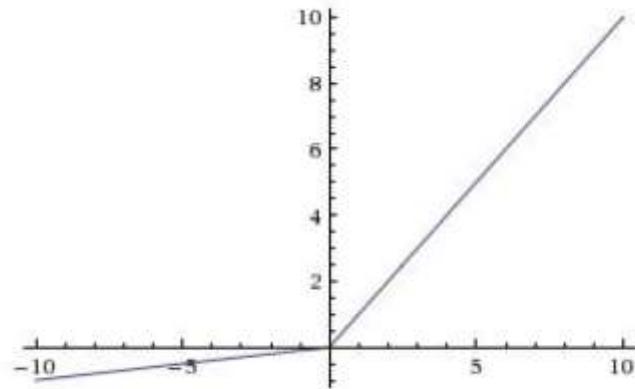
- $f(z) = \max(z, 0)$
- $f'(z) = \begin{cases} 1, & z \geq 0 \\ 0, & \text{o. w.} \end{cases}$
- Models are easier to optimize if their behavior is close to linear
- Converge much faster than sigmoid and tanh in practice (6x faster)
- Not differentiable at $z = 0$, but it is not a problem in practice
- Not zero-centered
- Fragile during training and can “die”



[Krizhevsky et al., 2012]

Gradient Computation: Leaky ReLU

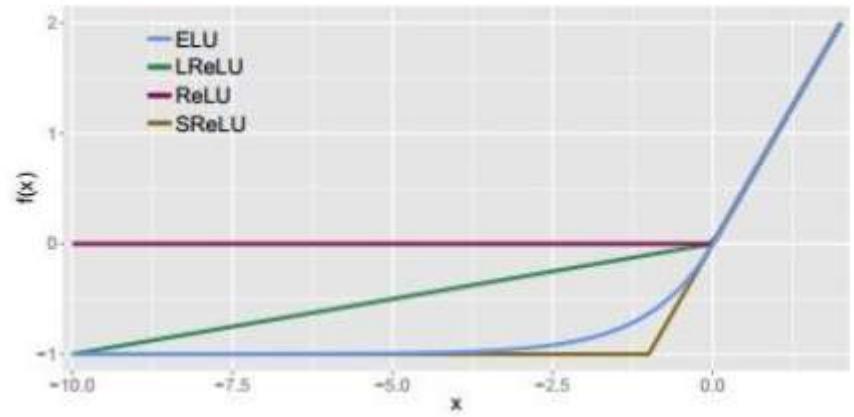
- $f(a) = \max(z, 0.01z)$
- $f'(a) = \begin{cases} 1, & z \geq 0 \\ 0.01, & \text{o. w.} \end{cases}$
- Will not die
- Parametric ReLU
 - $f(z) = \max(z, \mu z)$
 - $f'(z) = \begin{cases} 1, & z \geq 0 \\ \mu, & \text{o. w.} \end{cases}$
 - Update μ through backpropagation



[Mass et al., 2013]
[He et al., 2015]

Gradient Computation: Exponential Linear Unit (ELU)

- $f(z) = \begin{cases} z, & z \geq 0 \\ \mu(e^z - 1), & z < 0 \end{cases}$
- $f'(z) = \begin{cases} 1, & z \geq 0 \\ \mu e^z, & o.w. \end{cases}$
- All benefit of ReLU
- Almost zero-centered
- Compute $\exp()$ ☹



[Clevert et al., 2015]

Gradient Computation: Maxout

- $f(z) = \max(w_1 z + b_1, w_2 z + b_2)$
- Generalization of Leaky ReLU and ReLU
- Double the number of parameters

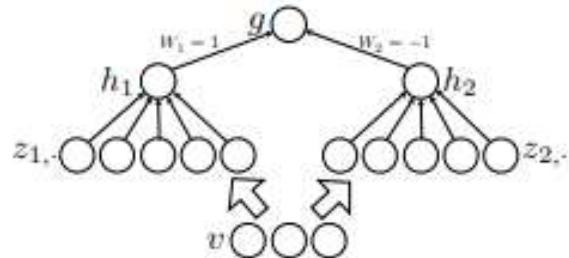


Figure 3. An MLP containing two maxout units can arbitrarily approximate any continuous function. The weights in the final layer can set g to be the difference of h_1 and h_2 . If z_1 and z_2 are allowed to have arbitrarily high cardinality, h_1 and h_2 can approximate any convex function. g can thus approximate any continuous function due to being a difference of approximations of arbitrary convex functions.

Softmax

- Cross-entropy:

$$H(p, q) = - \sum_x p(x) \log(q(x))$$

- In multiclass classification,

$$\mathbf{y} = [0, 0, 0, \dots, 0, 1, 0, 0, \dots, 0]^T \in \mathbb{R}^C,$$

Then

$$J = H(\mathbf{y}, \mathbf{h}) = - \log h_y$$

where $h_i = f_i(\mathbf{z}) = \frac{\exp(z_i)}{\sum_c^C \exp(z_c)} = P(y = i | \mathbf{z}), 1 \leq i \leq C.$

Softmax

- $\nabla_{\mathbf{h}} J = [0, \dots, 0, -h_y, 0, \dots, 0]^T$
- $\nabla_{\mathbf{z}} J = \left(\frac{\partial \mathbf{h}}{\partial \mathbf{z}}\right)^T \nabla_{\mathbf{h}} J = -\frac{1}{h_y} \nabla h_y(\mathbf{z}) = -(\mathbf{e}_y - \mathbf{h})$

In Practice

- For forward nn
 - Use ReLU, be careful with the learning rate
 - Tryout Leaky ReLU, ELU and Maxout
 - Tryout Tanh with low expectation
 - Never use sigmoid

Gradient Descent for Neural Network

$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

$$W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & \ddots & \cdots \end{bmatrix} \quad b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$\nabla J(\theta) = \left[\cdots, \frac{\partial J(\theta)}{\partial w_{ij}^l}, \cdots, \frac{\partial J(\theta)}{\partial b_i^l}, \cdots \right]^T$$

Algorithm

Initialization: start at θ^0

while($\theta^{(i+1)} \neq \theta^i$)

{

 compute gradient at θ^i

 update parameters

} $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta J(\theta^i)$

To compute the gradients of millions of parameters efficiently, we use **backpropagation**.

