

Contextualised Representation & Intro to Neural Nets

Intro to NLP

Rahul Mishra

IIIT-Hyderabad
Feb 13, 2024

INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

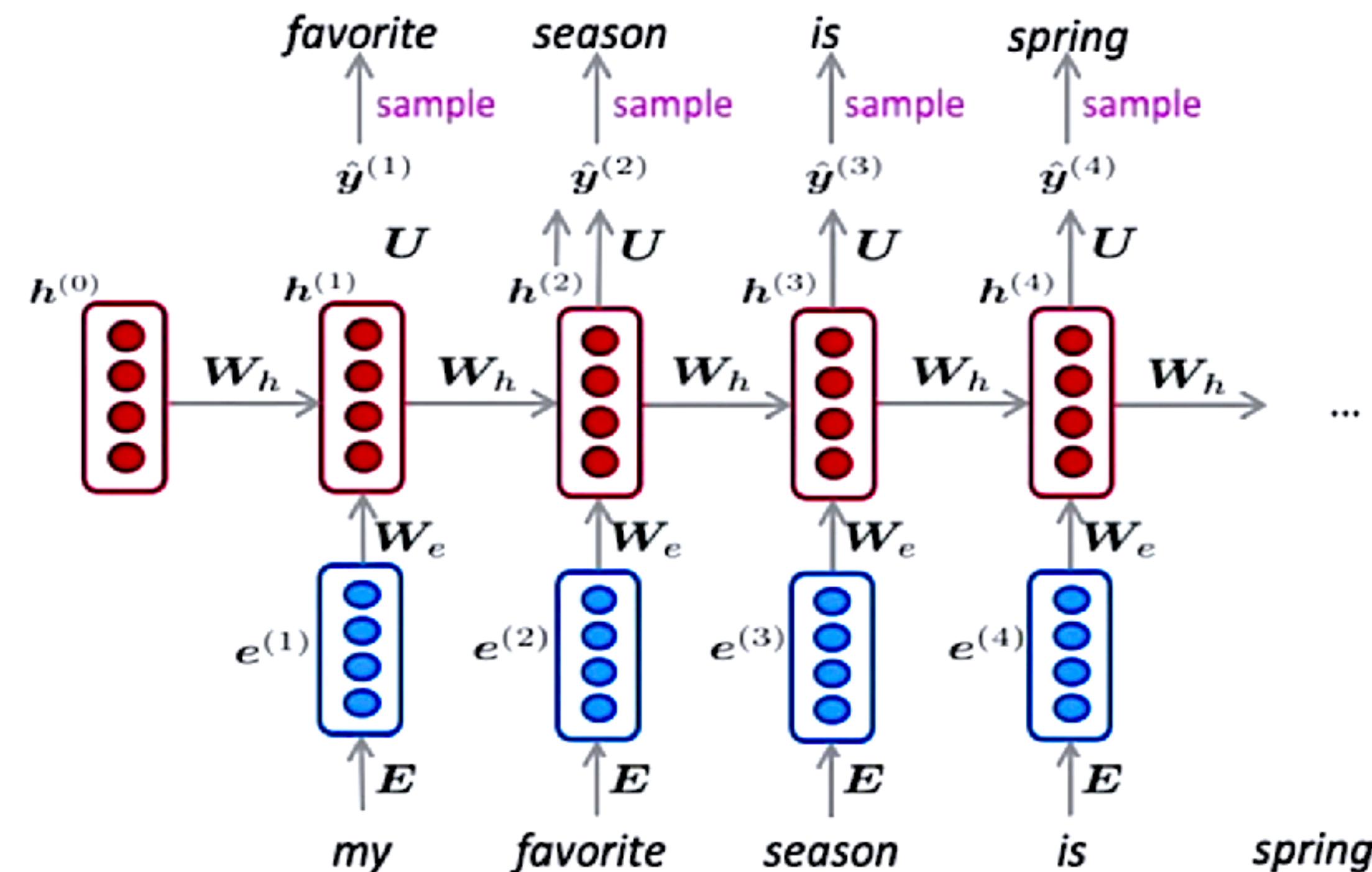
H Y D E R A B A D

Problem with Word2vec and similar methods

These have two problems:

- Always the same representation for a **word type** regardless of the context in which a **word token** occurs
 - We might want very fine-grained word sense disambiguation
- We just have one representation for a word, but words have different **aspects**, including semantics, syntactic behavior, and register/connotations

Do we have a simple solution?

