

CS 546 – Advanced Topics in NLP

Dilek Hakkani-Tür



UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN



Siebel School of
Computing
and Data Science

Topics for Today

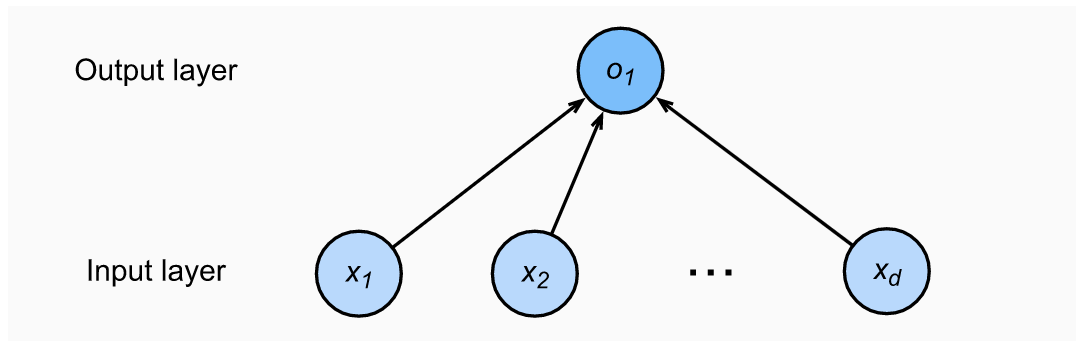


- Gradient Descent
- Softmax Regression
- Multi-layer Perceptron

From Linear Regression to Neural Networks



- Linear regression is a single-layer neural network, consisting of just a single neuron!
- Fully connected layer (also called *dense* layer): every input is connected to every output.



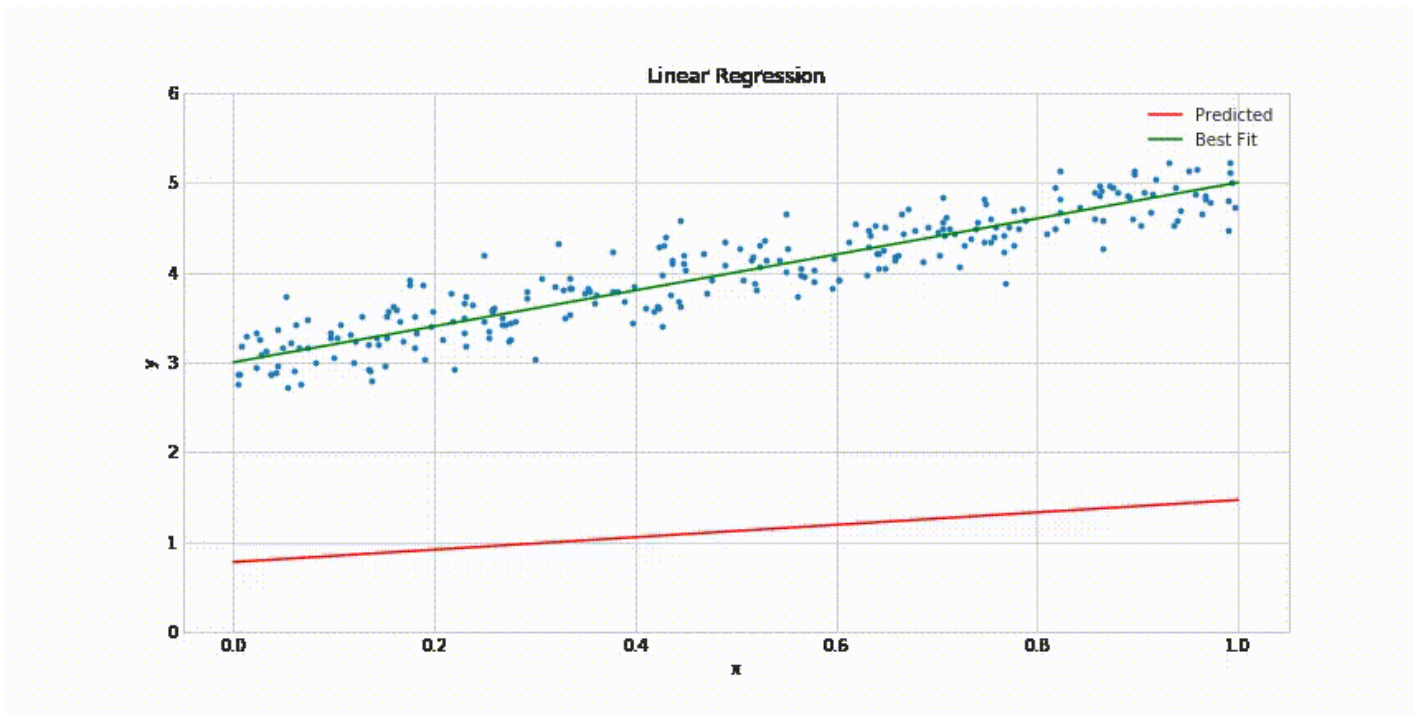
Gradient Descent



- Seen the solution to linear regression last Thursday!
- Even when we cannot solve the models analytically, we can still train models effectively by **iteratively reducing the error by updating the parameters in the direction that incrementally lowers the loss function**:
 - Derivative of the true loss (i.e., average of the losses computed over all examples in the training data).
 - Expressing the updates mathematically:

$$(\mathbf{w}, b) \leftarrow (\mathbf{w}, b) - \frac{\eta}{|K|} \sum_{i \in K} \partial_{(\mathbf{w}, b)} L^{(i)}(\mathbf{w}, b)$$

Case Study: $f(x) = mx + b$



Example from the web

Gradient Descent (cont.)



- Steps of the algorithm:
 1. initialize the values of the model parameters, typically at random
 2. iteratively update the parameters in the direction of the negative gradient.

Notation:

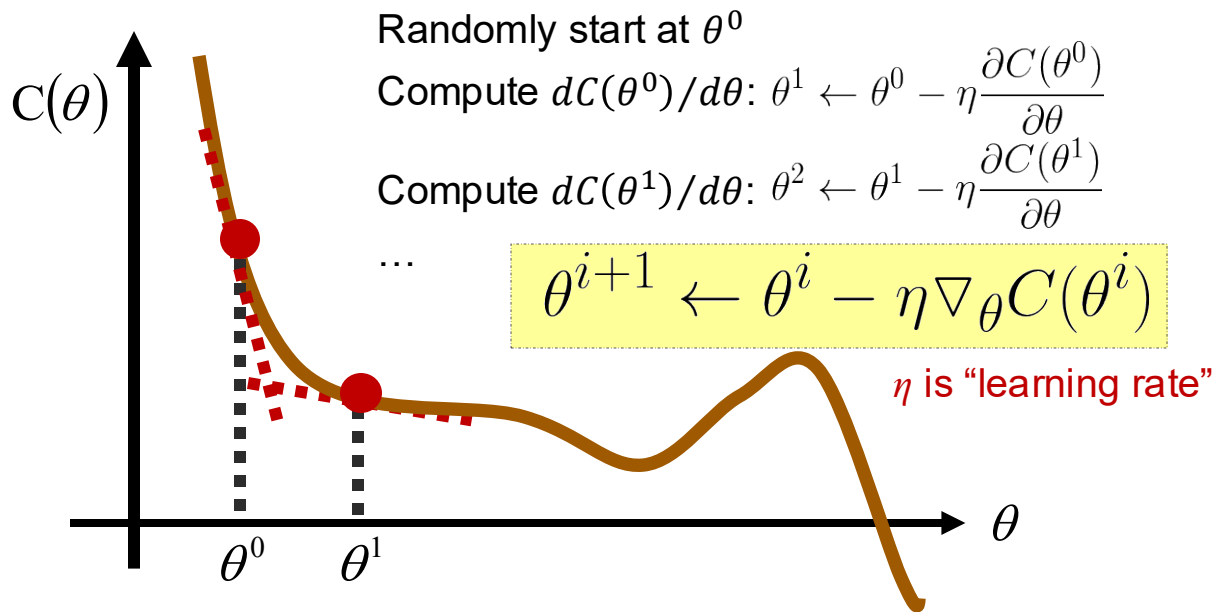
θ : model parameters (\mathbf{w}, b)

$C(\theta)$: $L(\mathbf{w}, b)$

Gradient Descent (cont.)



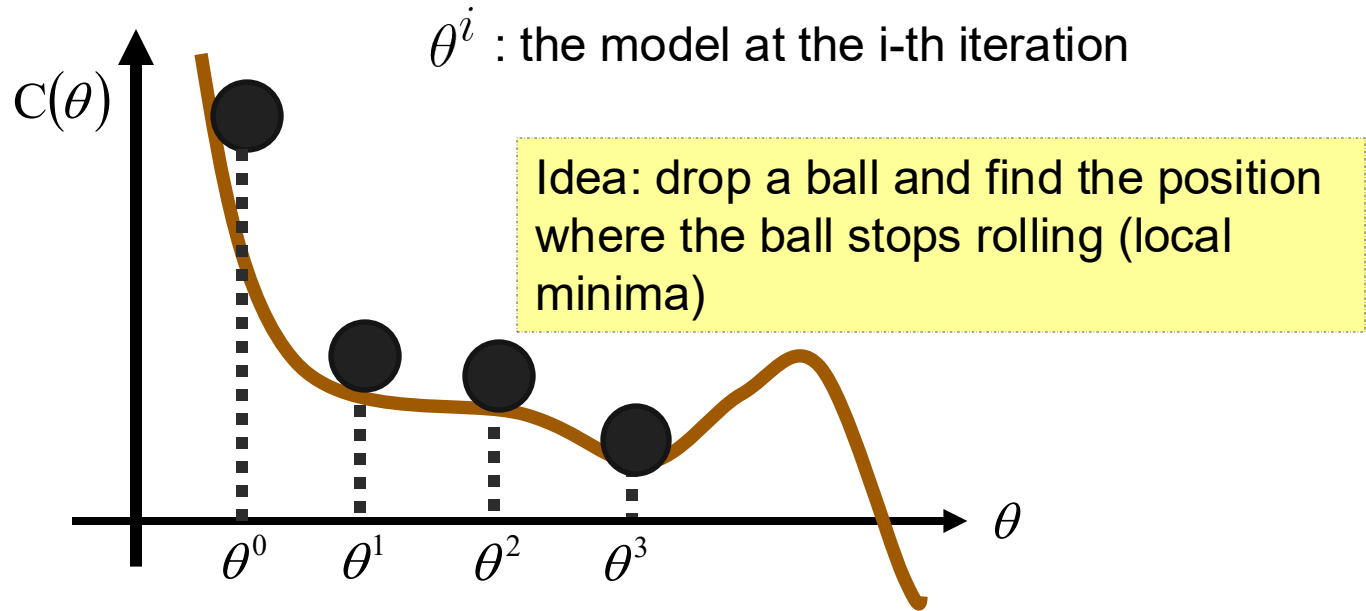
- Assume that θ has only one variable



Gradient Descent (cont.)