Gradient Descent Issue

- After see all training samples, the model can be updated **slowly**.
- It is too expensive to compute the full gradient

Thus, we have stochastic gradient descent (SGD)

Stochastic Gradient Descent

For $t = 1, 2, 3, \dots$ (epoch means one pass over the full training set)

sample $i \in \{1, 2, \dots, n\}$

$$\mathbf{s} = \nabla l(f(\mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)}), y^{(i)}) + \lambda \nabla \Omega(\boldsymbol{\theta}^{(t)})$$

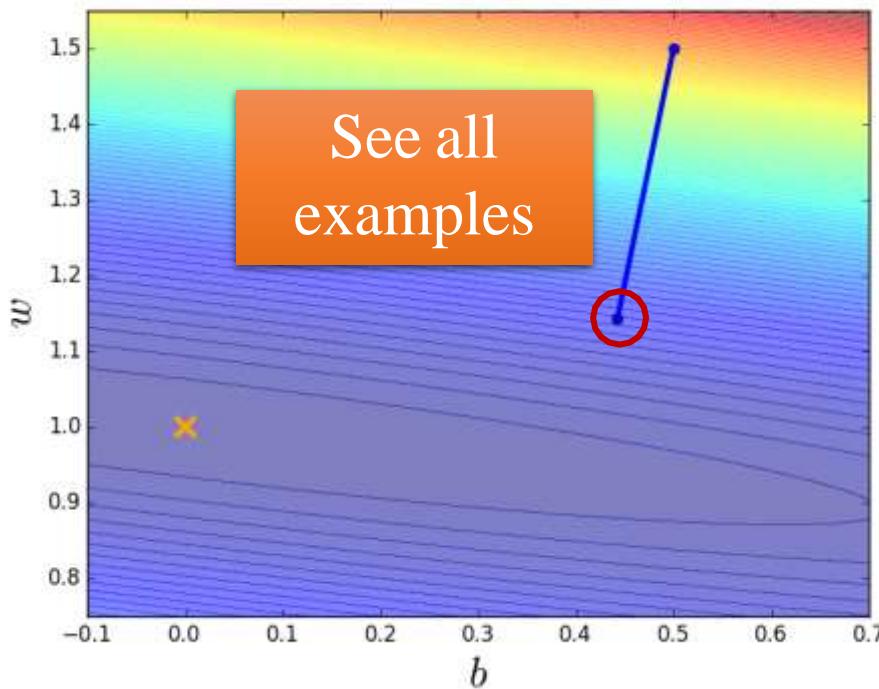
$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \cdot \mathbf{s}$$

- Computing stochastic gradient is much cheaper than full gradient

Gradient Descent v.s. SGD

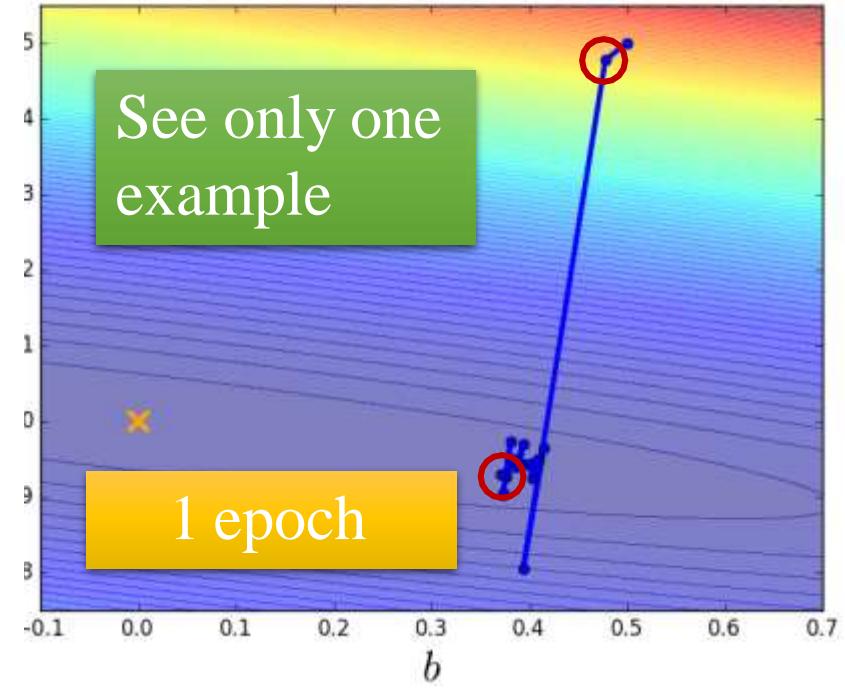
Gradient Descent

Update after seeing all examples



Stochastic Gradient Descent

If there are 20 examples, update 20 times in one epoch.



SGD approaches to the target point faster than gradient descent

Mini-Batch SGD

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

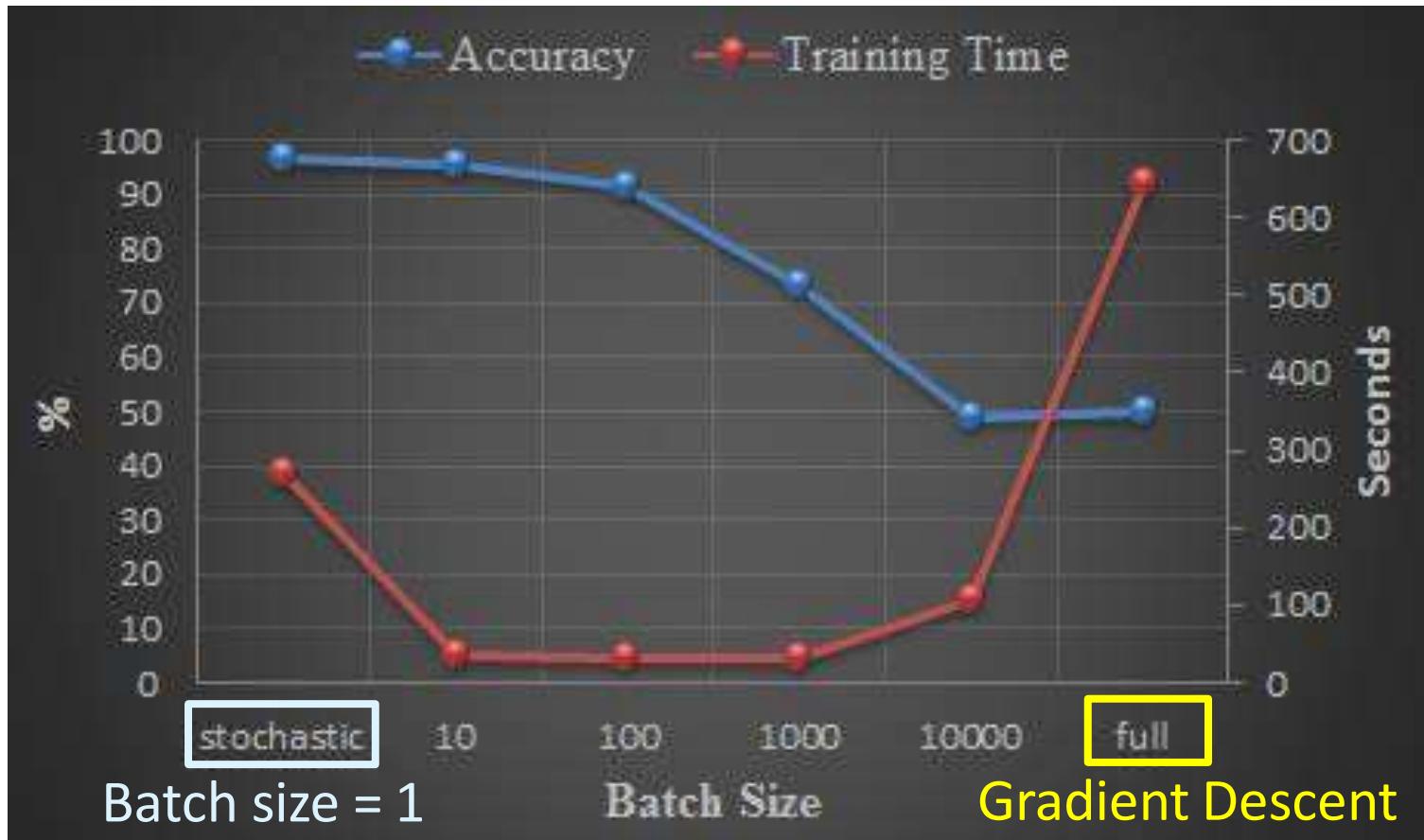
 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

