



Introduction to Machine Learning (CS419M)

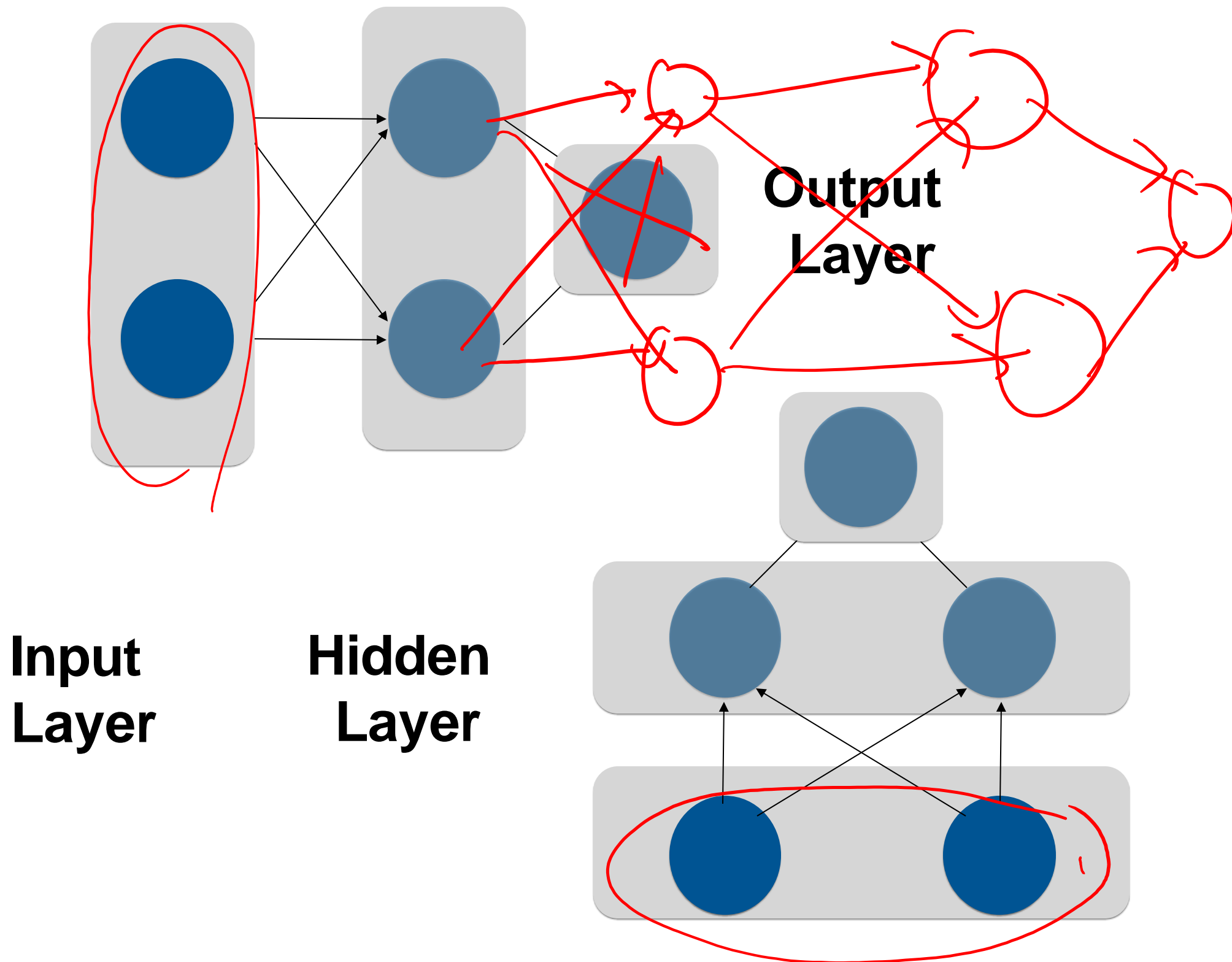
Lecture 17: Introduction to Neural Networks

Original author: Preethi Jyothi

Modified by: Sunita Sarawagi

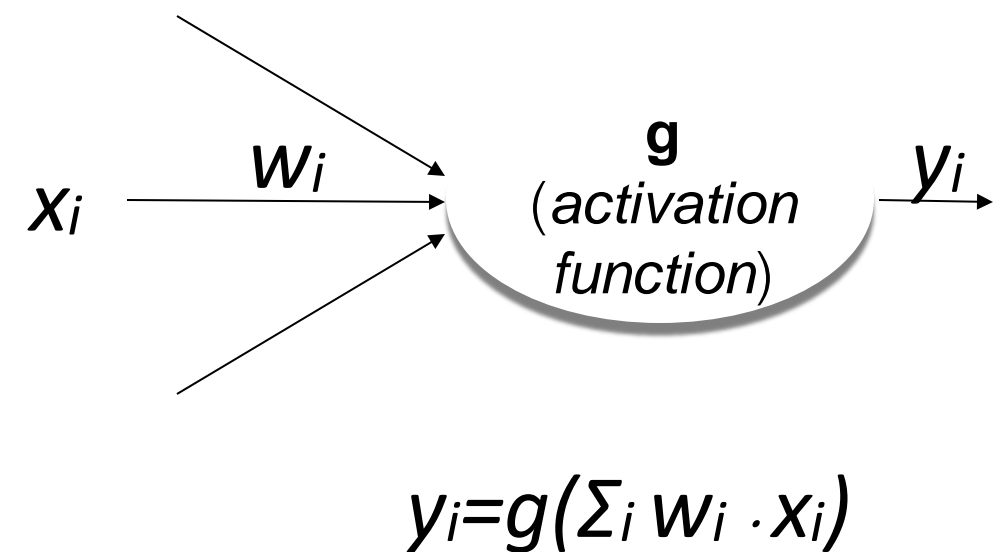
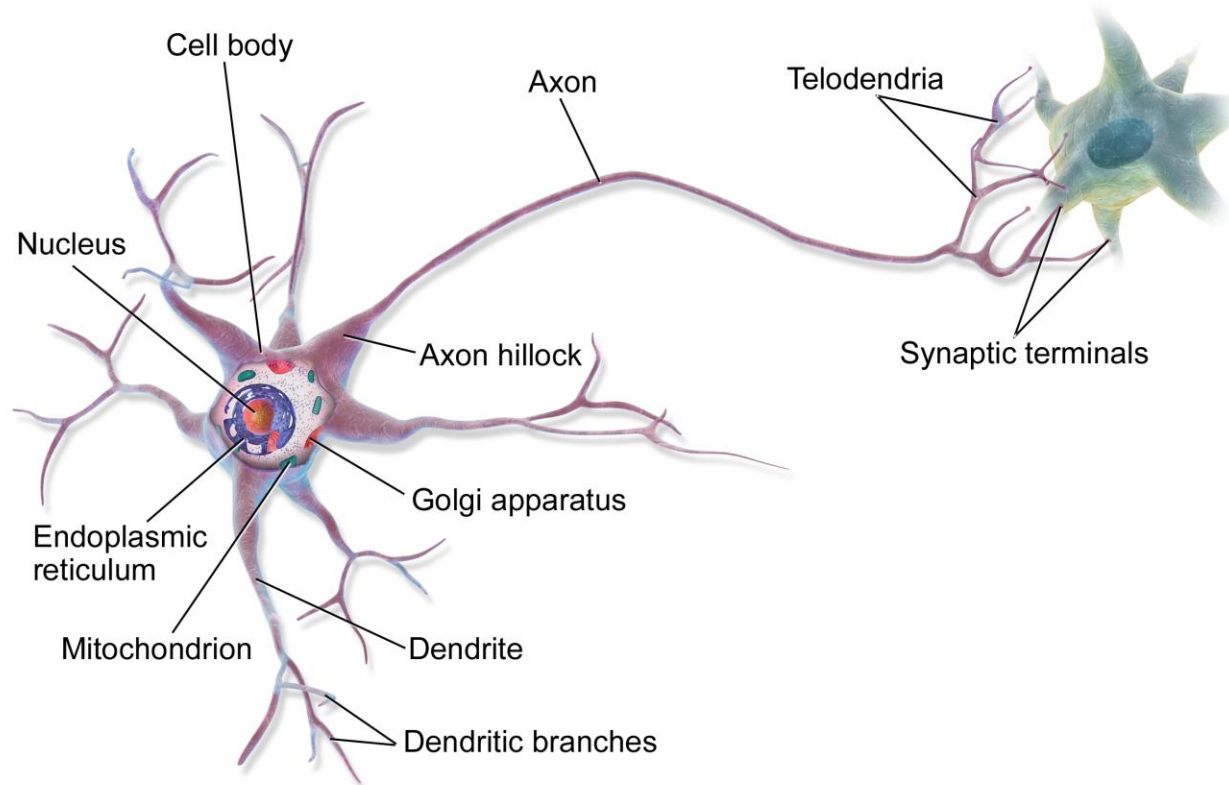
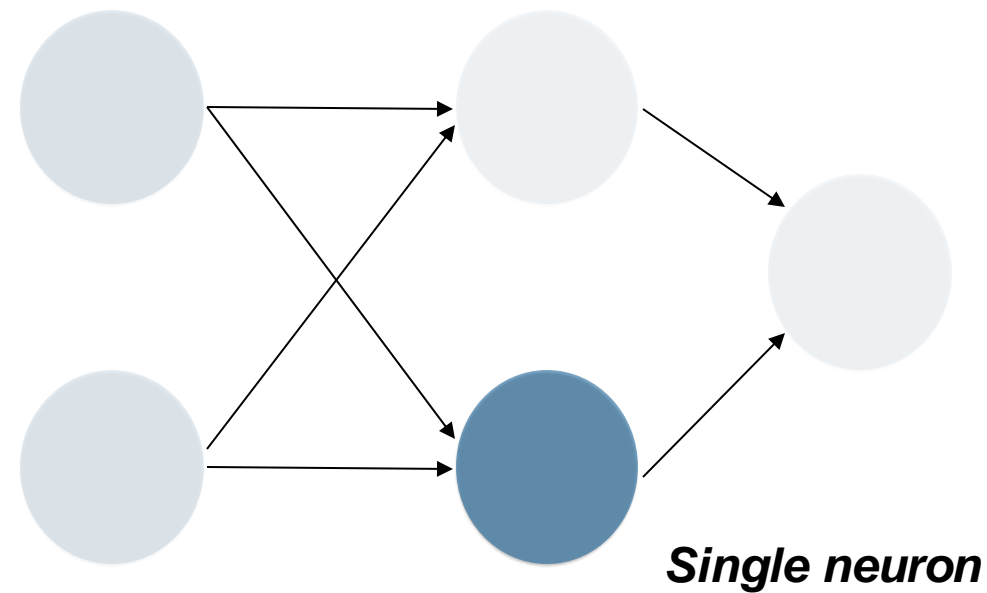
Mar 16, 2018