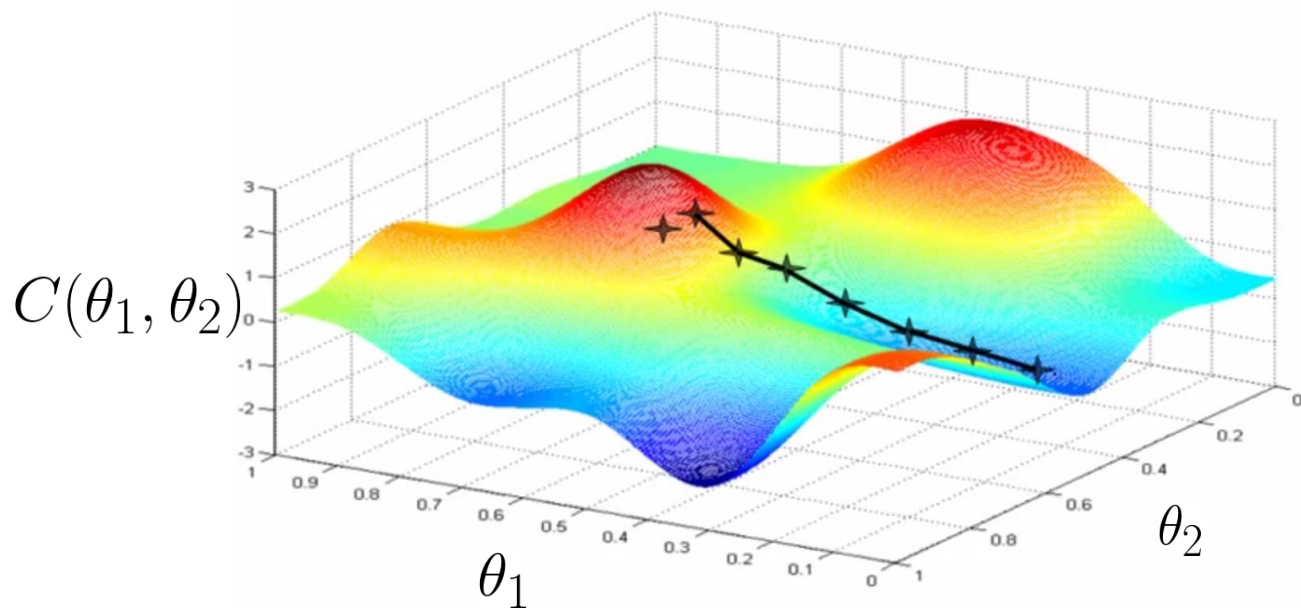
- Assume that θ has only one variable



Gradient Descent (cont.)



- Assume that θ has two variables $\{\theta_1, \theta_2\}$



Gradient Descent (cont.)



- Assume that θ has two variables $\{\theta_1, \theta_2\}$

- Randomly start at θ^0 : $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

- Compute the gradients of $C(\theta)$ at θ^0 : $\nabla_{\theta} C(\theta^0) = \begin{bmatrix} \frac{\partial C(\theta_1^0)}{\partial \theta_1} \\ \frac{\partial C(\theta_2^0)}{\partial \theta_2} \end{bmatrix}$

- Update parameters:

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial C(\theta_1^0)}{\partial \theta_1} \\ \frac{\partial C(\theta_2^0)}{\partial \theta_2} \end{bmatrix}$$

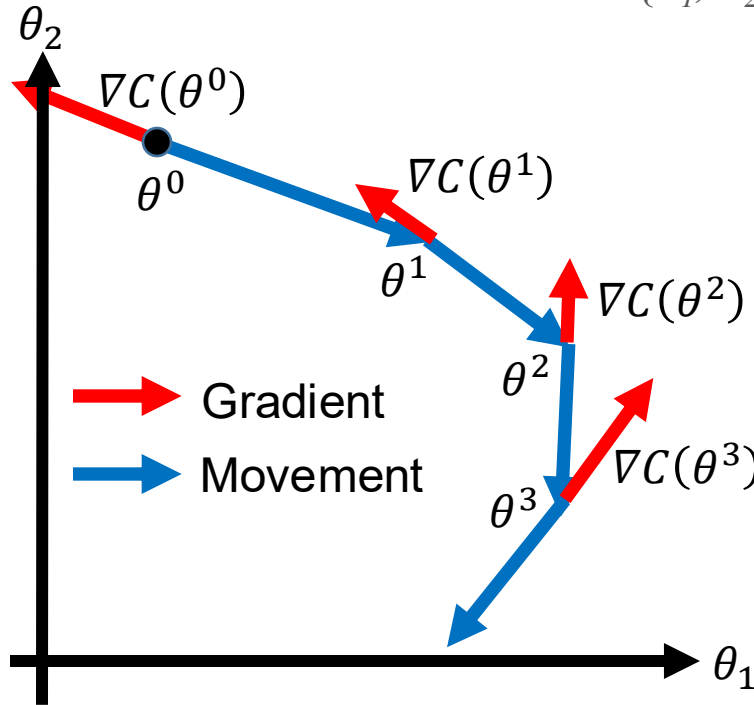
$$\theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} C(\theta^i)$$

- Compute the gradients of $C(\theta)$ at θ^1 : $\nabla_{\theta} C(\theta^1) = \begin{bmatrix} \frac{\partial C(\theta_1^1)}{\partial \theta_1} \\ \frac{\partial C(\theta_2^1)}{\partial \theta_2} \end{bmatrix}$
 - ...

Gradient Descent (cont.)



- Assume that θ has two variables $\{\theta_1, \theta_2\}$



Algorithm

Initialization: start at θ^0
while($\theta^{(i+1)} \neq \theta^i$)

{

compute gradient at θ^i
update parameters
 $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} C(\theta^i)$

}

An Issue with Gradient Descent



$$\theta^{i+1} = \theta^i - \eta \nabla C(\theta^i)$$

Training Data
 $\{(x_1, \hat{y}_1), (x_2, \hat{y}_2), \dots\}$

$$C(\theta) = \frac{1}{K} \sum_k \|f(x_k; \theta) - \hat{y}_k\| = \frac{1}{K} \sum_k C_k(\theta)$$

$$\nabla C(\theta^i) = \frac{1}{K} \sum_k \nabla C_k(\theta^i)$$

After seeing all training samples, the model can be updated → slow

Stochastic Gradient Descent (SGD)



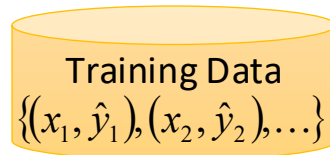
- Gradient Descent

$$\theta^{i+1} = \theta^i - \eta \nabla C(\theta^i) \quad \nabla C(\theta^i) = \frac{1}{K} \sum_k \nabla C_k(\theta^i)$$

- Stochastic Gradient Descent (SGD)

- Pick a training sample x_k

$$\theta^{i+1} = \theta^i - \eta \nabla C_k(\theta^i)$$



- If all training samples have the same probability to be picked, then

$$E[\nabla C_k(\theta^i)] = \frac{1}{K} \sum_k \nabla C_k(\theta^i)$$

The model can be updated after seeing one training sample → faster

Stochastic Gradient Descent (SGD) (cont.)



- When running SGD, the model starts θ^0

pick x_1 $\theta^1 = \theta^0 - \eta \nabla C_1(\theta^0)$

pick x_2 $\theta^2 = \theta^1 - \eta \nabla C_2(\theta^1)$

\vdots

pick x_k $\theta^k = \theta^{k-1} - \eta \nabla C_k(\theta^{k-1})$

\vdots

pick x_K $\theta^K = \theta^{K-1} - \eta \nabla C_K(\theta^{K-1})$

Training Data

$\{(x_1, \hat{y}_1), (x_2, \hat{y}_2), \dots\}$

see all
training
samples once

→ one epoch

pick x_1 $\theta^{K+1} = \theta^K - \eta \nabla C_1(\theta^K)$

Mini-Batch SGD



- Batch Gradient Descent

Use all K samples in each iteration

$$\theta^{i+1} = \theta^i - \eta \frac{1}{K} \sum_k \nabla C_k(\theta^i)$$

- Stochastic Gradient Descent (SGD)

- Pick a training sample x_k

Use 1 sample in each iteration

$$\theta^{i+1} = \theta^i - \eta \nabla C_k(\theta^i)$$

- Mini-Batch SGD

- Pick a set of B training samples as a batch b

Use all B samples in each iteration

B is the “batch size”

$$\theta^{i+1} = \theta^i - \eta \frac{1}{B} \sum_{x_k \in b} \nabla C_k(\theta^i)$$

Topics for Today



- Gradient Descent
- Softmax Regression
- Multi-layer Perceptron

Softmax Regression



- Regression is useful for “how much?” and “how many?” questions, where the target can take a real value.
- Many NLP problems are about “which one?”
 - Sentiment classification: positive, negative, neutral
 - Sentence boundary detection from text: is the punctuation mark defining a sentence boundary or not
 - What is the part-of-speech tag of the word “word” in this sentence? NOUN, VERB, etc.?
- We can represent target labels with one-hot encodings as well. For example, for sentiment classification:

$$y \in \{(1,0,0), (0,1,0), (0,0,1)\}.$$


Positive Negative Neutral

A Network Architecture for Multi-Class Problems



- To estimate the conditional probabilities associated with each class, we need a model with multiple outputs, one per class.
- As many linear functions as we have outputs.
- Example: Assume our input is represented with 4 features, then to compute the logits o_1, o_2, o_3 , for each output, we need:

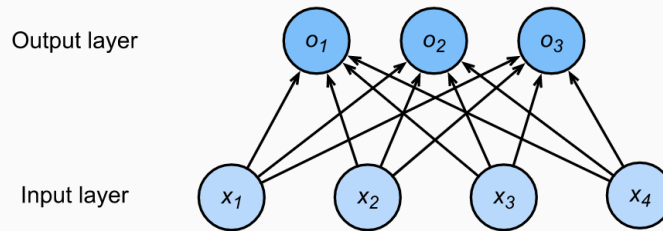
$$o_1 = x_1w_{11} + x_2w_{12} + x_3w_{13} + x_4w_{14} + b_1$$

$$o_2 = x_1w_{21} + x_2w_{22} + x_3w_{23} + x_4w_{24} + b_2$$

$$o_3 = x_1w_{31} + x_2w_{32} + x_3w_{33} + x_4w_{34} + b_3$$

$$\mathbf{o} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

3 x 4 matrix



Multi-Class Problems



- Need to interpret model outputs as probabilities and optimize our parameters to produce probabilities that maximize the likelihood of the observed data
- To generate predictions, we can set a threshold or choose the *argmax* of the predicted probabilities
- Interpret the logits o directly as our outputs of interest, however:
 - Nothing constrains these numbers to sum to 1.
 - Depending on the inputs, they can take negative values.

Softmax



- Transforms logits such that they become nonnegative and sum to 1, while requiring that the model remains differentiable.

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \quad \text{where} \quad \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}$$

- $\hat{y}_1 + \hat{y}_2 + \hat{y}_3 = 1$,
- $0 \leq \hat{y}_i \leq 1$ for all i , and
- The ordering of the logits has not changed; hence we can still pick the output using `argmax`!

Maximum Likelihood Estimation



- The softmax function gives us a vector $\hat{\mathbf{y}}$, interpreted as estimated conditional probabilities of each class given the input x , e.g., $y_1^{\hat{}} = P(y=\text{cat}|\mathbf{x})$.

$$P(Y | X) = \prod_{i=1}^n P(y^{(i)} | x^{(i)})$$

$$-\log P(Y | X) = \sum_{i=1}^n -\log P(y^{(i)} | x^{(i)})$$

- Maximizing $P(Y|X)$ (and thus equivalently minimizing $-\log P(Y|X)$) corresponds to predicting the label well!

Cross-Entropy Loss and Softmax



$$l(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{j=1}^q y_j \log \hat{y}_j.$$

q : the number of classes

- If we add \mathbf{o} into the definition of the loss l and use the definition of the softmax we obtain:

$$\begin{aligned} l(\mathbf{y}, \hat{\mathbf{y}}) &= - \sum_{j=1}^q y_j \log \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} \\ &= \sum_{j=1}^q y_j \log \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j \\ &= \log \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j. \end{aligned}$$

- If we compute its derivative with respect to o_i , we get:

$$\partial_{o_j} l(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} - y_j = \text{softmax}(\mathbf{o})_j - y_j.$$

The gradient is the difference between the probability assigned to the true class by our model. This makes computing gradients very easy in practice!