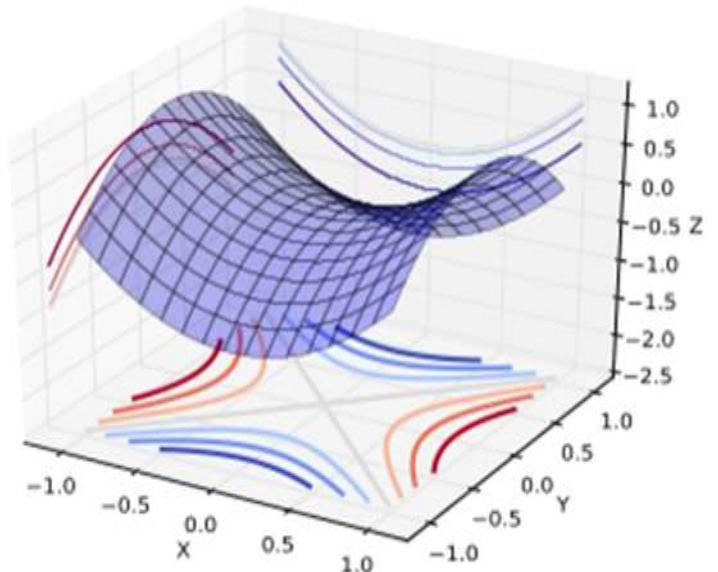
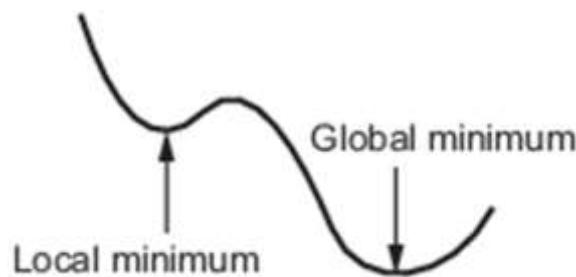
Handwriting Digit Classification

Training speed: mini-batch > SGD > Gradient Descent



Big Issue: Local Optima

- Neural networks has no guarantee for obtaining global optimal solution
- Saddle points



Advanced practical tips (to be presented in the last lecture)

Summary: How to Train Multilayer Neural Nets?

- Define the loss function $l(\cdot)$ properly
- A procedure to compute loss $l(\cdot)$ (*forward propagation*)
- A procedure to compute gradient $\nabla l(\cdot)$ (*back propagation*)
- Regularizer and its gradient $\Omega(\cdot)$ and $\nabla \Omega(\cdot)$
- Perform gradient based optimization method

Forward/Backward Propagation

```
class ComputationalGraph(object):

    #...

    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss

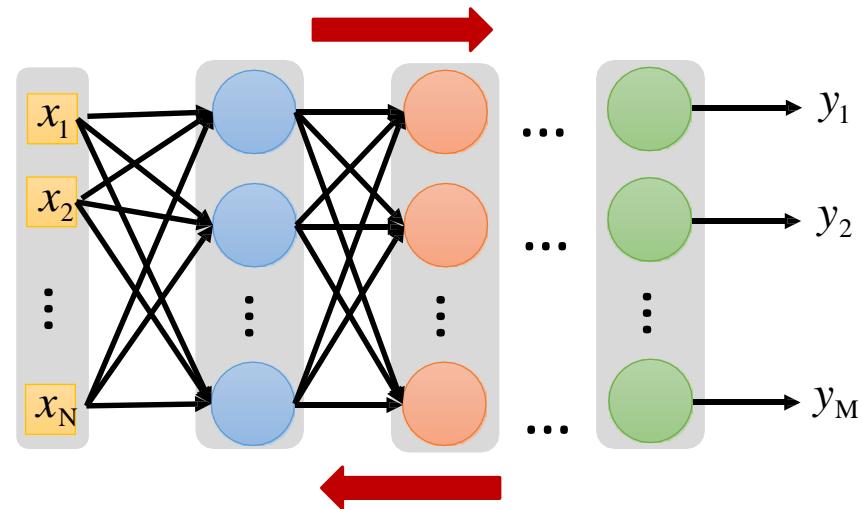
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

Overview

- Model Architectures
 - Artificial neurons
 - Activation function and saturation
 - Feedforward neural nets
- How to train a neural net
 - Loss Function Design
 - Optimization
 - Gradient Descent and Stochastic Gradient Descent
 - Back-propagation

Forward v.s. Back Propagation

- In a feedforward neural network
 - forward propagation
 - from input x to output y information flows forward through the network
 - during training, forward propagation can continue onward until it produces a scalar cost $J(\theta)$
 - back-propagation
 - allows the information from the cost to then flow backwards through the network, in order to compute the **gradient**
 - can be applied to any function



Chain Rule

$$\Delta w \rightarrow \Delta x \rightarrow \Delta y \rightarrow \Delta z$$

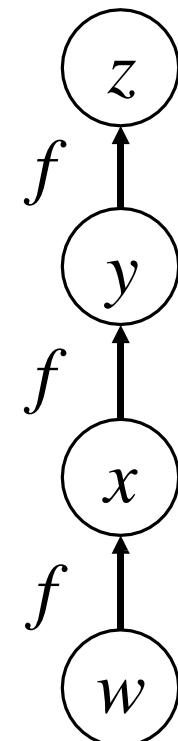
$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w}$$

$$= f'(y) f'(x) f'(w)$$

forward propagation for loss (cost)

$$= f'(f(f(w))) | f'(f(w)) | f'(w)$$

back-propagation for gradient



Gradient Descent for Optimization

$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$

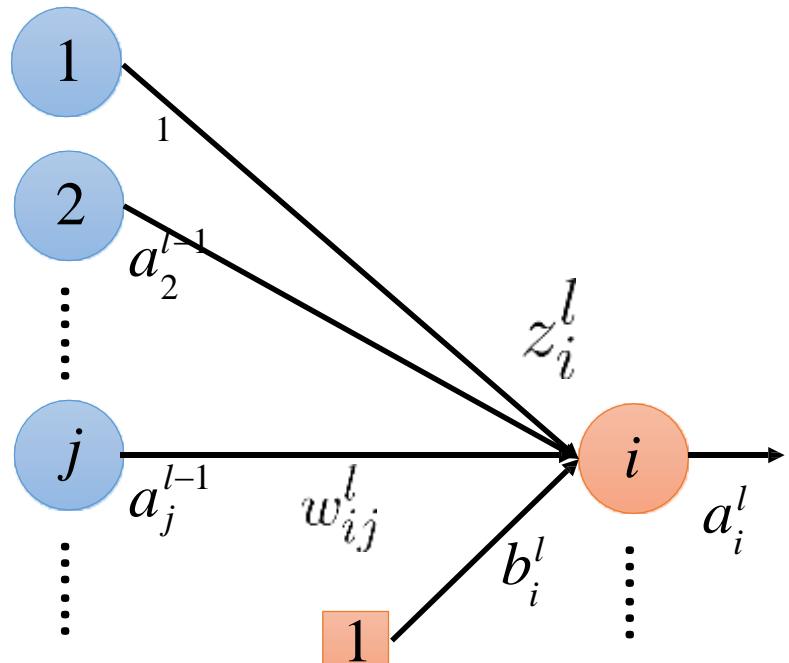
$$W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \dots \\ w_{21}^l & w_{22}^l & \dots \\ \vdots & & \dots \end{bmatrix} \quad b^l = \begin{bmatrix} \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$\nabla J(\theta) = \begin{bmatrix} \vdots \\ \frac{\partial J(\theta)}{\partial w_{ij}^l} \\ \vdots \\ \frac{\partial J(\theta)}{\partial b_i^l} \\ \vdots \end{bmatrix}$$

Algorithm