Feed-forward Neural Network



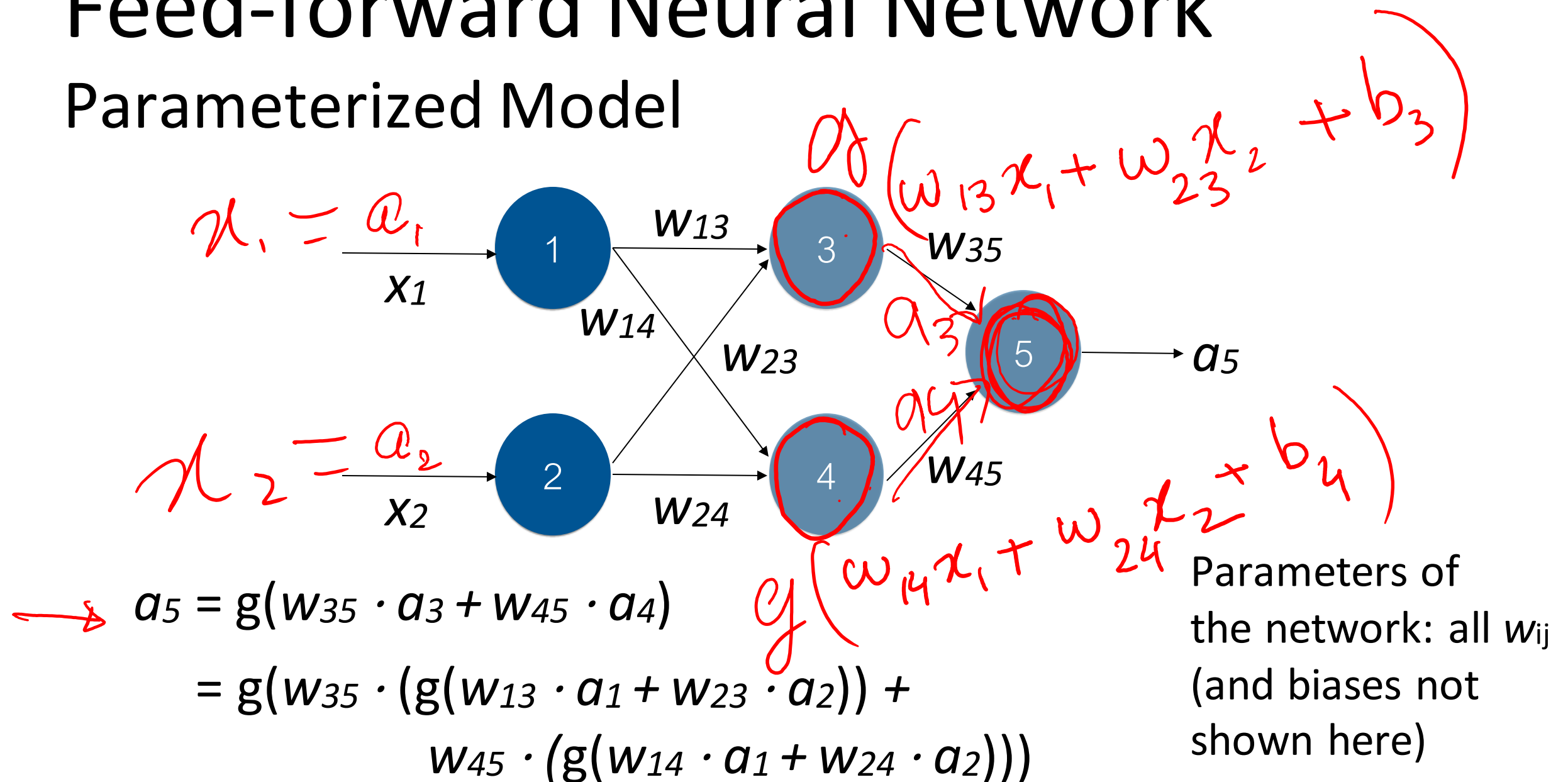
Feed-forward Neural Network

Brain Metaphor



Feed-forward Neural Network

Parameterized Model



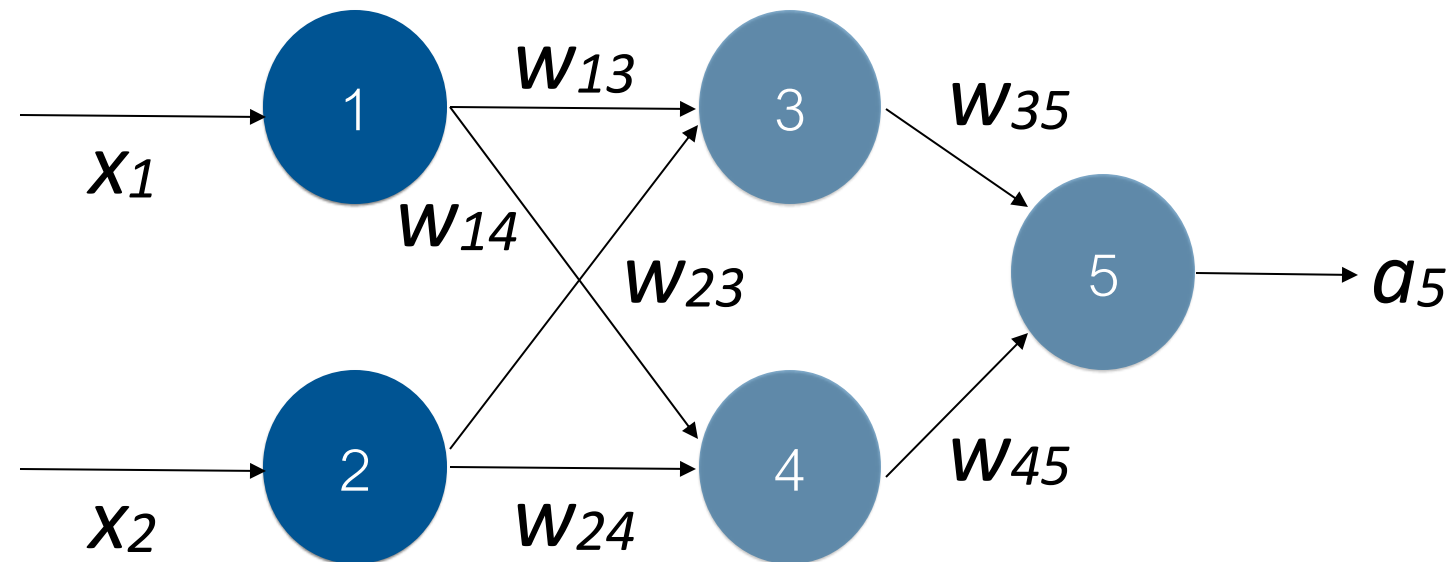
If \mathbf{x} is a 2-dimensional vector and the layer above it is a 2-dimensional vector \mathbf{h} , a fully-connected layer is associated with:

$$\mathbf{h} = \mathbf{xW} + \mathbf{b}$$

where w_{ij} in \mathbf{W} is the weight of the connection between i^{th} neuron in the input row and j^{th} neuron in the first hidden layer and \mathbf{b} is the bias vector

Feed-forward Neural Network

Parameterized Model



$$\begin{aligned}
 a_5 &= g(w_{35} \cdot a_3 + w_{45} \cdot a_4) \\
 &= g(w_{35} \cdot (g(w_{13} \cdot a_1 + w_{23} \cdot a_2)) + \\
 &\quad w_{45} \cdot (g(w_{14} \cdot a_1 + w_{24} \cdot a_2)))
 \end{aligned}$$

The simplest neural network is the perceptron:

$$\text{Perceptron}(x) = \mathbf{xW} + \mathbf{b}$$

A 1-layer feedforward neural network has the form:

$$\text{MLP}(x) = g(\mathbf{xW}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2$$

Activation Functions (g)

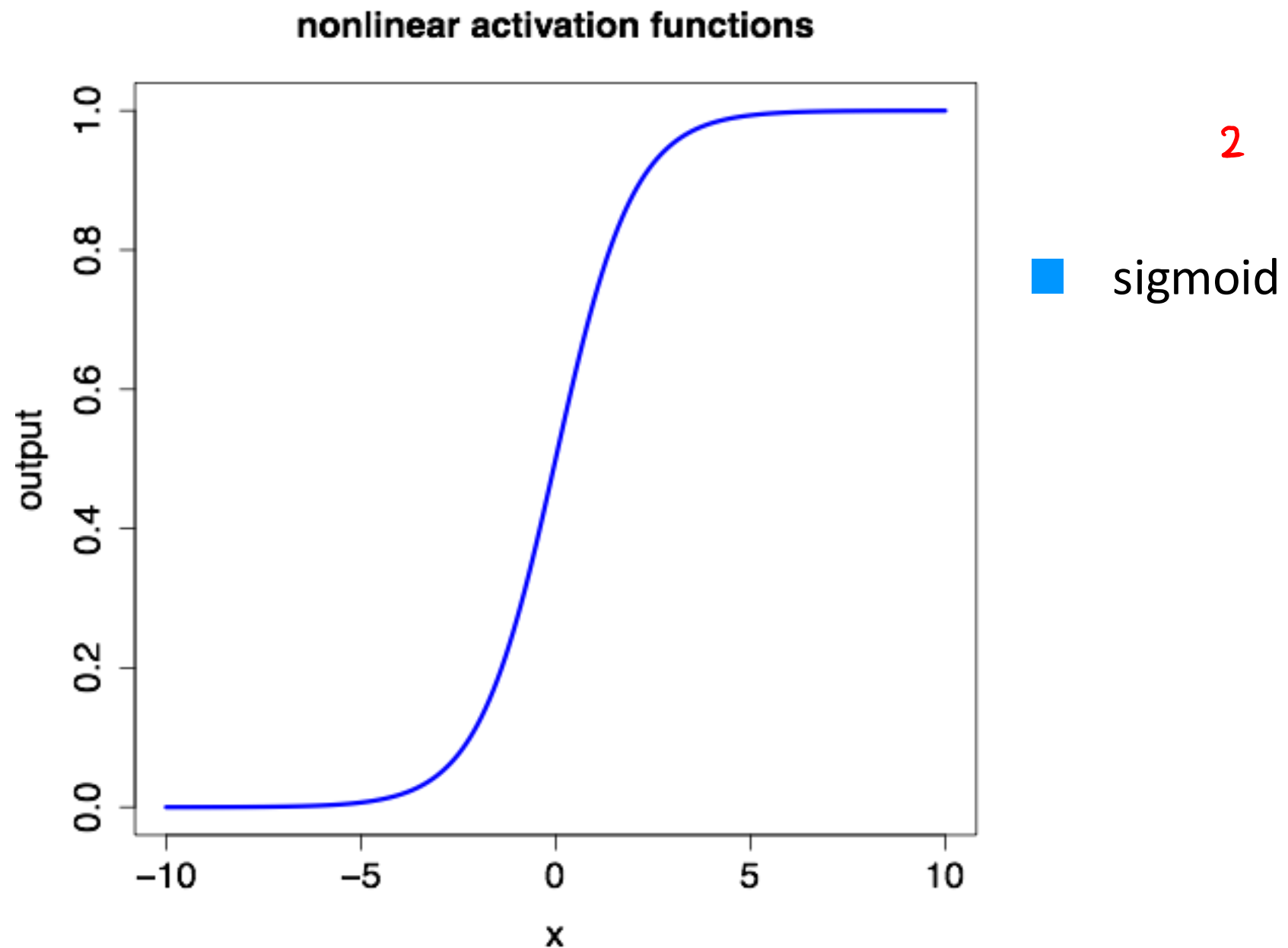
- Cannot be linear: Will get back a linear classifier otherwise: Show.

$$\begin{aligned} &\rightarrow g(xW_1 + b_1)W_2 + b_2 \\ &(xW_1 + b_1)W_2 + b_2 \\ &\begin{bmatrix} w_{11}x_1 + w_{13}x_2 \\ w_{12}x_1 + w_{14}x_2 \end{bmatrix}^T \begin{bmatrix} w_{21} \\ w_{22} \end{bmatrix} \\ &= \underbrace{w_{11}w_{21}}_{\underline{\quad}} x_1 + \underbrace{w_{13}w_{21}}_{\underline{\quad}} x_2 + \underline{\quad} \end{aligned}$$
$$\begin{aligned} g(z) &= z \\ x &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad b_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ w_2 &= \begin{bmatrix} w_{21} \\ w_{22} \end{bmatrix} \quad b_2 = 0 \\ w_1 &= \begin{bmatrix} w_{11} & w_{12} \\ w_{13} & w_{14} \end{bmatrix} \end{aligned}$$

- Want a function that is efficient to compute, easy to optimizer (informative gradient), almost linear

Common Activation Functions (g)