Topics for Today

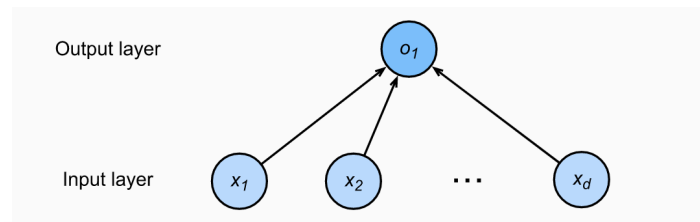
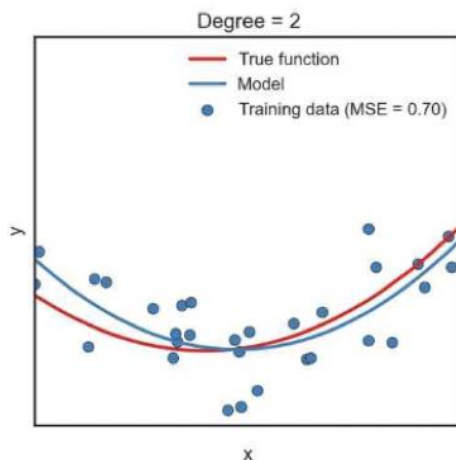
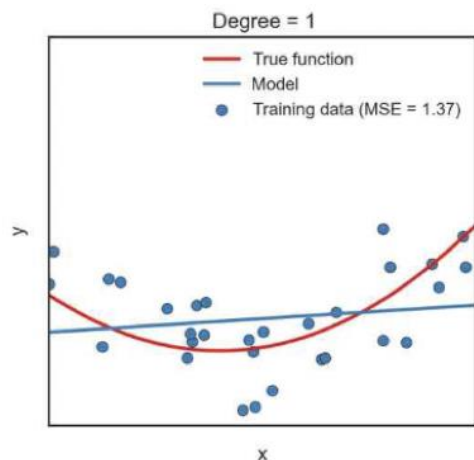


- Gradient Descent
- Softmax Regression
- Multi-layer Perceptron

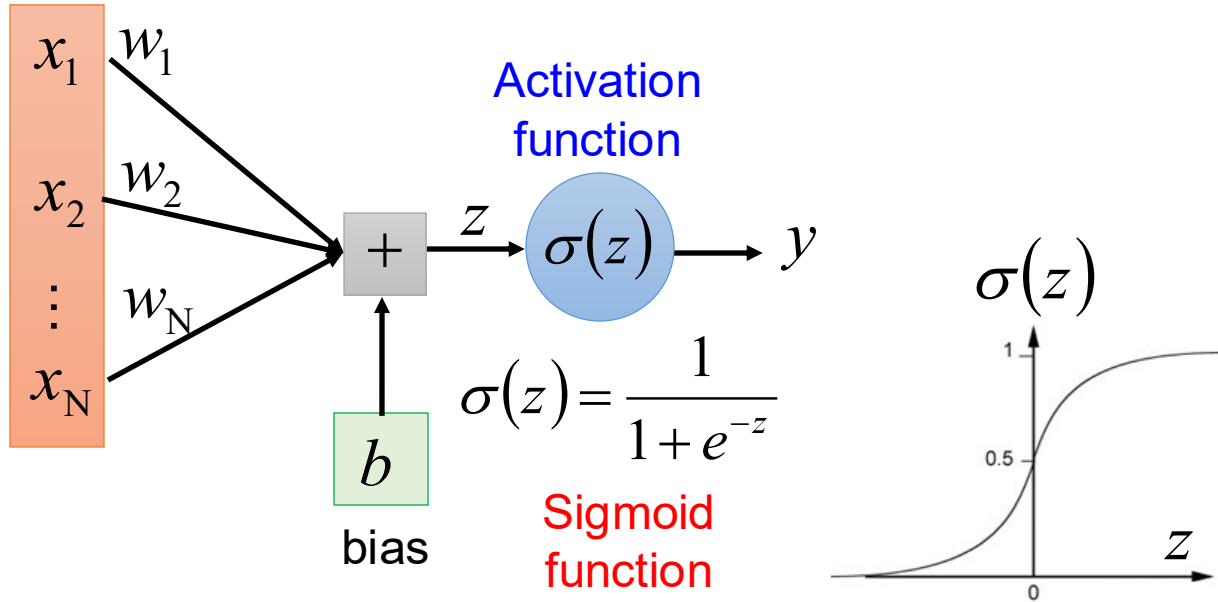
Non-Linearity



- Linear regression inputs directly to our outputs via a single linear transformation.
- But what if our data is not linear?



Perceptron - A Single Neuron



Each neuron is a very simple function

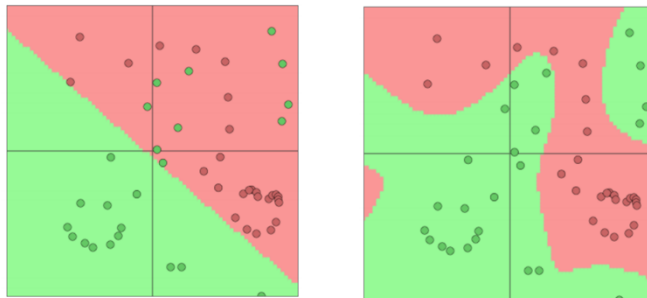
Why non-linearity?



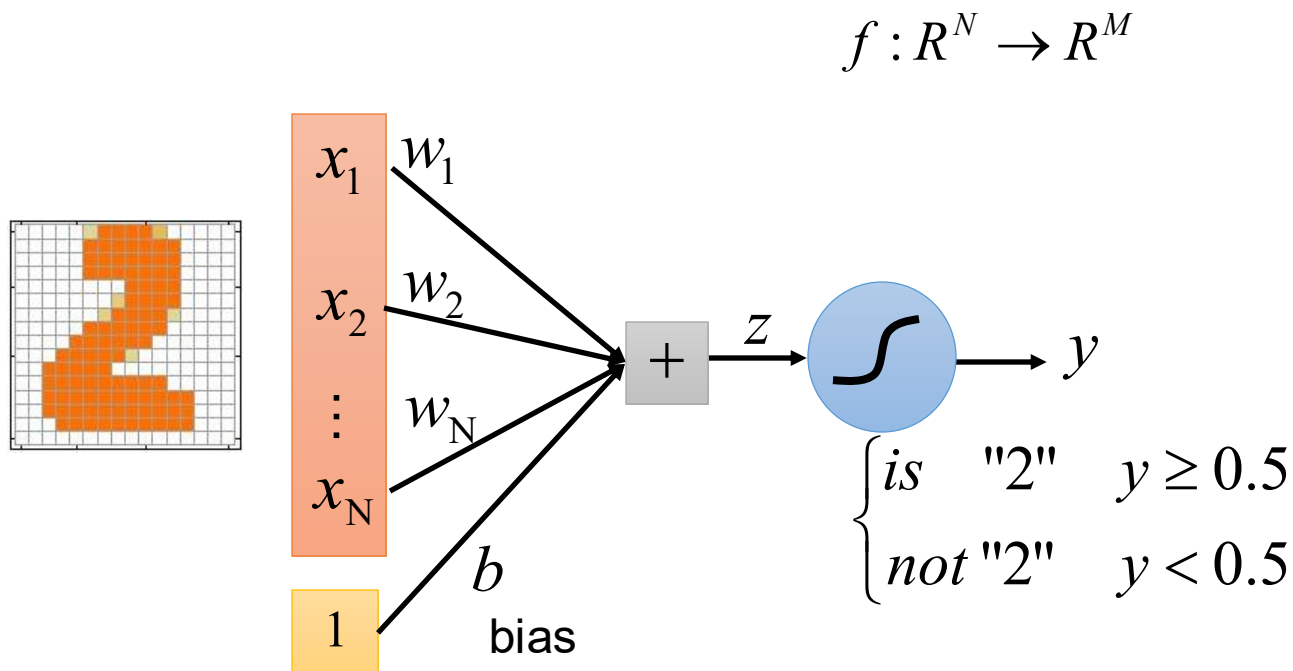
- Function approximation
 - ***Without non-linearity***, deep neural networks work the same as linear transform

$$W_1(W_2 \cdot x) = (W_1 W_2)x = Wx$$

- ***With non-linearity***, networks with more layers can approximate more complex functions



Perceptron - A Single Neuron

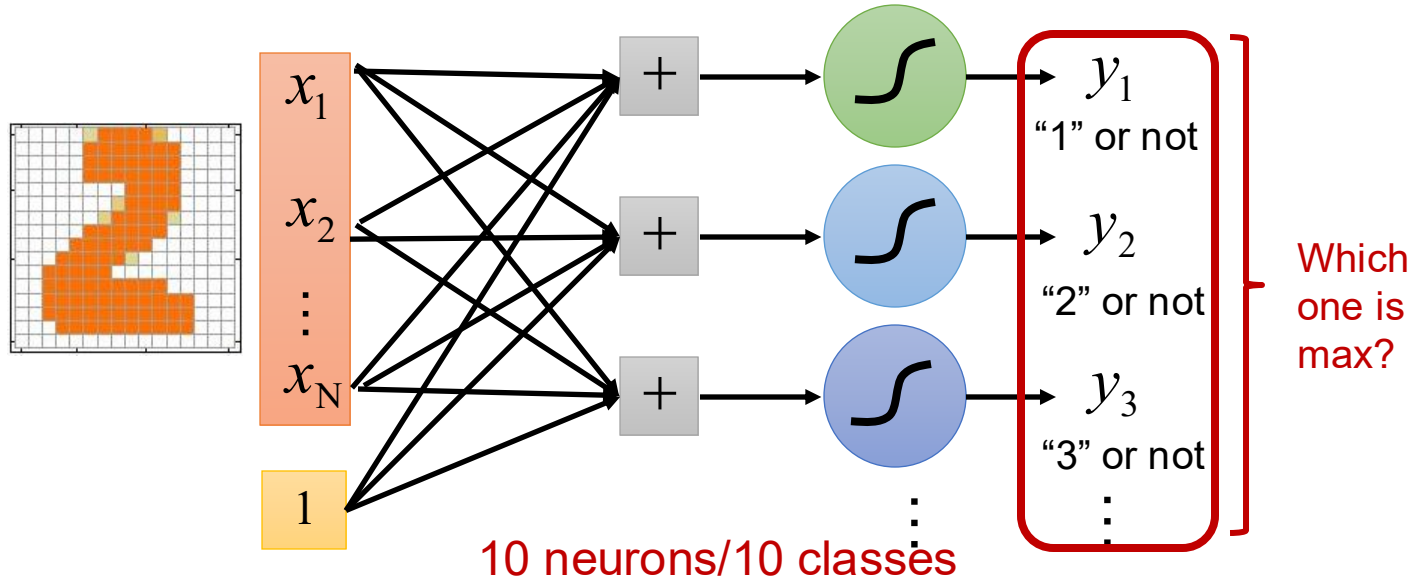


A single neuron can only handle binary classification

A Layer of Neurons

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

- Handwriting digit classification

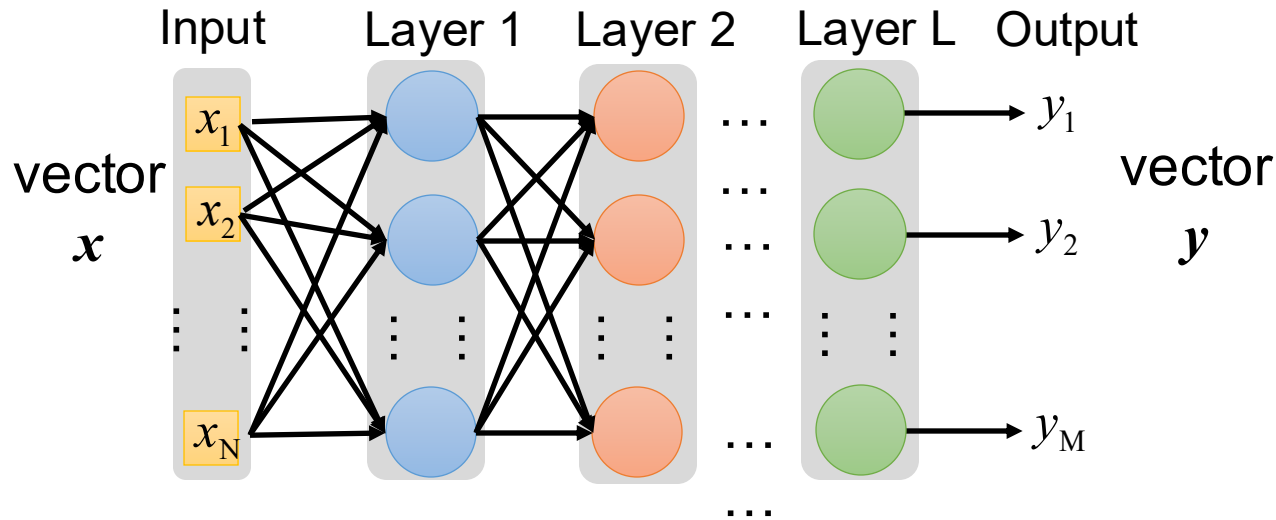


A layer of neurons can handle multiple possible output,
and the result depends on the max one