```
Initialization: start at  $\theta^0$ 
while( $\theta^{(i+1)} \neq \theta^i$ )
{
    compute gradient at  $\theta^i$ 
    update parameters
}
 $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} J(\theta^i)$ 
```

$$\frac{\partial J(\theta)}{\partial w_{ij}^l}$$

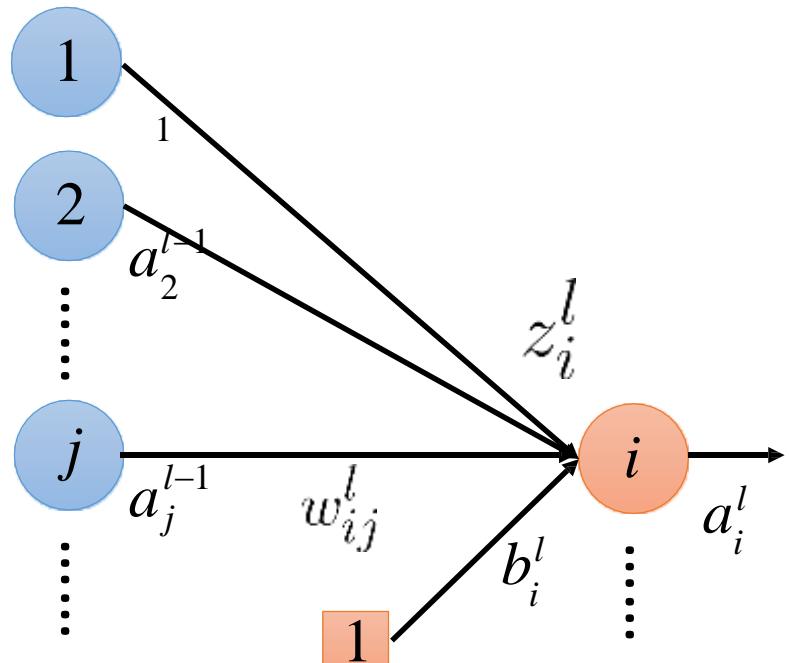


Layer $l-1$

Layer l

$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} \quad (l > 1)$$

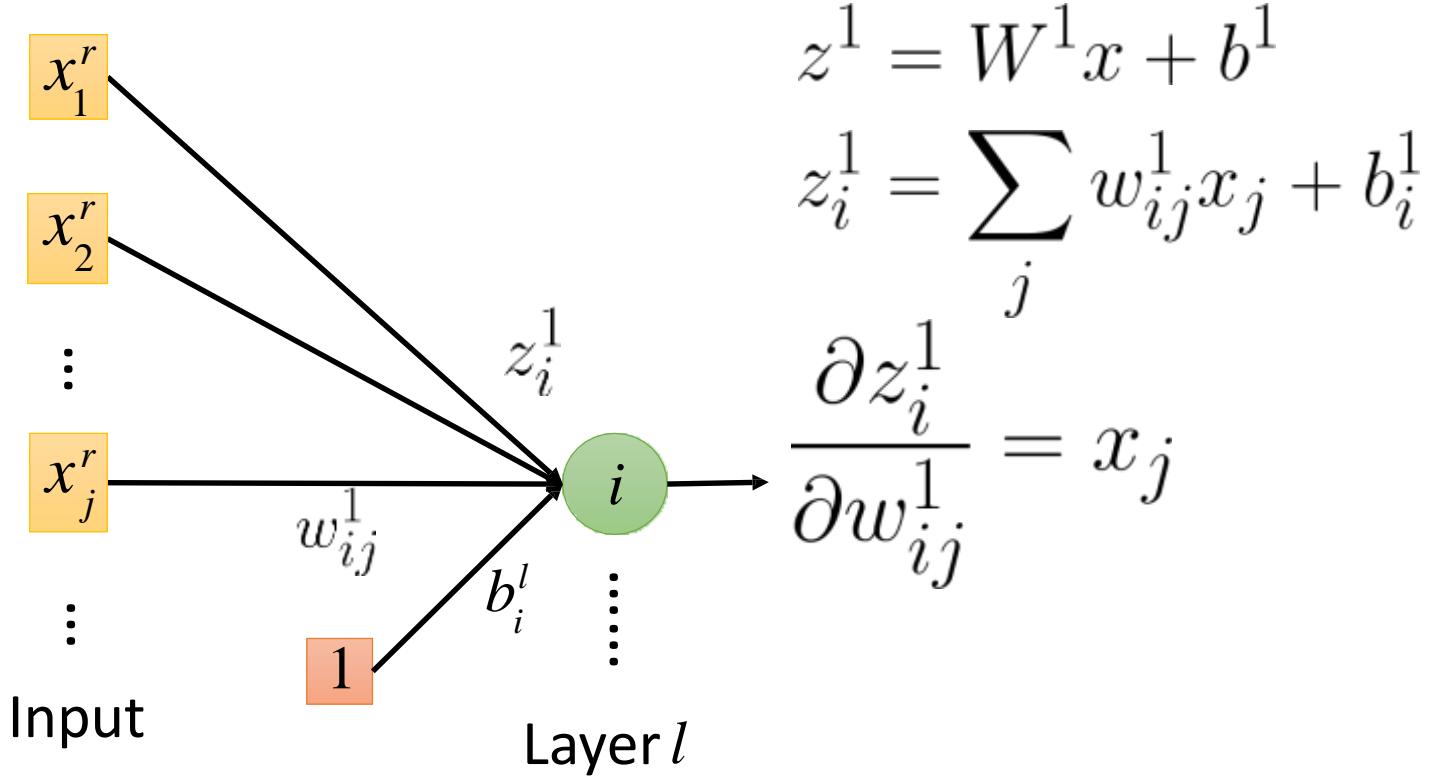


$$z^l = W^l a^{l-1} + b^l$$

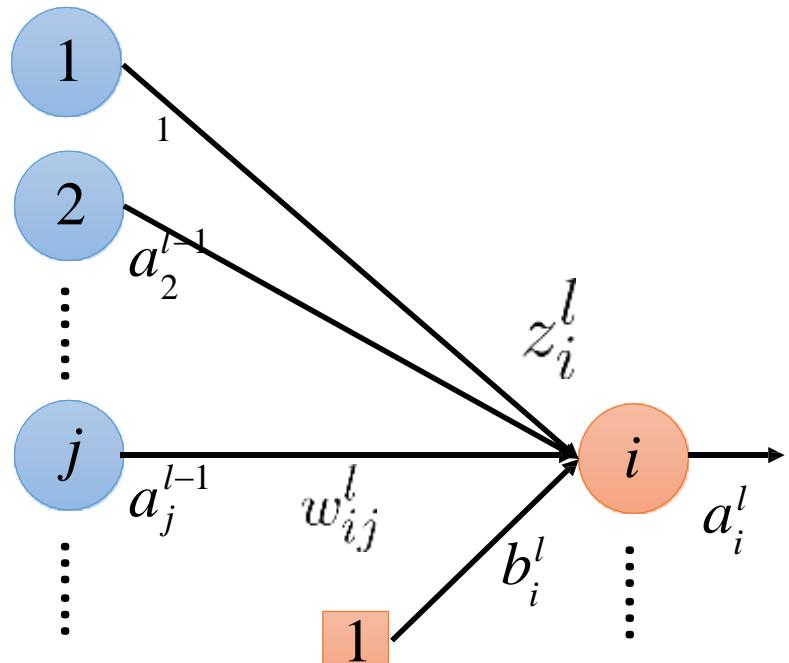
$$z_i^l = \sum_j w_{ij}^l a_j^{l-1} + b_i^l$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = a_j^{l-1}$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} \quad (l = 1)$$



$$\frac{\partial J(\theta)}{\partial w_{ij}^l}$$

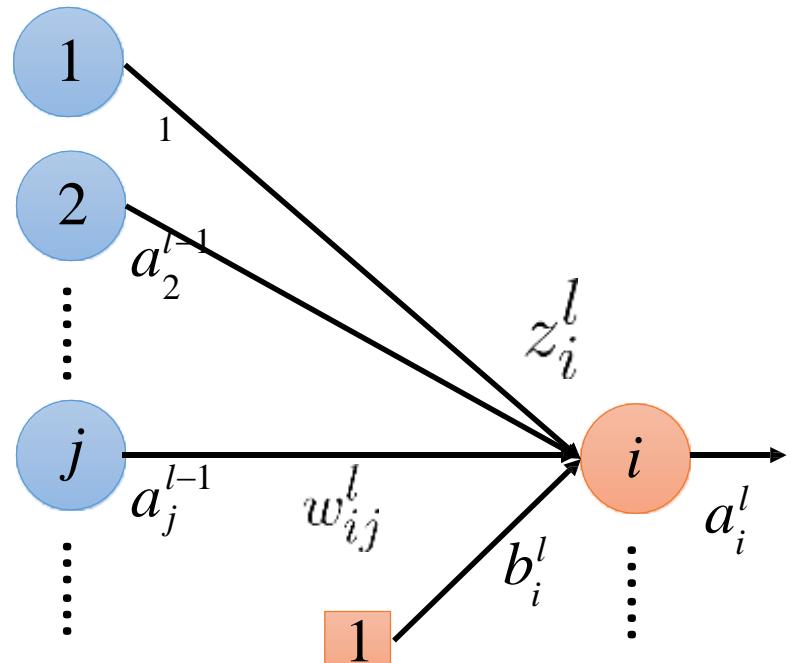


Layer l

$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = \begin{cases} a_j^{l-1}, & l > 1 \\ x_j, & l = 1 \end{cases}$$