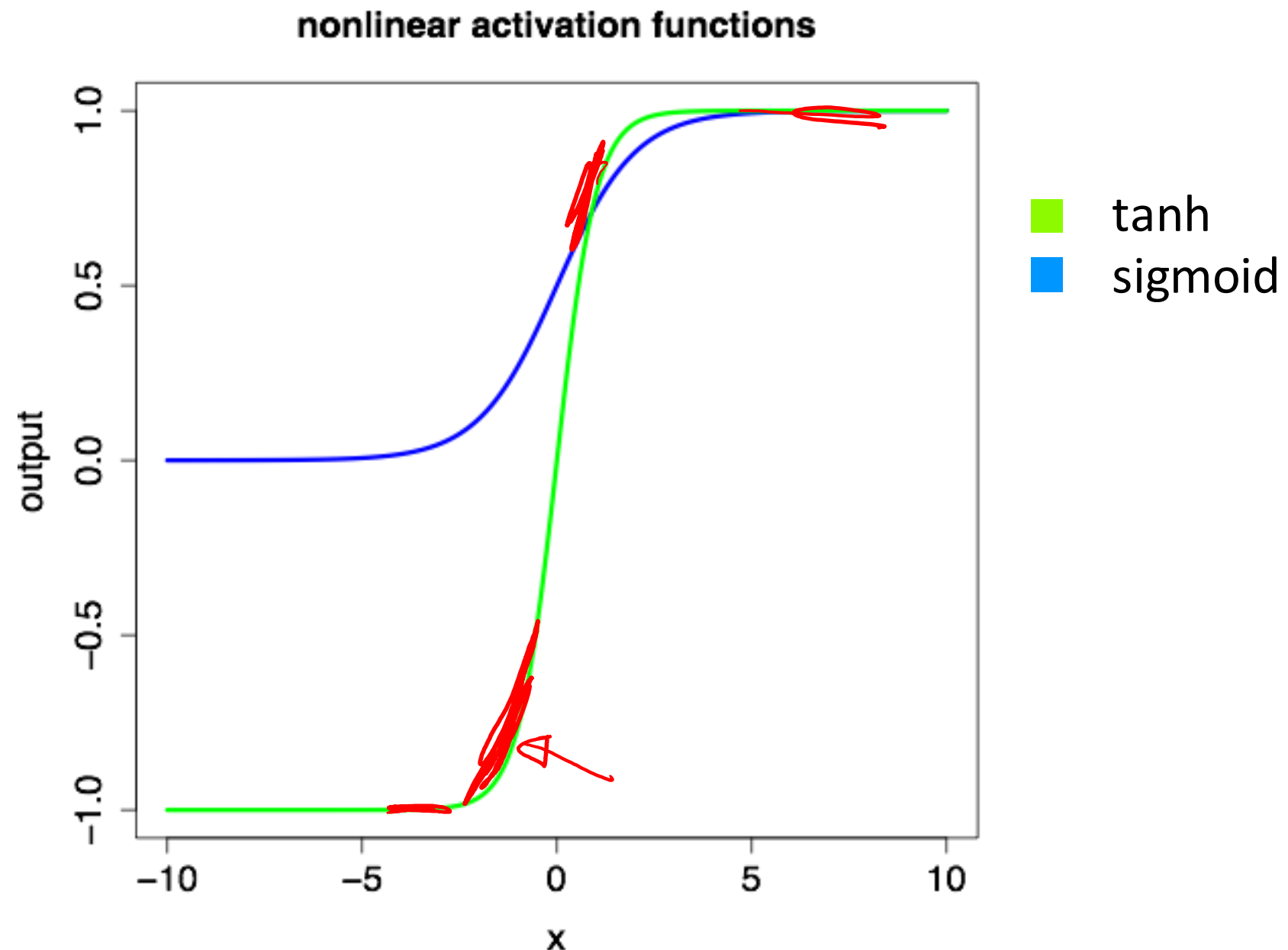
Sigmoid: $\sigma(x) = 1/(1 + e^{-x})$



Common Activation Functions (g)

Sigmoid: $\sigma(x) = 1/(1 + e^{-x})$

Hyperbolic tangent (tanh): $\tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$

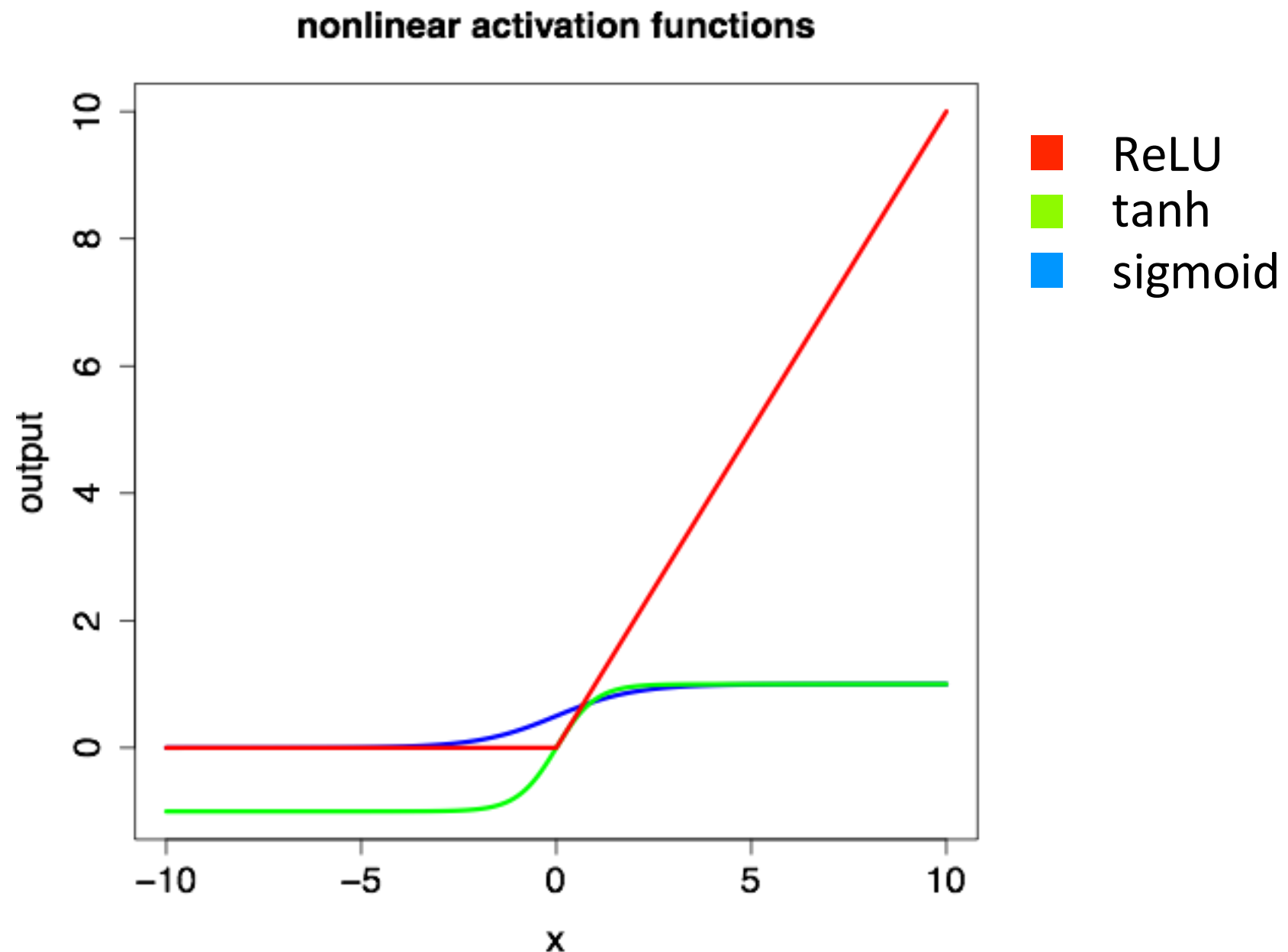


Common Activation Functions (g)