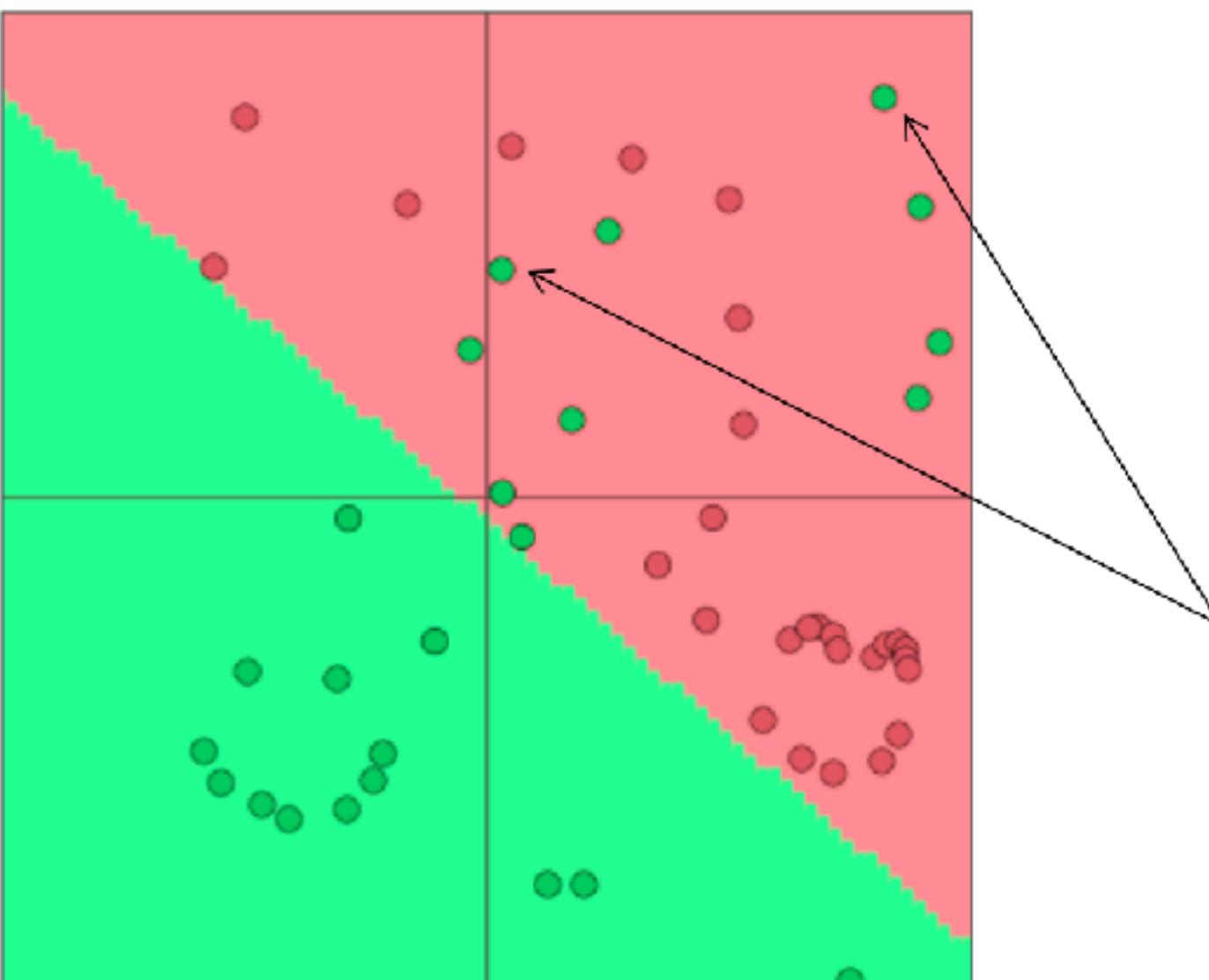


ELMo
ULMfit
BERT
GPT

We will revisit this slide after covering basics of
neural nets and RNNs

Neural Networks

- Softmax (\approx logistic regression) alone not very powerful
- Softmax gives only linear decision boundaries