$$\frac{\partial J(\theta)}{\partial w_{ij}^l}$$

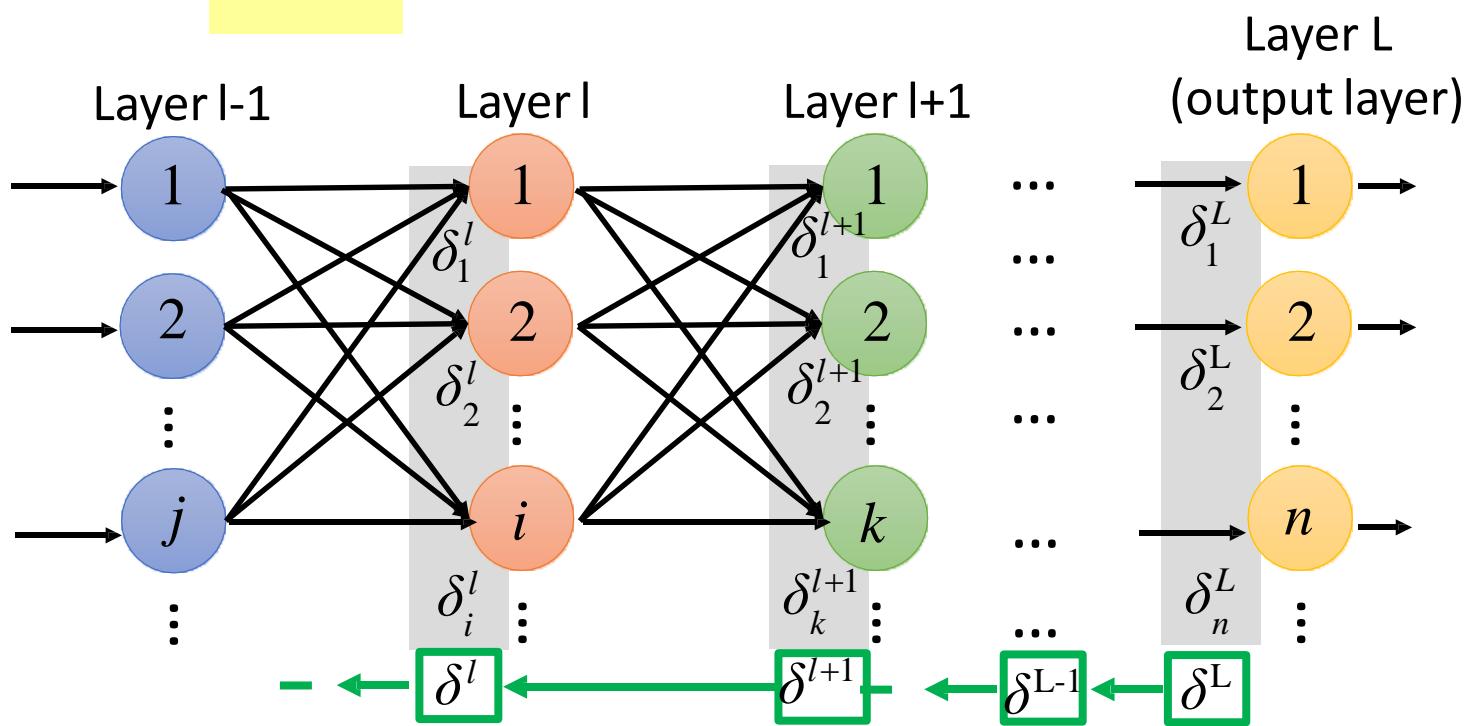


$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\frac{\partial J(\theta)}{\partial z_i^l}$$

$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

δ_i^l : the propagated gradient
 · corresponding to the l -th layer



Idea: computing δ^l layer by layer (from δ^L to δ^1) is more efficient

$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

Idea: from L to 1

- (1) Initialization: compute δ^L
- (2) Compute δ^l based on δ^{l+1}

$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

Idea: from L to 1

(1) Initialization: compute δ^L

(2) Compute δ^l based on δ^{l+1}

$$\delta_i^L = \frac{\partial J}{\partial z_i^L} = \boxed{\frac{\partial J}{\partial y_i}} \frac{\partial y_i}{\partial z_i^L}$$

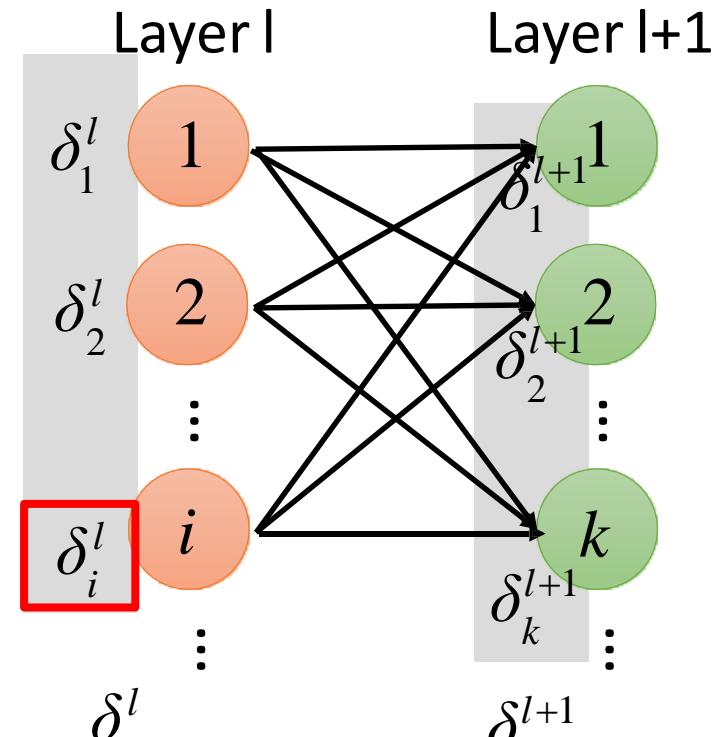
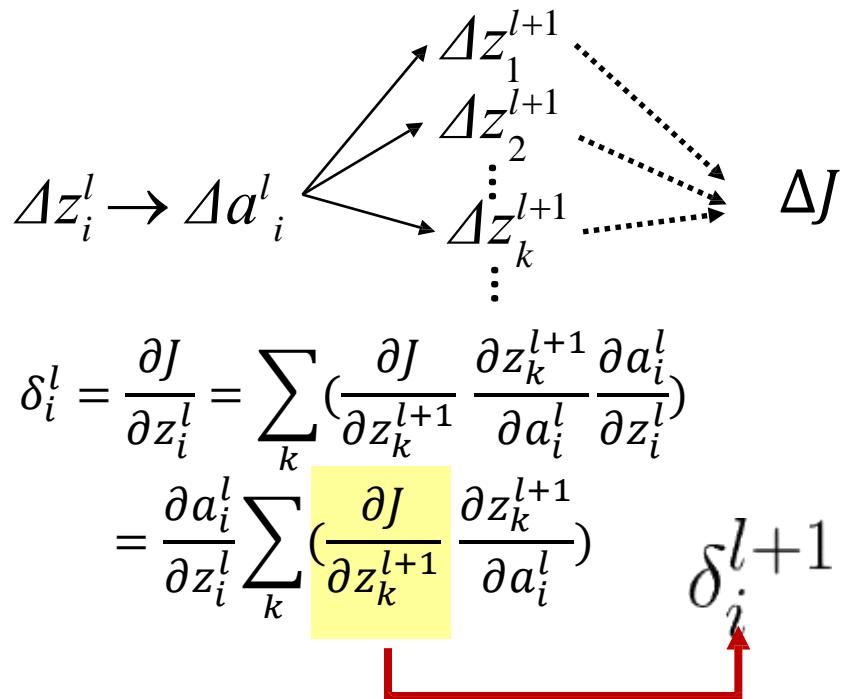
$\frac{\partial J}{\partial y_i}$ depends on the loss function