Sigmoid: $\sigma(x) = 1/(1 + e^{-x})$

Hyperbolic tangent (tanh): $\tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$

Rectified Linear Unit (ReLU): $\text{ReLU}(x) = \max(0, x)$



Choosing $g()$

Considerations: want some non-linearity, informative gradient (e.g. when convex), fast computation, close to linear

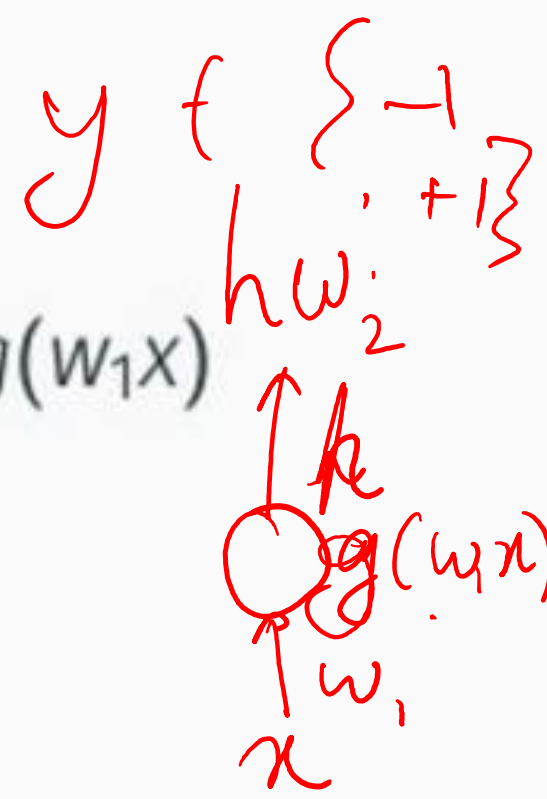
Role of the gradient of g during training

Training objective of DNN with one hidden unit $h = g(w_1x)$

$$J(w_1, w_2, x, y) = L(hw_2y) = L(g(w_1x)w_2y)$$

Gradient of above w.r.t w_1 is $L'w_2yg'x$

If $g' = 0$, the gradient becomes zero and we do not know in what direction to move w_1 .



Choosing $g()$

- RELU: not differential but okay since gradient is informative. second-derivative zero in most places (useful for optimization)
 - Caution: watch out for inactive RelU: initialize affine input bias parameter to small positives. Gradient zero \implies information flow to lower layers is blocked.
- Sigmoid/Tanh: $\tanh(z) = 2 \text{ sigmoid}(2z)$. Non-convex. Well-behaved (linear) only for small values of z , gradients very small for small or large z , problem for multi-layer network.

Example XOR

Neural networks can model decisions that conventional linear classifiers cannot.

$$y = f^*(x) = x_1 \oplus x_2$$

Training data = all four combinations.

Linear classifier $\hat{y} = w_1x_1 + w_2x_2 + b$ trained with least square loss yields $w_1 = w_2 = 0, b = 1/2$

Cannot discriminate

Non-linear classifier such as one with x_1x_2 as feature

$(\hat{y} = w_1x_1 + w_2x_2 + w_3x_1x_2 + b)$ can discriminate but the burden is on us to create the useful non-linear features.

Example: XOR

0000

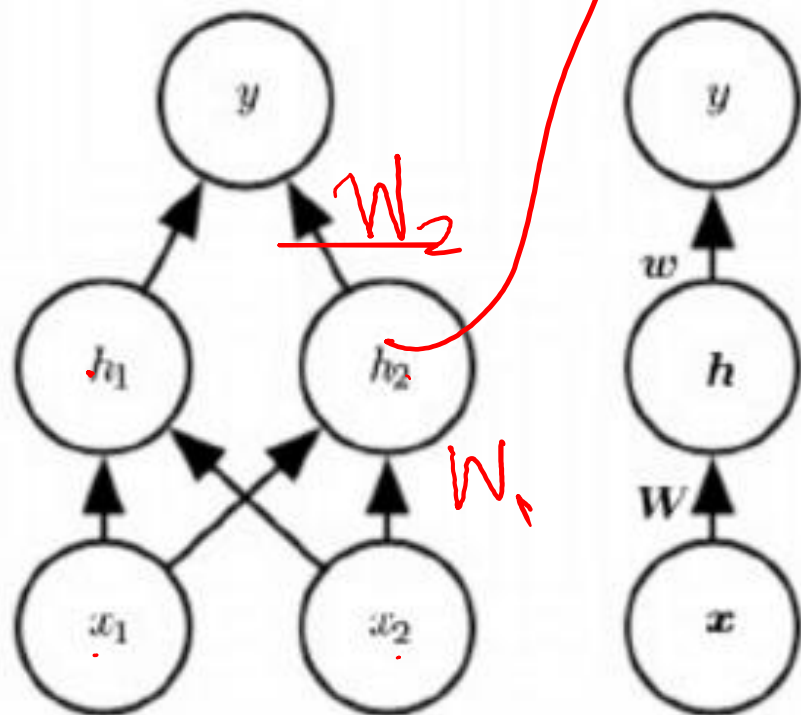
A generic two layer neural network with ReLU:

$$y = f(x) = W^2 \max(0, W^1 x + b^1) + b^2$$

Role of non-linear transform.

$$W^1 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$b^1 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, W^2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b^2 = 0$$



$g(W^1 x + b^1)$

$[h_1, h_2]$

$$= g\left(\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix}\right)$$

$$\max\left(\begin{bmatrix} x_1 + x_2 & x_1 + x_2 - 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix}\right)$$

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$W^2 h + b^2 = 0$$

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$2 - 2 = 0$$

Output layers

- Depends on the output type

- Binary class labels: sigmoid function transforms arbitrary reals to probability of Bernoulli

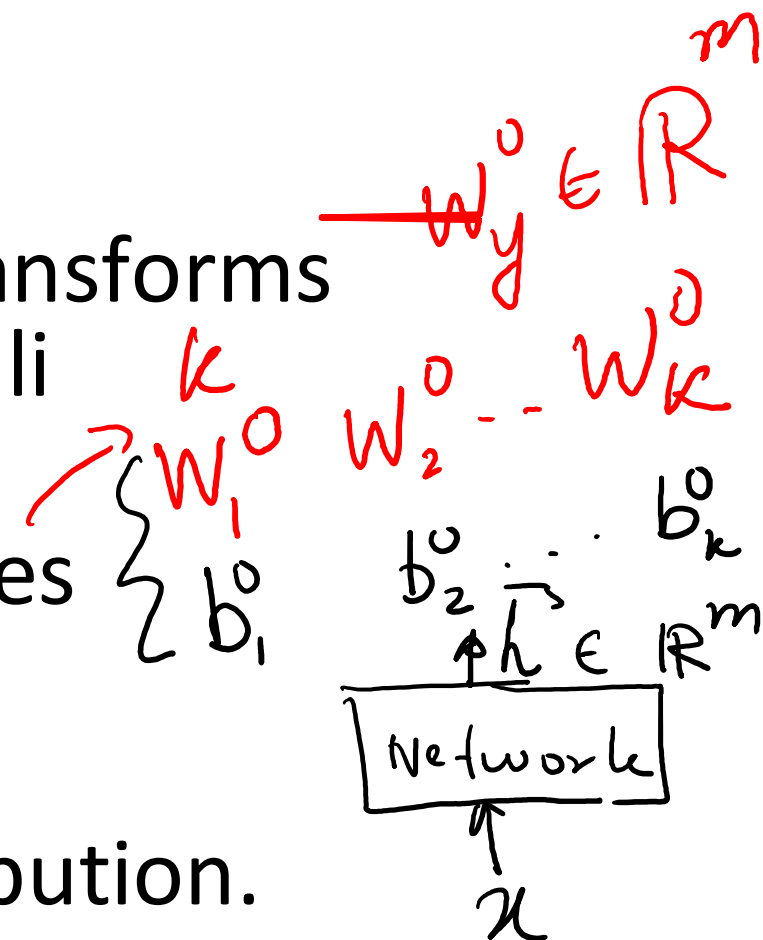
- Multi-class class labels: softmax provides multinomial probabilities