Deep Neural Networks (DNNs)

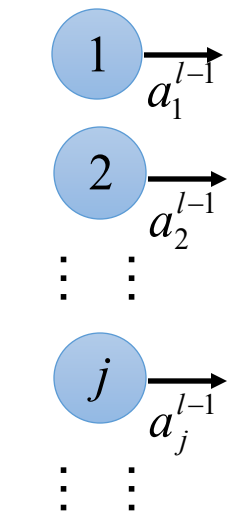


$$f: R^N \rightarrow R^M$$

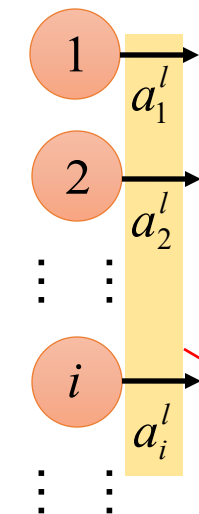


Deep NN: multiple hidden layers

DNNs – Notation



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



Layer l
 N_l nodes

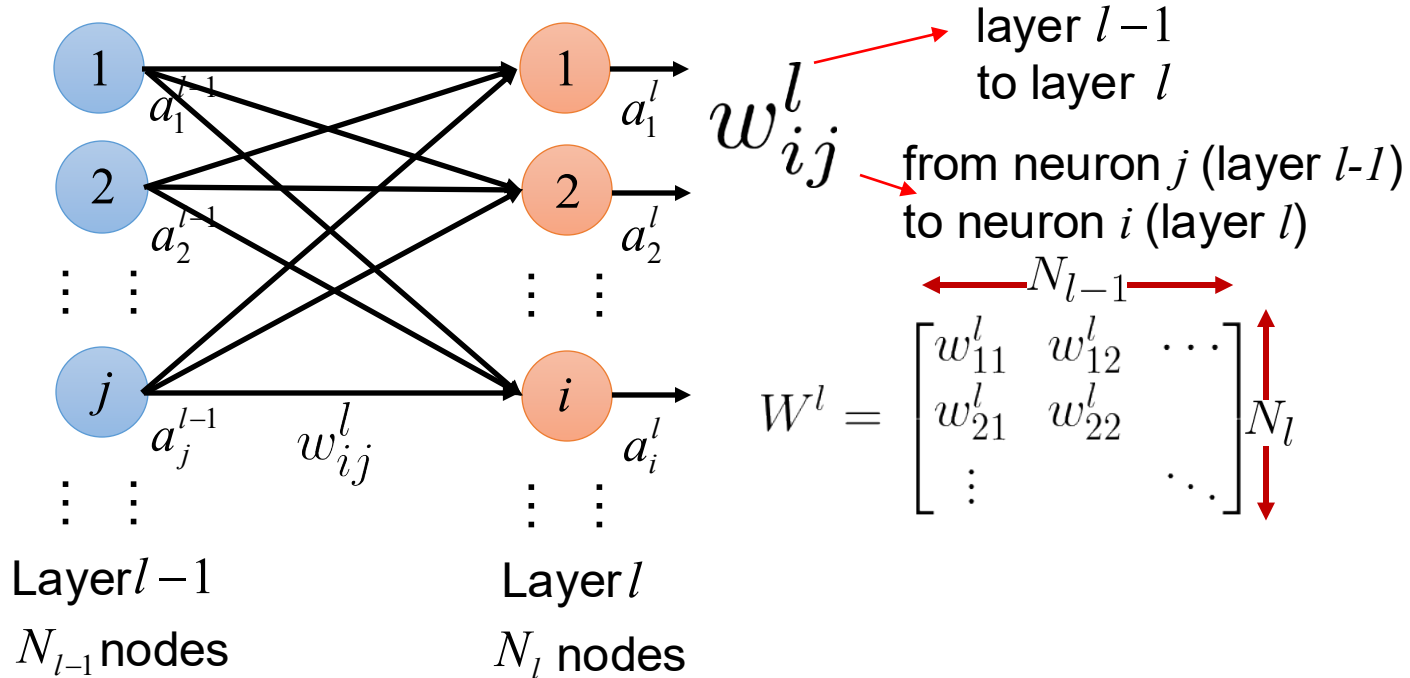
Output of a neuron:

a_i^l → layer l
 a_i^l → neuron i

$$a^l = \begin{bmatrix} \vdots \\ a_i^l \\ \vdots \end{bmatrix}$$

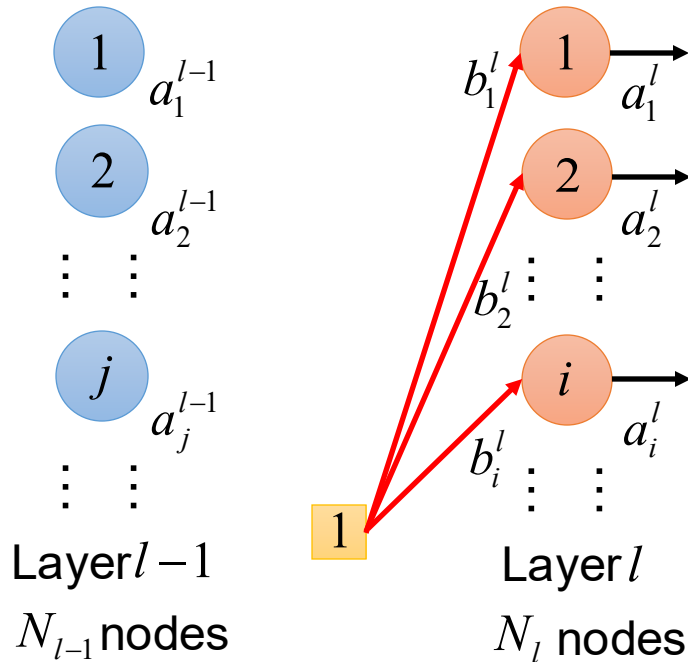
output of one layer → a vector

DNNs – Notation (cont.)



weights between two layers \rightarrow a matrix

DNNs – Notation (cont.)

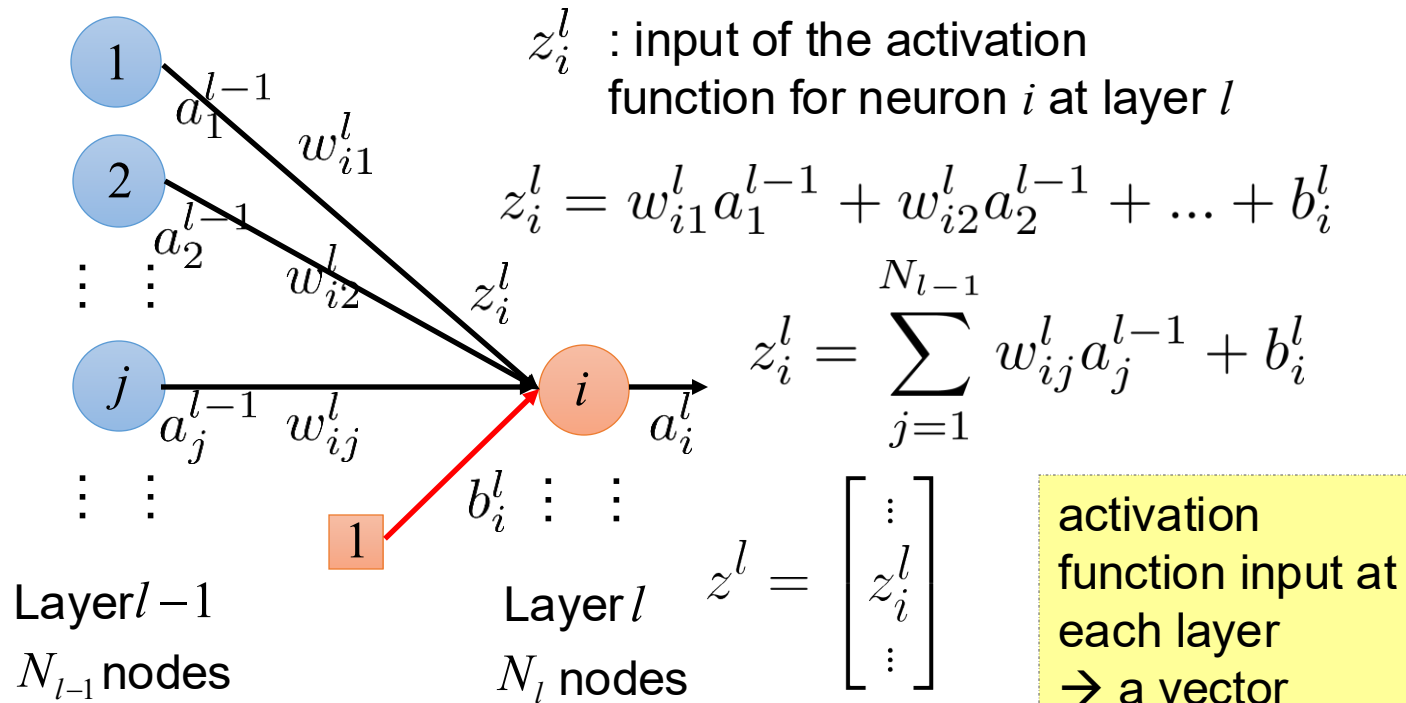


b_i^l : bias for neuron i
at layer l

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

bias of all neurons at each layer \rightarrow a vector

DNNs – Notation (cont.)