$$\delta^L = \nabla J(y) \odot \nabla a(z^L)$$

$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

(1) Initialization: compute δ^L

(2) Compute δ^l based on δ^{l+1}

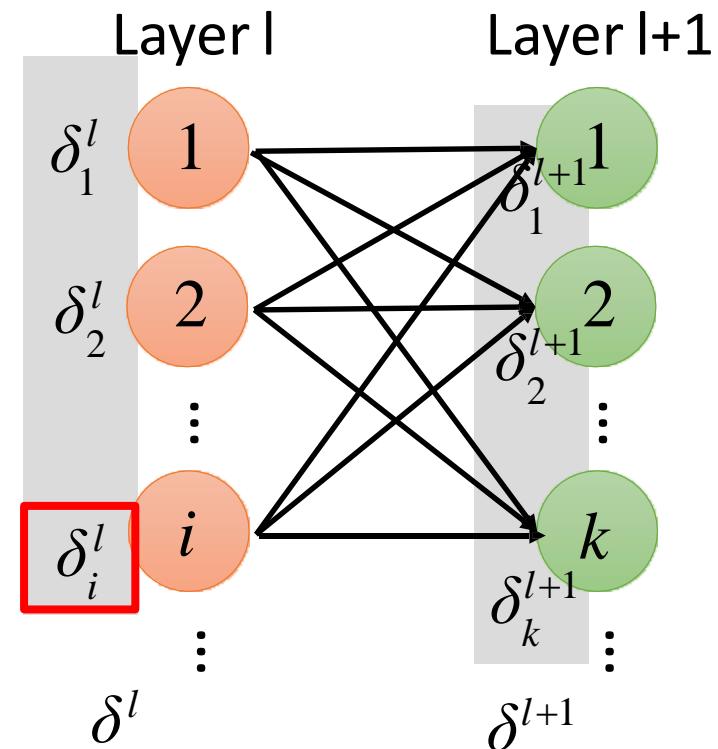


$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

(1) Initialization: compute δ^L

(2) Compute δ^l based on δ^{l+1}

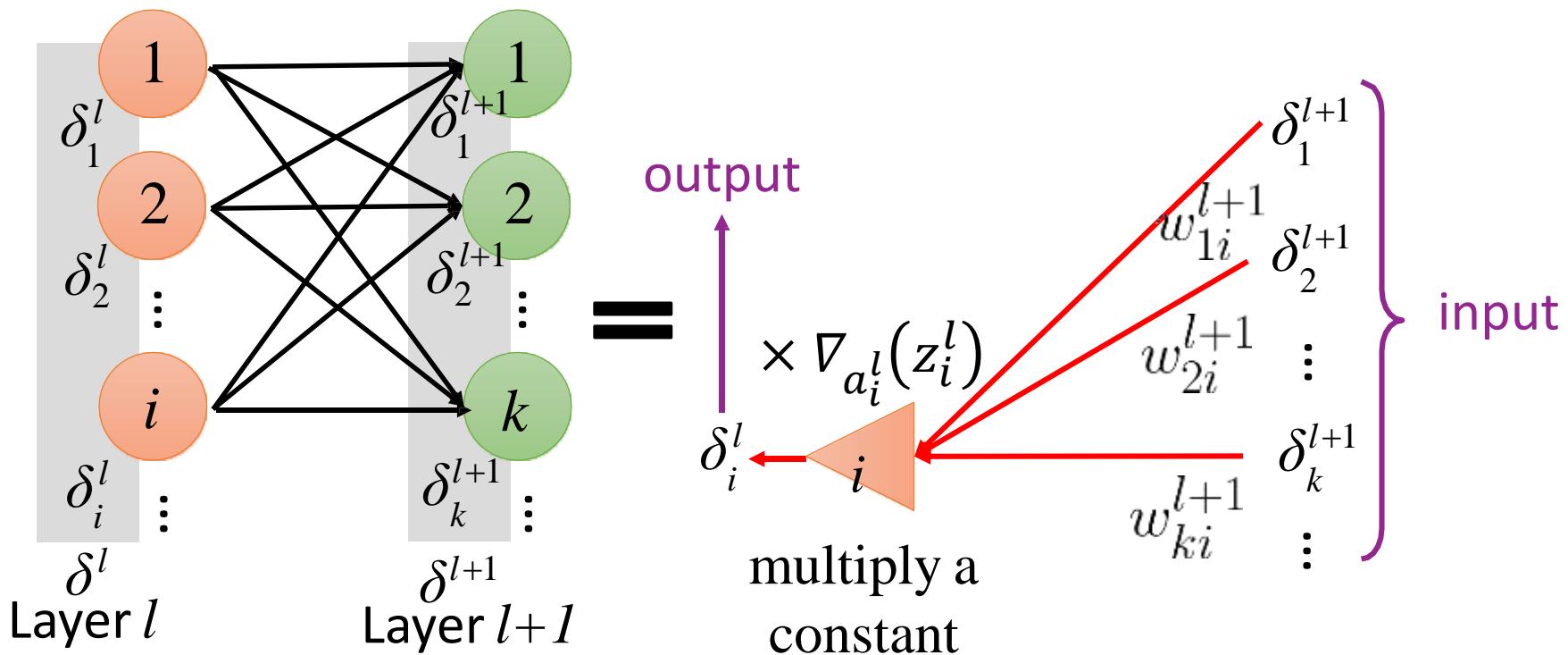
$$\begin{aligned}\delta_i^l &= \frac{\partial a_i^l}{\partial z_i^l} \sum_k \left(\frac{\partial J}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial a_i^l} \right) \\ &= \frac{\partial a_i^l}{\partial z_i^l} \sum_k \delta_k^{l+1} w_{ki}^{l+1} \\ &= \nabla a_i^l(z_i^l) \sum_k \delta_k^{l+1} w_{ki}^{l+1}\end{aligned}$$



$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

Rethink the propagation

$$\delta_i^l = \nabla a_i^l(z_i^l) \sum_k \delta_k^{l+1} w_{ki}^{l+1}$$



$$\delta^l = \nabla a(z^l) \odot (W^{l+1})^T \delta^{l+1}$$