$$P(y|x) = \frac{e^{w_{yh} + b_y}}{\sum_{y'} e^{w_{y'h} + b_{y'}}$$

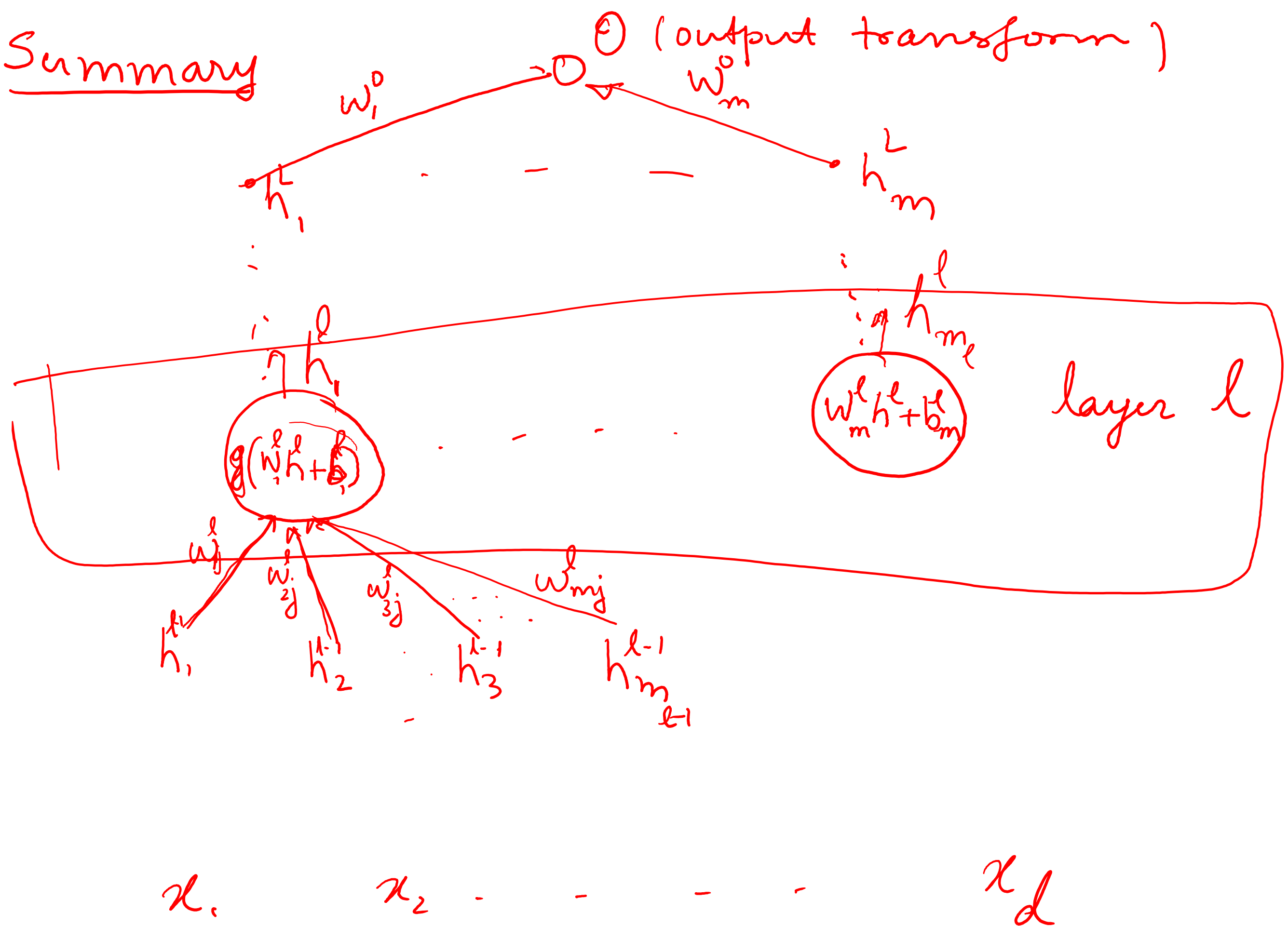
- Real: output is mean of Gaussian distribution.

- Advantage of all of above: Probability distribution over output. Maximum likelihood training loss is convex in parameters of outer-most layer.

- Similar to conventional training.



Summary



Training a Feed-forward network

- To train a neural network, define a loss function $L(y, \tilde{y})$:
a function of the true output y and the predicted output \tilde{y}

- $L(y, \tilde{y})$ assigns a non-negative numerical score to the neural network's output, \tilde{y}

- The parameters of the network are set to minimise L over the training examples (i.e. a sum of losses over different training samples)

- L is typically minimised using a *gradient-based method*

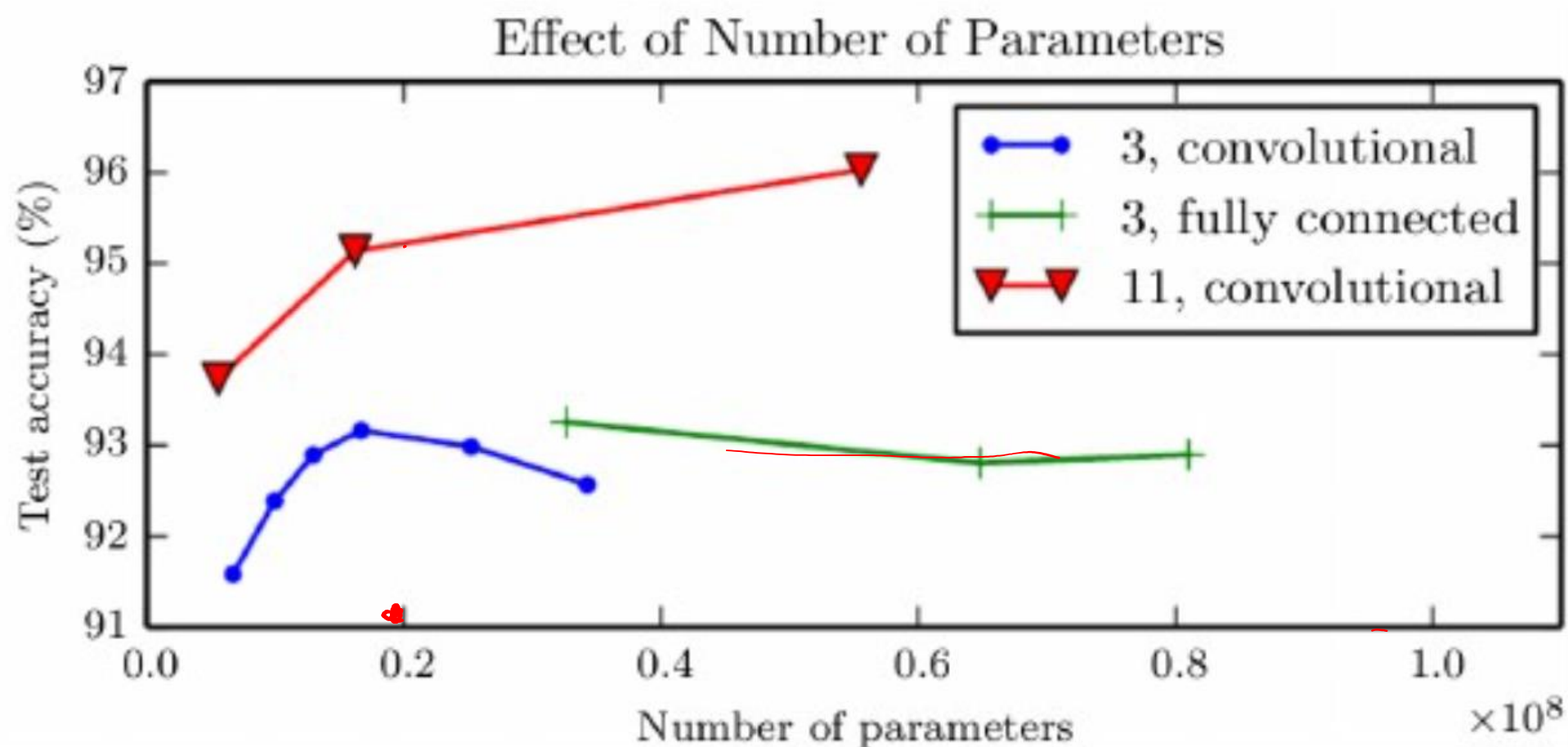
Network Architecture

Choosing the number of layers and width of the network and connection between layer

- **Universal approximation theorem:** A network with one hidden layer (sigmoid type activation) can approximate any continuous function from a closed and bounded set given enough hidden units.
- Proof also extended to work for RELU activations.
- Not useful in practice:
 - number of hidden units required may be exponentially large,
 - the parameters of the network may not be easily learnable: might overfit on a wrong function.