DNNs – Notation Summary



a_i^l : output of a neuron

a^l : output vector of a layer

z_i^l : input of activation
function

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

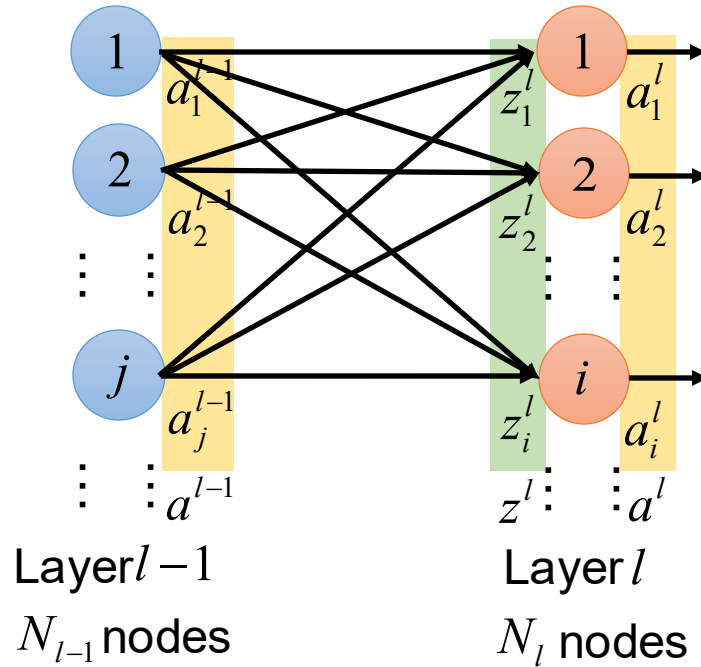
w_{ij}^l : a weight

W^l : a weight matrix

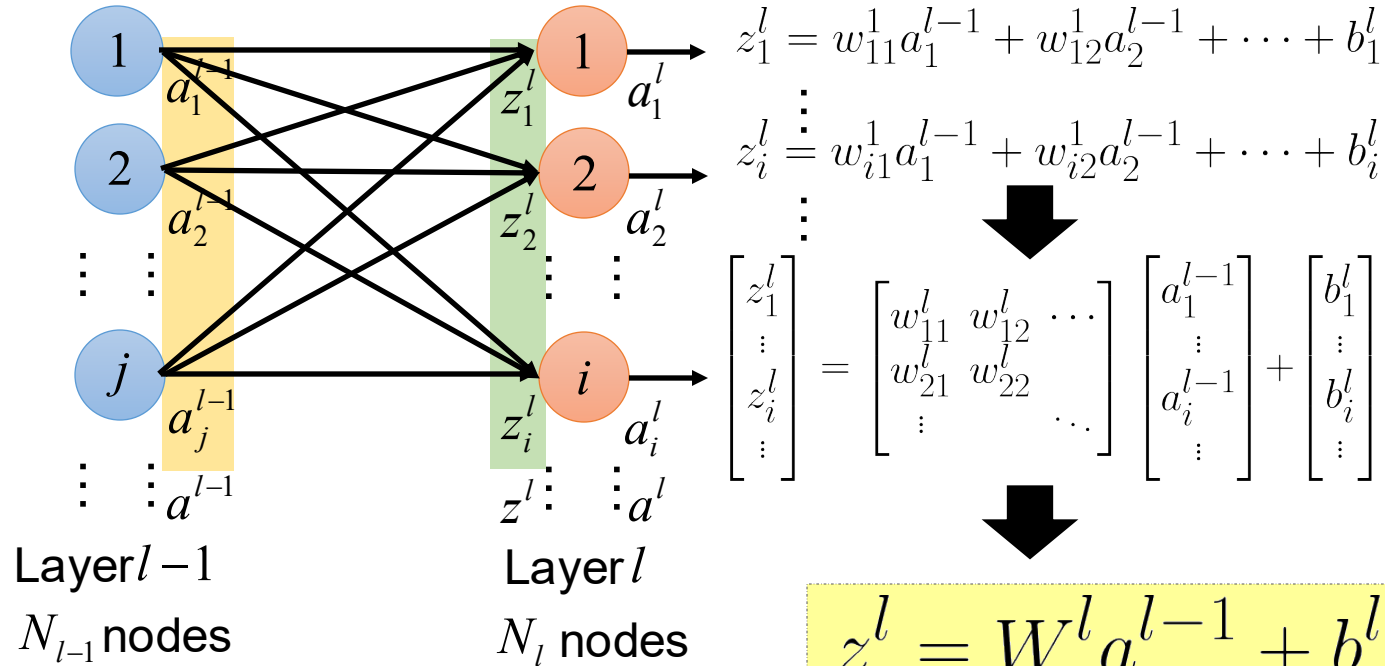
b_i^l : a bias

b^l : a bias vector

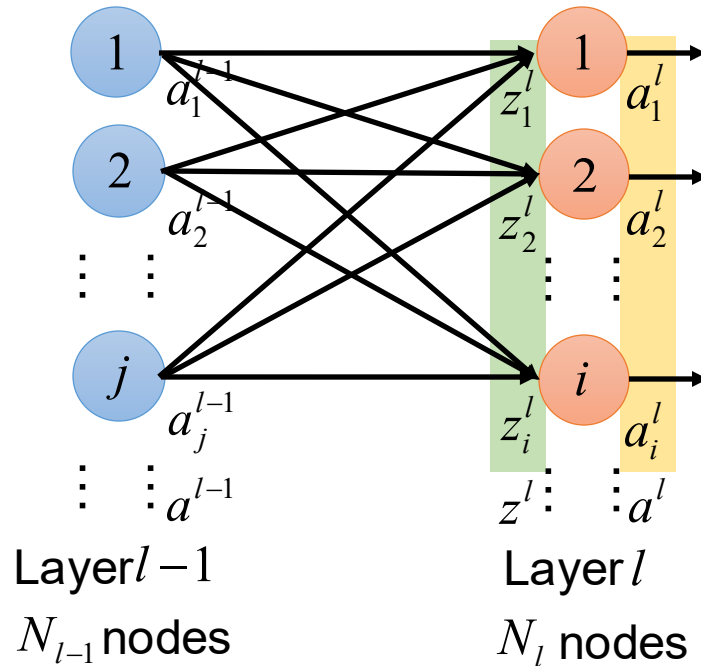
Output of a Layer



Output of a Layer – From a to z



Output of a Layer – From z to a

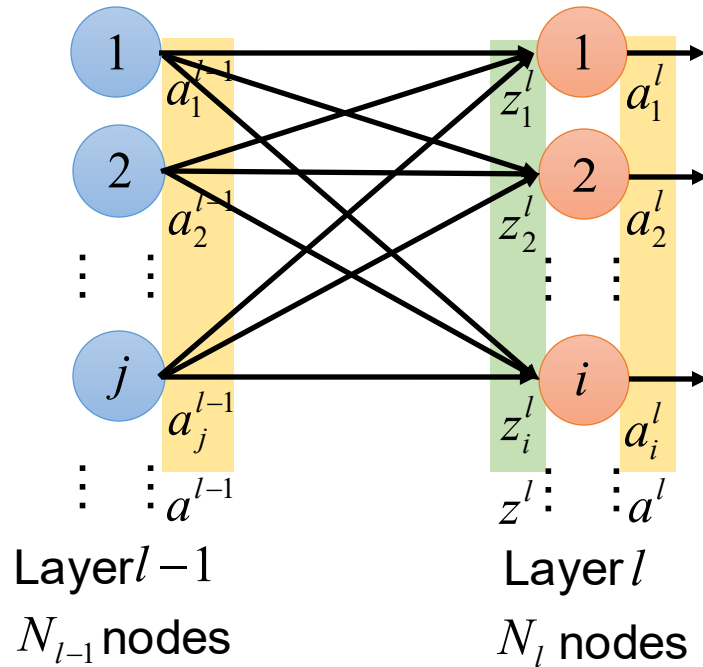


$$a_i^l = \sigma(z_i^l)$$

$$\begin{bmatrix} a_1^l \\ a_2^l \\ \vdots \\ a_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} \sigma(z_1^l) \\ \sigma(z_2^l) \\ \vdots \\ \sigma(z_i^l) \\ \vdots \end{bmatrix}$$