$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

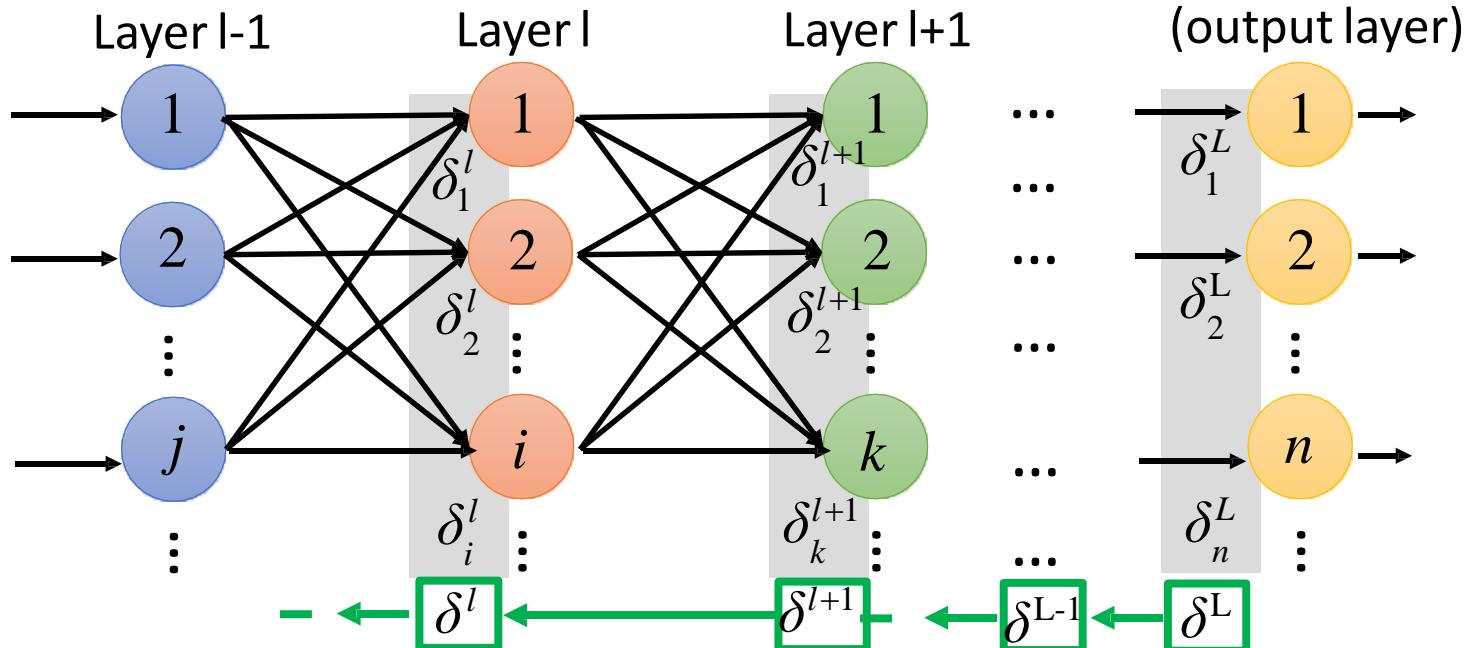
$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

Idea: from L to 1

- (1) Initialization: compute δ^L
- (2) Compute δ^l based on δ^{l+1}

$$\delta^L = \nabla J(y) \odot \nabla a(z^L)$$

$$\delta^l = \nabla a(z^l) \odot (W^{l+1})^T \delta^{l+1}$$



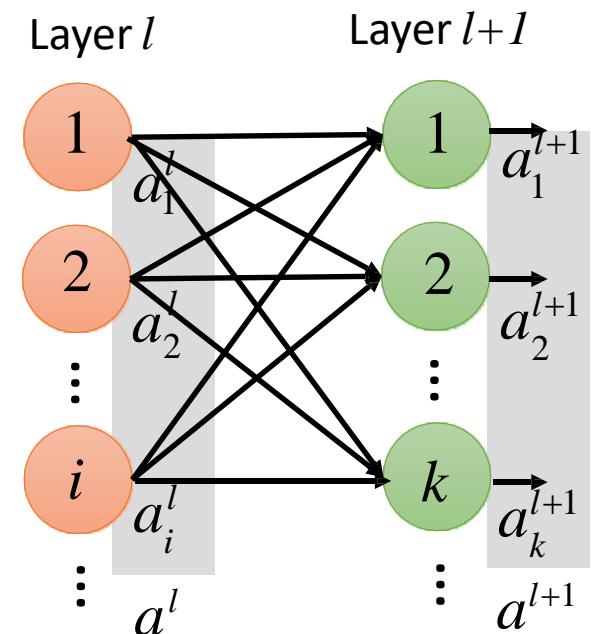
Backpropagation

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = \begin{cases} a_j^{l-1}, & l > 1 \\ x_j, & l = 1 \end{cases}$$

Forward Pass

$$\begin{aligned} z^1 &= W^1 x + b^1 & a^1 &= \sigma(z^1) \\ &\vdots && \\ z^l &= W^l a^{l-1} + b^l & a^l &= \sigma(z^l) \\ &\vdots && \end{aligned}$$

$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$



Backpropagation

$$\frac{\partial J(\theta)}{\partial z_i^l} = \delta_i^l$$

Backward Pass

$$\delta^L = \nabla J(y) \odot \nabla a(z^L)$$

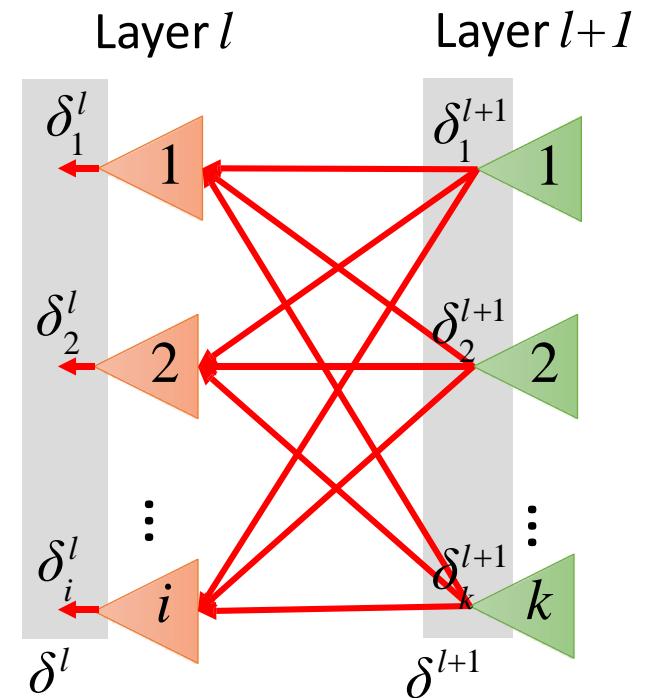
$$\delta^{L-1} = \nabla a(z^l) \odot (W^{l+1})^T \delta^{l+1}$$

⋮

$$\delta^l = \nabla a(z^l) \odot (W^{l+1})^T \delta^{l+1}$$

⋮

$$\frac{\partial J(\theta)}{\partial w_{ij}^l} = \boxed{\frac{\partial J(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}}$$



Reading Materials

- Automatic Differentiation in Machine Learning: a Survey (2015)

Summary: How to Train Multilayer Neural Nets?

- Define the loss function $l(\cdot)$ properly
- A procedure to compute loss $l(\cdot)$ (*forward propagation*)
- A procedure to compute gradient $\nabla l(\cdot)$ (*back propagation*)
- Regularizer and its gradient $\Omega(\cdot)$ and $\nabla \Omega(\cdot)$
- Perform gradient based optimization method

Summary

- Model Architectures
 - Artificial neurons
 - Activation function and saturation
 - Feedforward neural nets
- How to train a neural net
 - Loss Function Design
 - Optimization
 - Gradient Descent and Stochastic Gradient Descent
 - Back-propagation