Effect of depth

- Many functions can be efficiently represented with multiple hidden layers but require exponential width with single hidden layer
- The number of linear regions carved out via d inputs, $l+1$ depth, $m = c$ units per hidden layer is $O((c+1)^d c^d)$
- Empirically too, larger depth leads to better generalization and lower error.



Stochastic Gradient Descent (SGD)

SGD Algorithm

Inputs:

Function NN(x; θ), Training examples, $x_1 \dots x_n$ and outputs, $y_1 \dots y_n$ and Loss function L .

do until stopping criterion

Pick a training example x_i, y_i

Compute the loss $L(\text{NN}(x_i; \theta), y_i)$

Compute gradient of L , ∇L with respect to θ

$\theta \leftarrow \theta - \eta \nabla L$

done

Return: θ

$\rightarrow N$

\rightarrow pick a batch of b examples

\uparrow learning rate

Training a Neural Network

Define the **Loss function** to be minimised as a node L

Goal: Learn weights for the neural network which minimise L

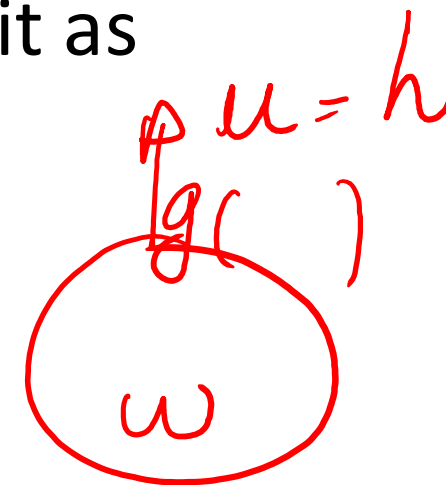
Gradient Descent: Find $\frac{\partial L}{\partial w}$ for every weight w , and update it as
 $w \leftarrow w - \eta \frac{\partial L}{\partial w}$

How do we efficiently compute $\frac{\partial L}{\partial w}$ for all w ?

Will compute $\frac{\partial L}{\partial u}$ for every node u in the network!

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \cdot \frac{\partial u}{\partial w}$$

easy, local (pointing to $\frac{\partial u}{\partial w}$)
compute recursively (pointing to $\frac{\partial L}{\partial u}$)



Training a Neural Network

New goal: compute $\partial L / \partial u$ for every node u ^{*h*} in the network

Simple algorithm: Backpropagation

Key fact: Chain rule of differentiation

If L can be written as a function of variables v_1, \dots, v_n , which in turn depend (partially) on another variable u , then

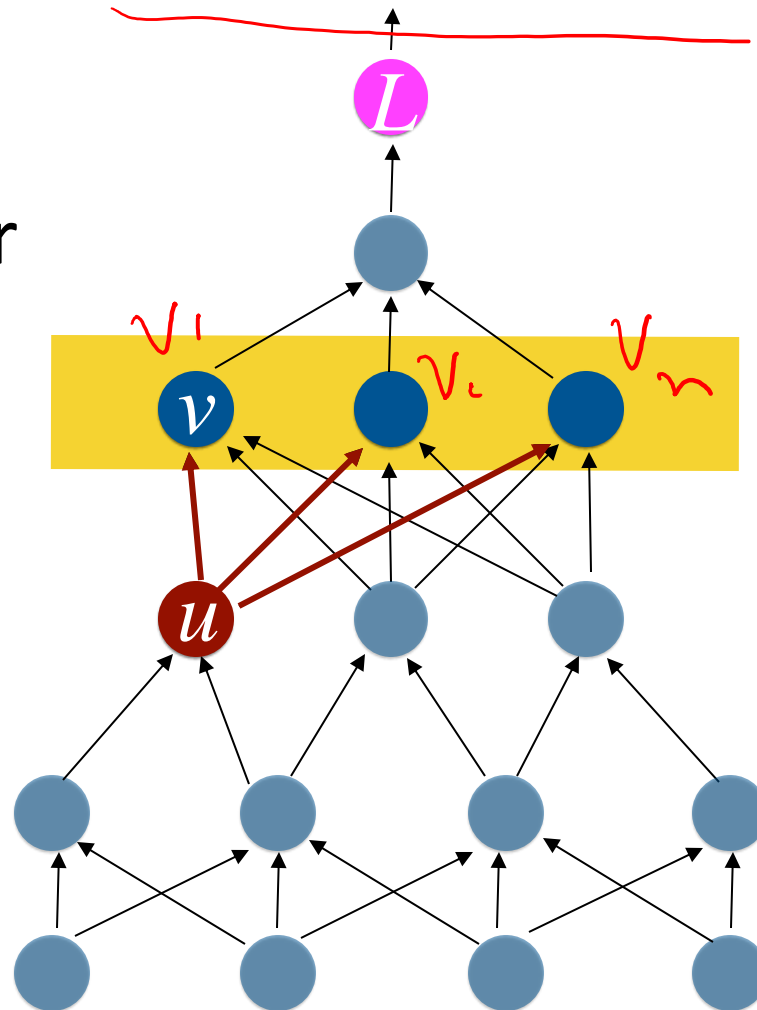
$$\partial L / \partial u = \sum_i \partial L / \partial v_i \cdot \partial v_i / \partial u$$

Backpropagation

If L can be written as a function of variables v_1, \dots, v_n , which in turn depend (partially) on another variable u , then

$$\partial L / \partial u = \sum_i \partial L / \partial v_i \cdot \partial v_i / \partial u$$

Consider v_1, \dots, v_n as the layer above u , $\Gamma(u)$



Then, the chain rule gives

$$\partial L / \partial u = \sum_{v \in \Gamma(u)} \partial L / \partial v \cdot \partial v / \partial u$$

Backpropagation

$$\partial L / \partial u = \sum_{v \in \Gamma(u)} \partial L / \partial v \cdot \partial v / \partial u$$

Backpropagation

Base case: $\partial L / \partial L = 1$

For each u (top to bottom):

For each $v \in \Gamma(u)$:

Inductively, have computed $\partial L / \partial v$

Directly compute $\partial v / \partial u$

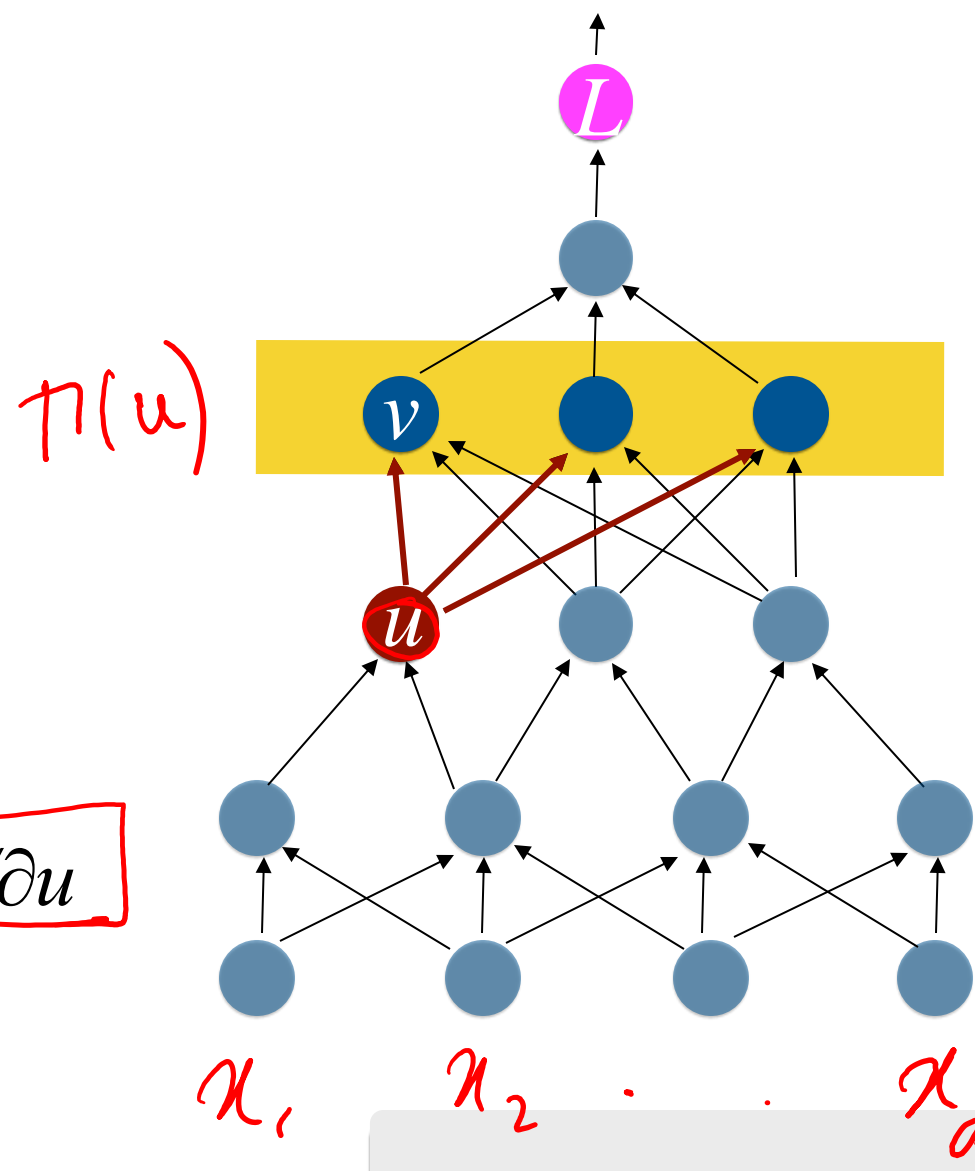
Compute $\partial L / \partial u$

Compute $\partial L / \partial w$

where $\partial L / \partial w = \partial L / \partial u \cdot \partial u / \partial w$

Forward Pass

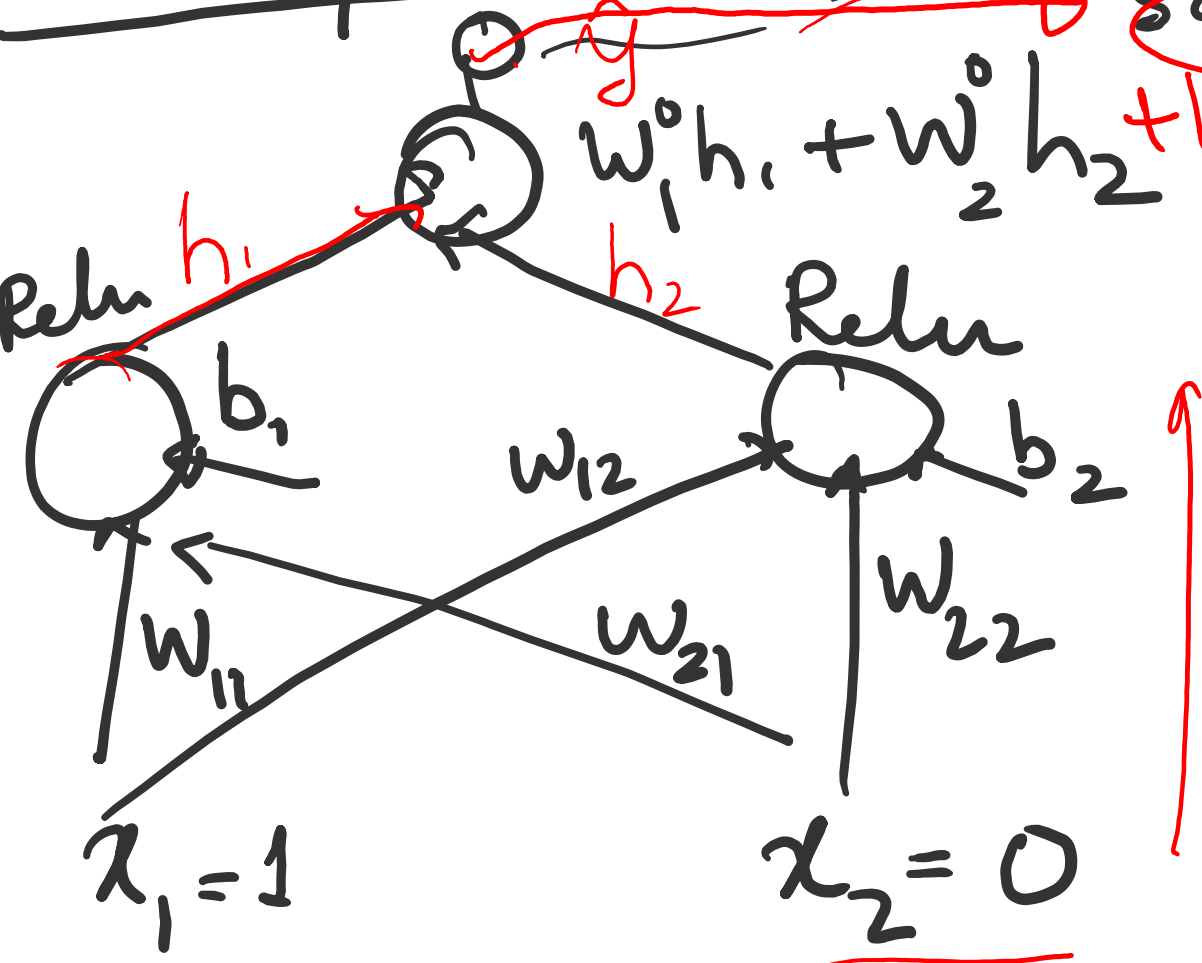
First, in a forward pass, compute values of all nodes given an input (The values of each node will be needed during backprop)



Where values computed in the forward pass are needed

Example XOR

square loss



Initially: $w = 0$

all $b = 1$
 $b^0 = 0, b_1 = 1, b_2 = -1$

Fwd:

$h^1 = 1, h^2 = 0$

$\hat{y} = 0 = w^0_1 h_1 + w^0_2 h_2 + b^0$

$L = \text{loss} : (y - \hat{y})^2 = 1.$

Bwd

$\frac{\partial L}{\partial g} = -2$

$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial g} \frac{\partial g}{\partial h_1} = -2 \cdot w^0_1$

$h_1 = \text{Relu}(w_{11}x_1 + w_{21}x_2 + b_1)$

$\frac{\partial L}{\partial w^0_1} = \frac{\partial L}{\partial g} \frac{\partial g}{\partial w^0_1} = -2h_1 = -2$

$\Delta w^0_1 = 0 - \eta(-2) = 2\eta$

$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial h_1} \frac{\partial h_1}{\partial w_{11}} = 0 \cdot 1 \cdot x_1$