$$a^l = \sigma(z^l)$$

Output of a Layer



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

$$a^l = \sigma(z^l)$$



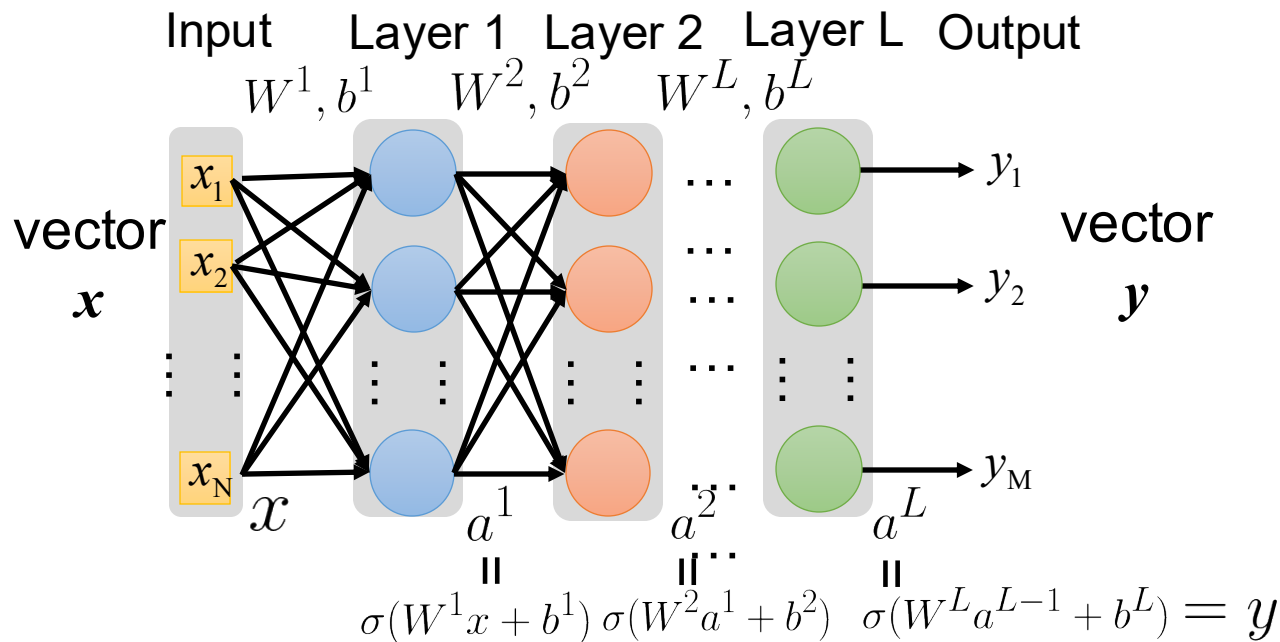
$$a^l = \sigma(W^l a^{l-1} + b^l)$$

DNN Formulation



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

- Fully connected feedforward network

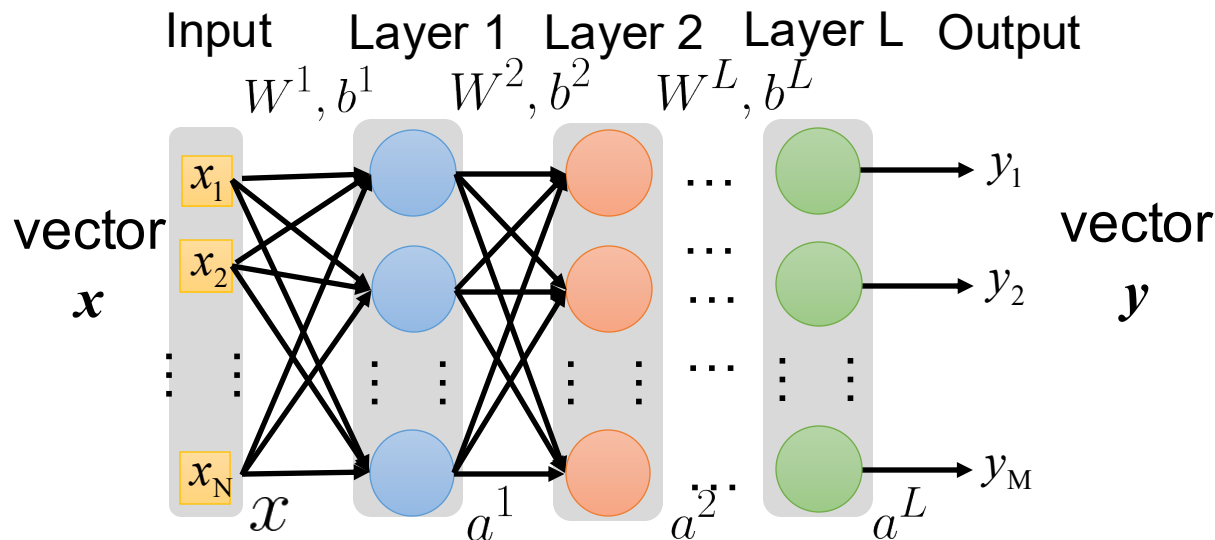


DNN Formulation



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

- Fully connected feedforward network

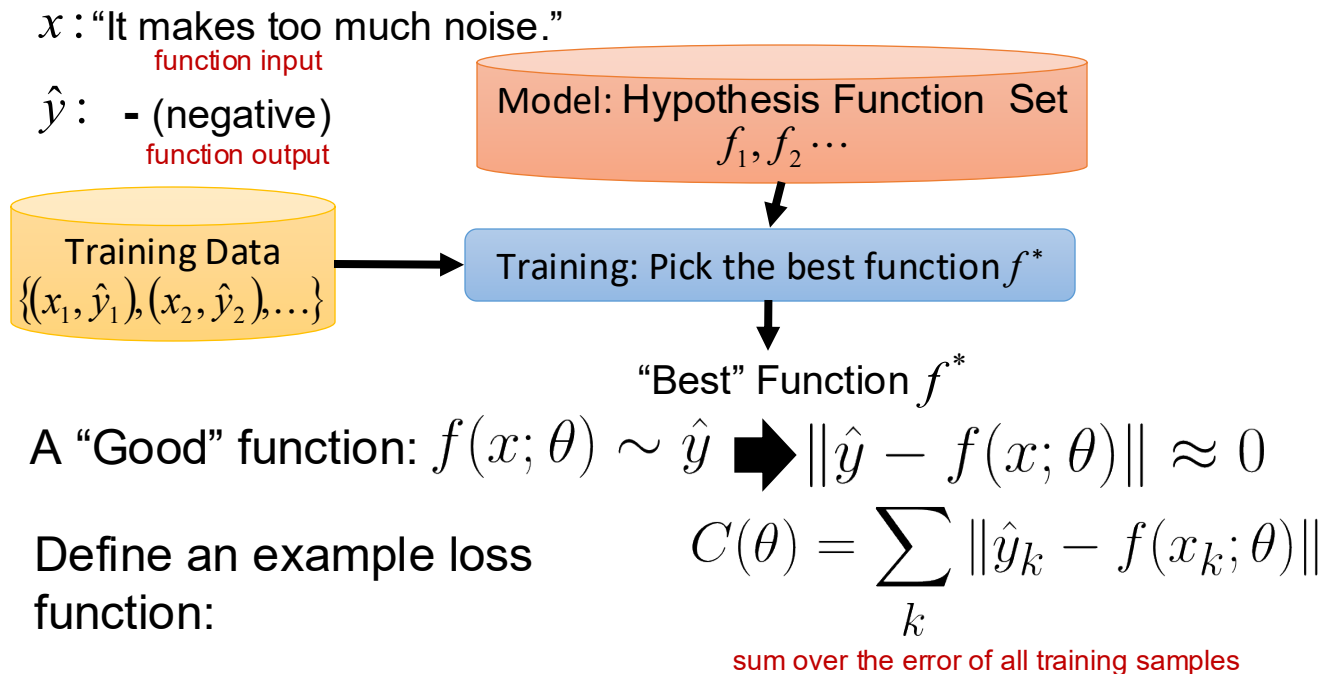


$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

Loss Function for Training



[Back to an earlier slide!](#)



Gradient Descent for DNNs



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

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

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

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

Algorithm

Initialization: start at θ^0

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

{

 compute gradient at θ^i

 update parameters

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

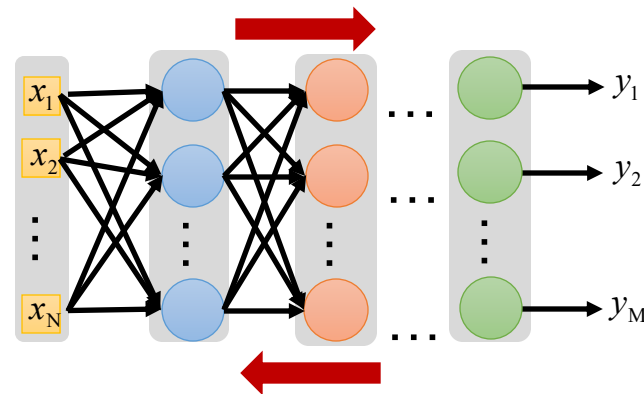
}

To update weights efficiently, we use **backpropagation**.

Forward vs. Backward 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 loss $C(\theta)$
 - 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



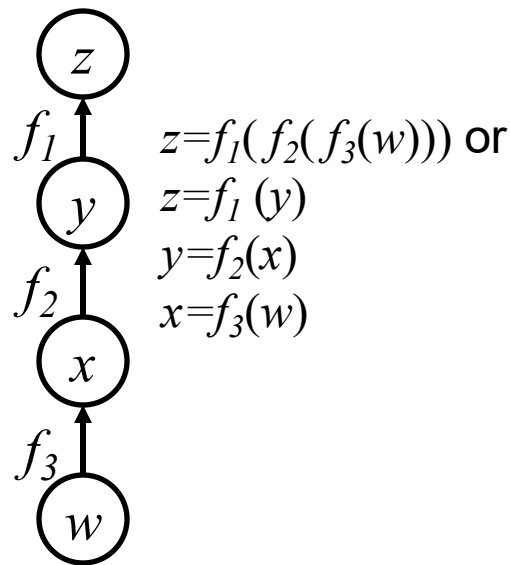
$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w}$$

$$= f'_1(y) f'_2(x) f'_3(w)$$

forward propagation for computing loss

$$= f'_1(f_2(f_3(w))) f'_2(f_3(w)) f'_3(w)$$

back-propagation of gradients for updating weights



Gradient Descent for DNNs



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

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

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

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

Algorithm

Initialization: start at θ^0

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

{

 compute gradient at θ^i

 update parameters

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

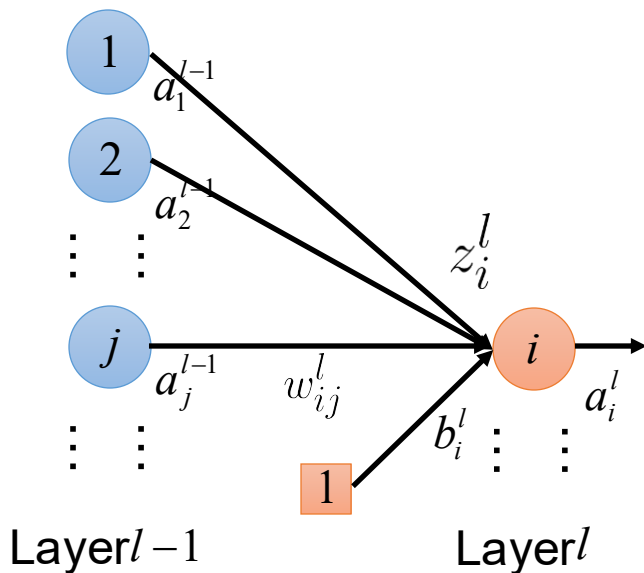
}

To update weights efficiently, we use **backpropagation**.

Gradient Descent for DNNs (cont.)



$$\partial C(\theta) / \partial w_{ij}^l$$

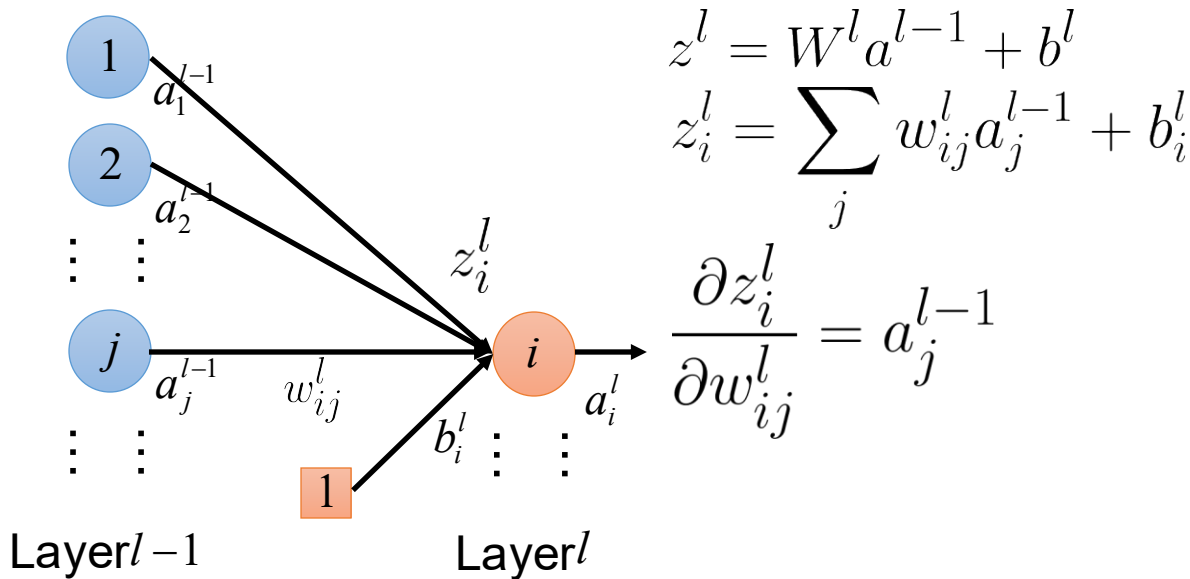


$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

Gradient Descent for DNNs (cont.)



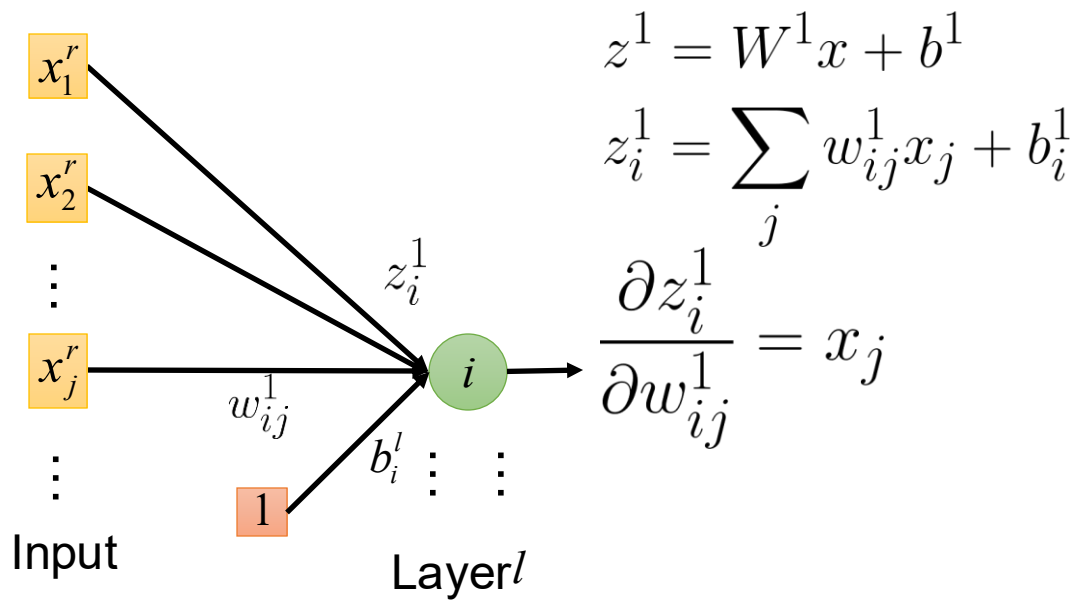
$$\partial z_i^l / \partial w_{ij}^l \quad (l > 1)$$



Gradient Descent for DNNs (cont.)



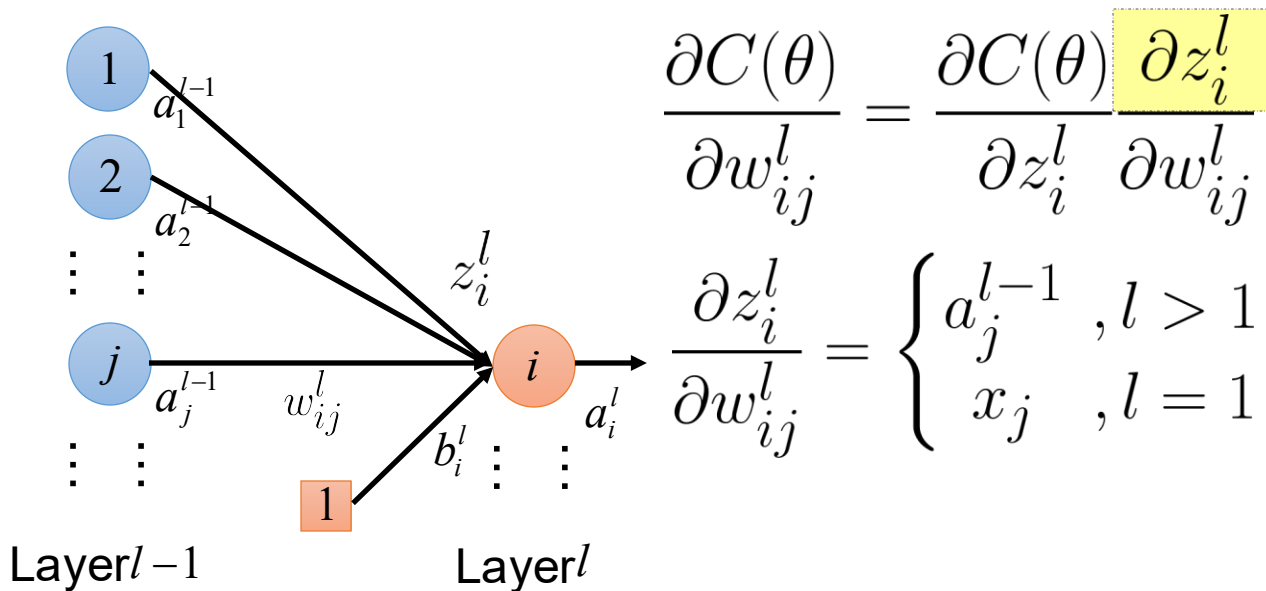
$$\partial z_i^l / \partial w_{ij}^l \quad (l = 1)$$



Gradient Descent for DNNs (cont.)



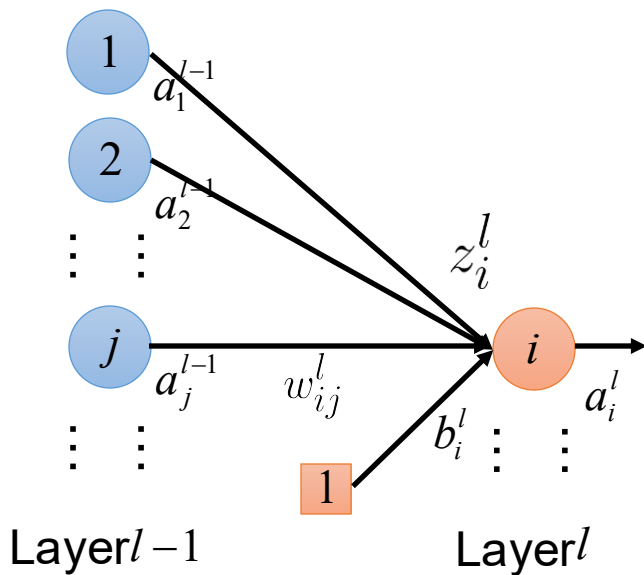
$$\partial C(\theta) / \partial w_{ij}^l$$



Gradient Descent for DNNs (cont.)

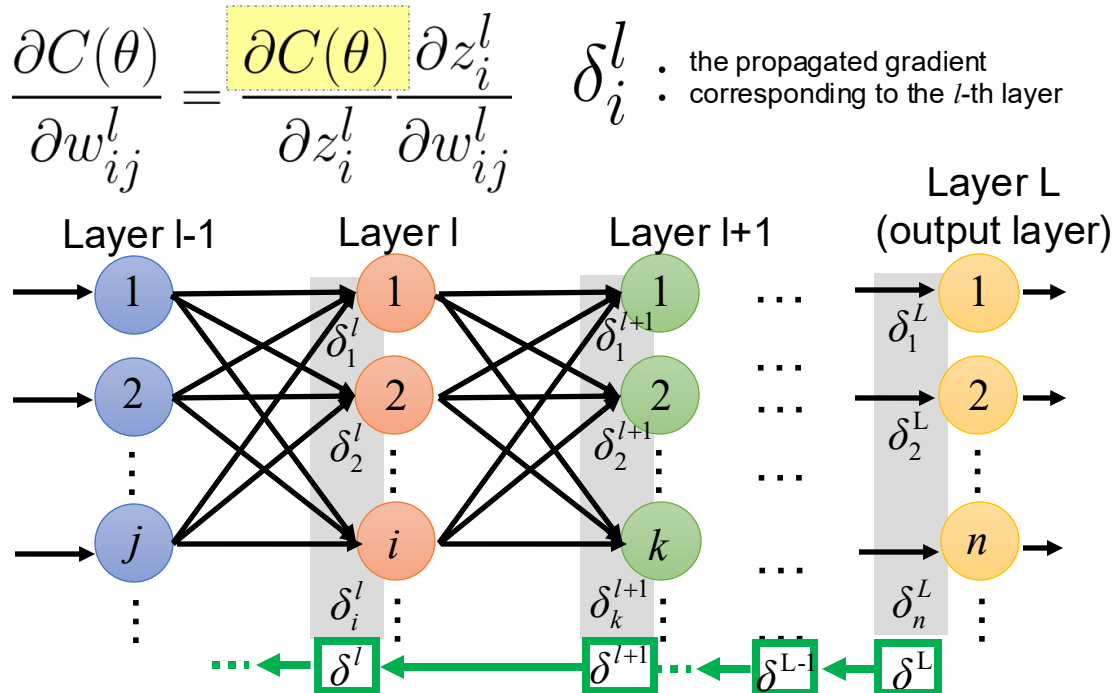


$$\partial C(\theta) / \partial w_{ij}^l$$



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

Gradient Descent for DNNs (cont.)



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

Gradient Descent for DNNs (cont.)



$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

- Idea: from L to 1
- ① Initialization: compute δ^L
- ② Compute δ^l based on δ^{l+1}

Gradient Descent for DNNs (cont.)



$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

- Idea: from L to 1

① **Initialization: compute δ^L**

② Compute δ^l based on δ^{l+1}

$$\begin{aligned} \delta_i^L &= \frac{\partial C}{\partial z_i^L} & \Delta z_i^L &\rightarrow \Delta a_i^L = \Delta y_i \rightarrow \Delta C \\ &= \frac{\partial C}{\partial y_i} \frac{\partial y_i}{\partial z_i^L} \end{aligned}$$

$\partial C / \partial y_i$ depends on the loss function

Gradient Descent for DNNs (cont.)

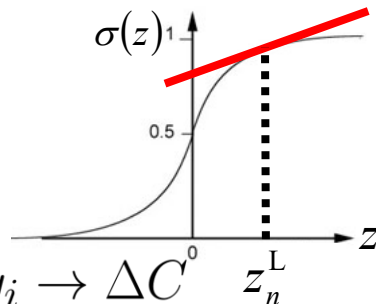


$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

- Idea: from L to 1

① **Initialization: compute δ^L**

② Compute δ^l based on δ^{l+1}



$$\begin{aligned} \delta_i^L &= \frac{\partial C}{\partial z_i^L} & \Delta z_i^L \rightarrow \Delta a_i^L = \Delta y_i \rightarrow \Delta C^0 & z_n^L \\ &= \frac{\partial C}{\partial y_i} \frac{\partial y_i}{\partial z_i^L} = a_i^L = \sigma(z_i^L) & \sigma'(z^L) = \begin{bmatrix} \sigma'(z_1^L) \\ \sigma'(z_2^L) \\ \vdots \\ \sigma'(z_i^L) \\ \vdots \end{bmatrix} & \nabla C(y) = \begin{bmatrix} \frac{\partial C}{\partial y_1} \\ \frac{\partial C}{\partial y_2} \\ \vdots \\ \frac{\partial C}{\partial y_i} \\ \vdots \end{bmatrix} \\ &= \frac{\partial C}{\partial y_i} \sigma'(z_i^L) & \delta^L = \sigma'(z^L) \odot \nabla C(y) \end{aligned}$$

Gradient Descent for DNNs (cont.)

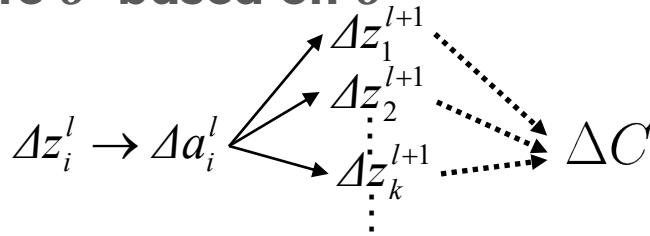


$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

- Idea: from L to 1

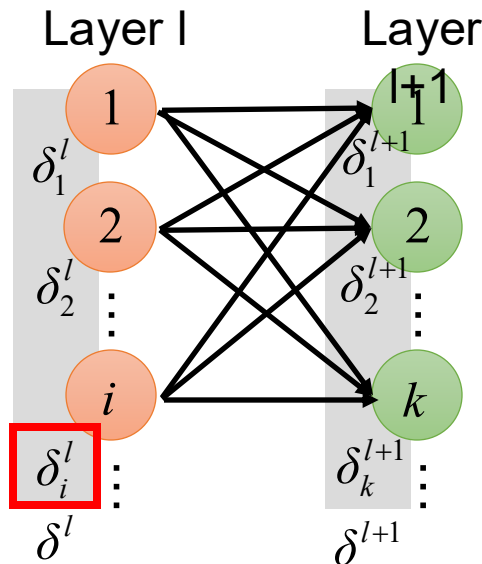
① Initialization: compute δ^L

② Compute δ^l based on δ^{l+1}



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

δ_i^{l+1}



Gradient Descent for DNNs (cont.)



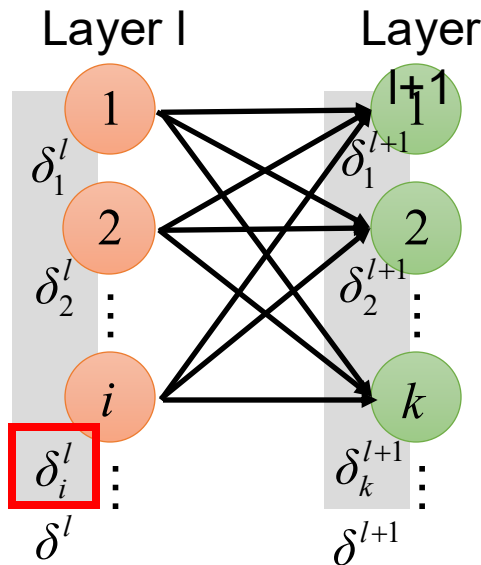
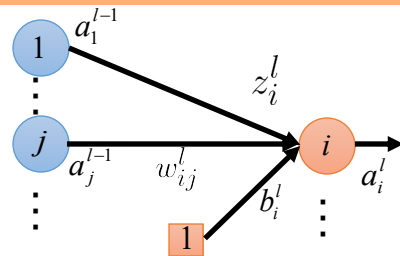
$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

- Idea: from L to 1

① Initialization: compute δ^L

② Compute δ^l based on δ^{l+1}

$$\begin{aligned} \delta_i^l &= \frac{\partial a_i^l}{\partial z_i^l} \sum_k \frac{\partial z_k^{l+1}}{\partial a_i^l} \delta_k^{l+1} \\ &= \sum_k w_{ki}^{l+1} a_i^l + b_k^{l+1} \\ &= \sigma'(z_i) \sum_k \frac{\partial z_k^{l+1}}{\partial a_i^l} \delta_k^{l+1} \\ &= \sigma'(z_i) \sum_k w_{ki}^{l+1} \delta_k^{l+1} \end{aligned}$$



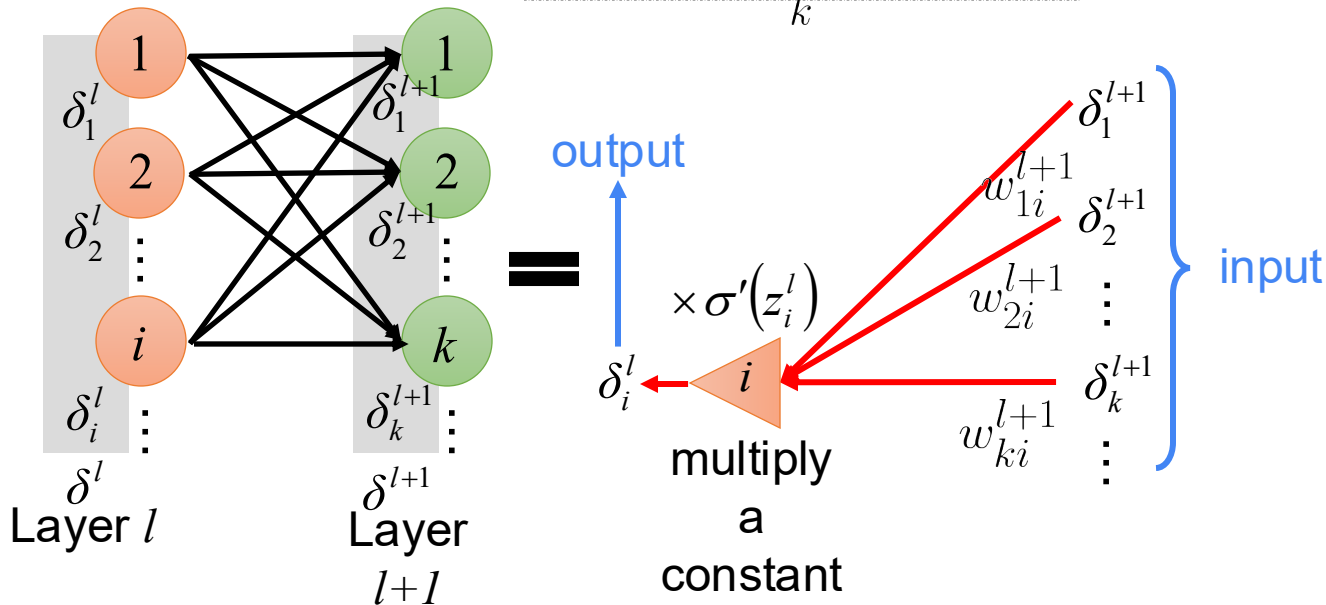
Gradient Descent for DNNs (cont.)



$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

- Rethink the propagation

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



Gradient Descent for DNNs (cont.)

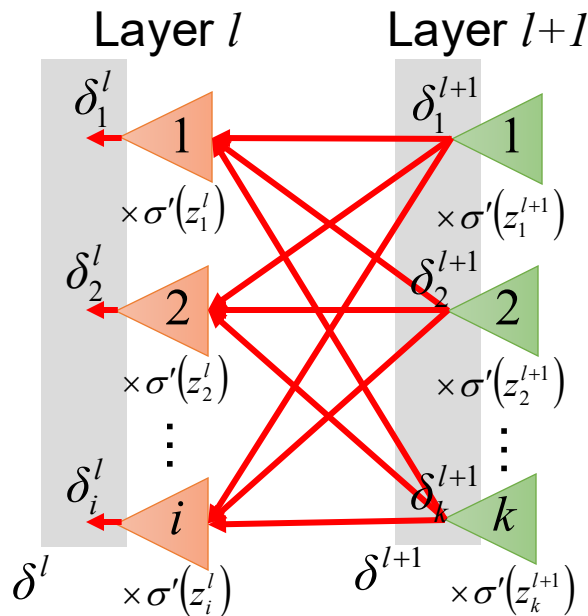


$$\partial C(\theta) / \partial z_i^l = \delta_i^l$$

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

$$\sigma'(z^l) = \begin{bmatrix} \sigma'(z_1^l) \\ \sigma'(z_2^l) \\ \vdots \\ \sigma'(z_i^l) \\ \vdots \end{bmatrix}$$

$$\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$$



Gradient Descent for DNNs (cont.)



$$\frac{\partial C(\theta)}{\partial z_i^l} = \delta_i^l \quad \frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

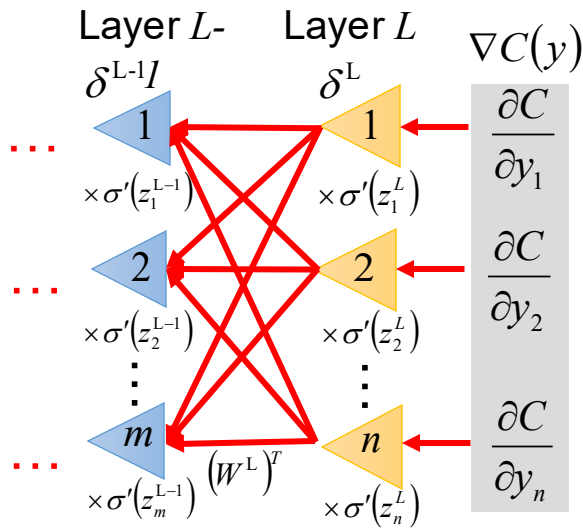
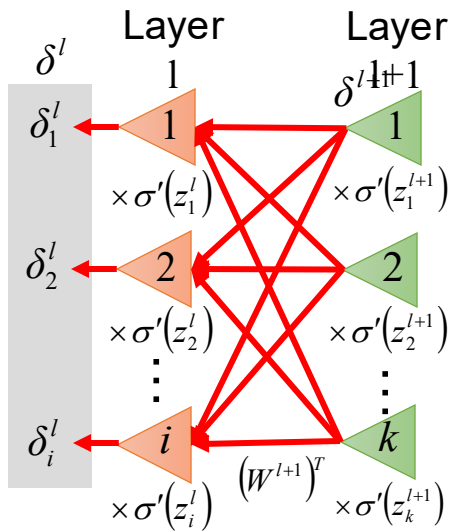
- Idea: from L to 1

① Initialization: compute δ^L

② Compute δ^{l-1} based on δ^l

$$\delta^L = \sigma'(z^L) \odot \nabla C(y)$$

$$\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$$



Backpropagation

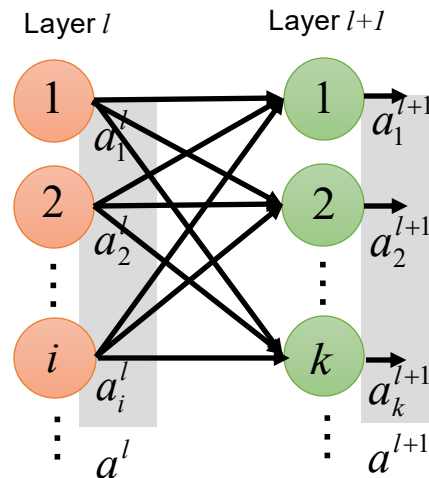


$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

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

Forward Pass

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



Backpropagation

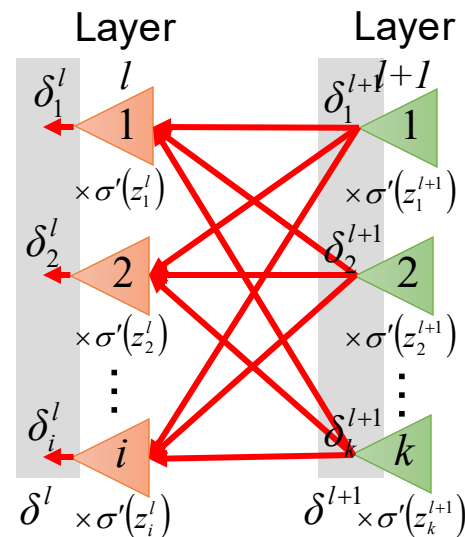


$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

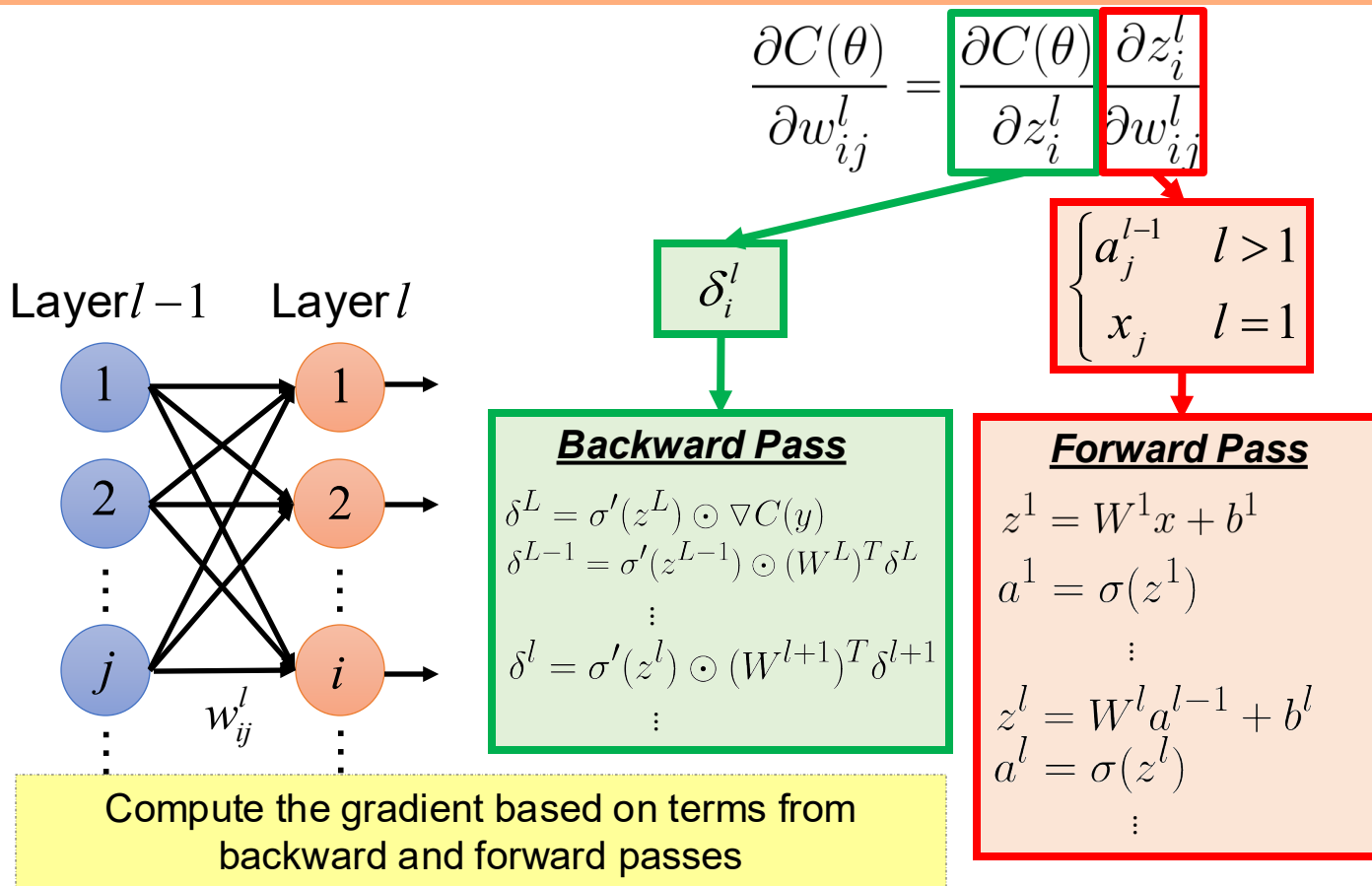
$$\frac{\partial C(\theta)}{\partial z_i^l} = \delta_i^l$$

Backward Pass

$$\begin{aligned}\delta^L &= \sigma'(z^L) \odot \nabla C(y) \\ \delta^{L-1} &= \sigma'(z^{L-1}) \odot (W^L)^T \delta^L \\ &\vdots \\ \delta^l &= \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1} \\ &\vdots\end{aligned}$$



Gradient Descent - Summary



Topics for Thursday



- Distributional Similarity
- Sparse Word Representations
- Word Embeddings and Word2Vec
- Language Modeling
- Unexpected things we learn with word embeddings