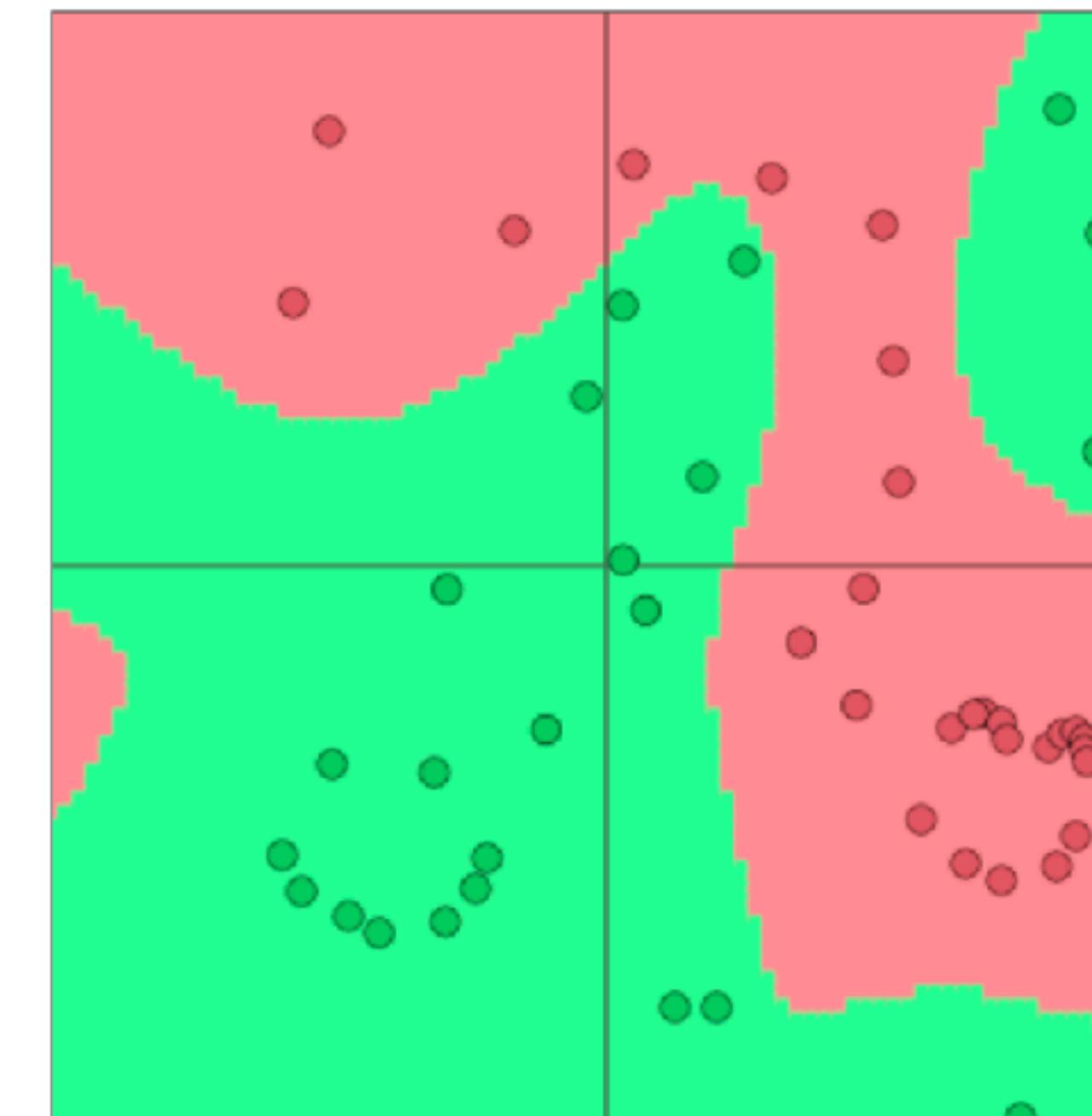
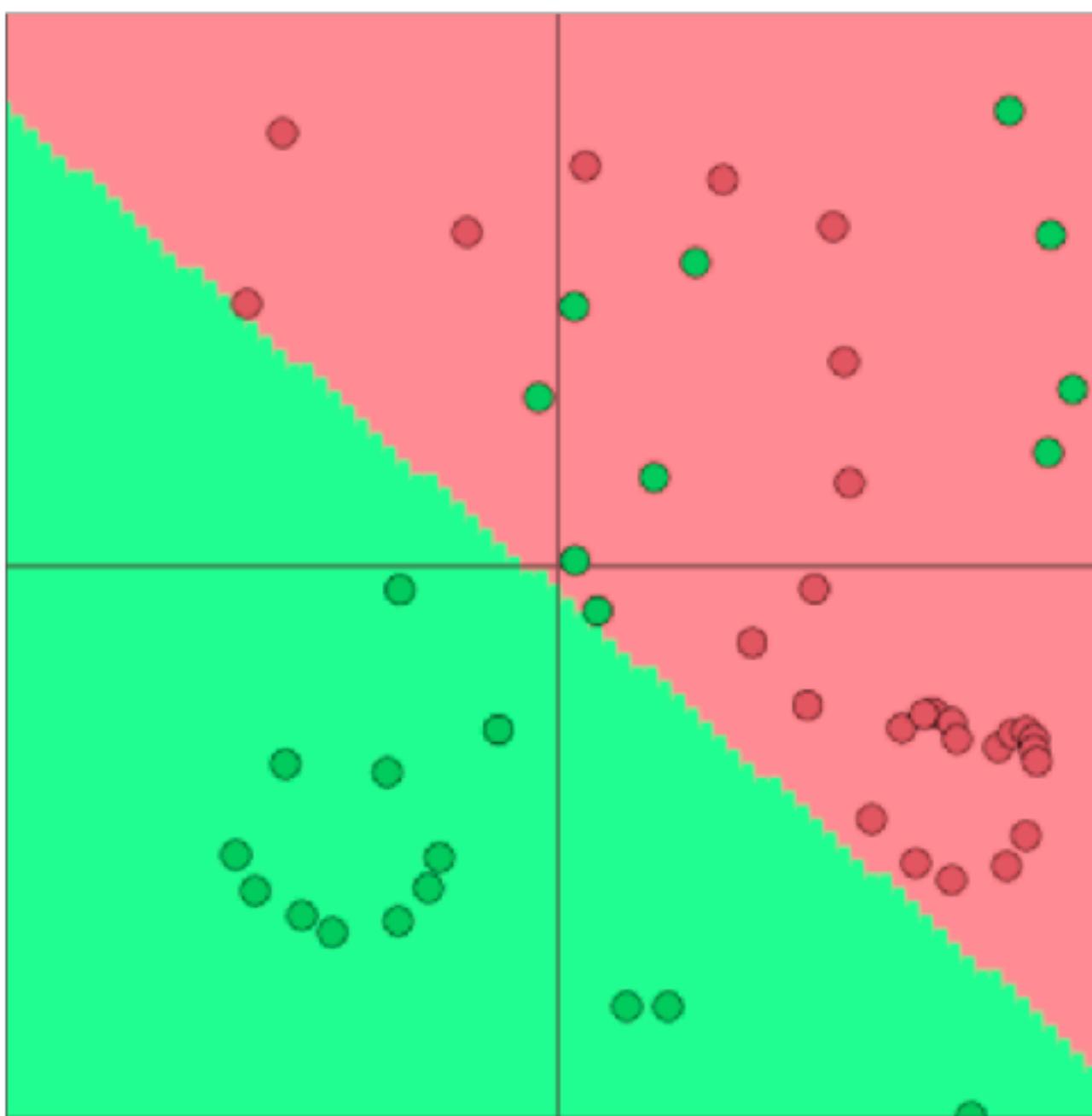
This can be quite limiting

→ Unhelpful when a
problem is complex

Wouldn't it be cool to
get these correct?

NNs to the rescue

- Neural networks can learn much more complex functions and nonlinear decision boundaries!
 - In original space

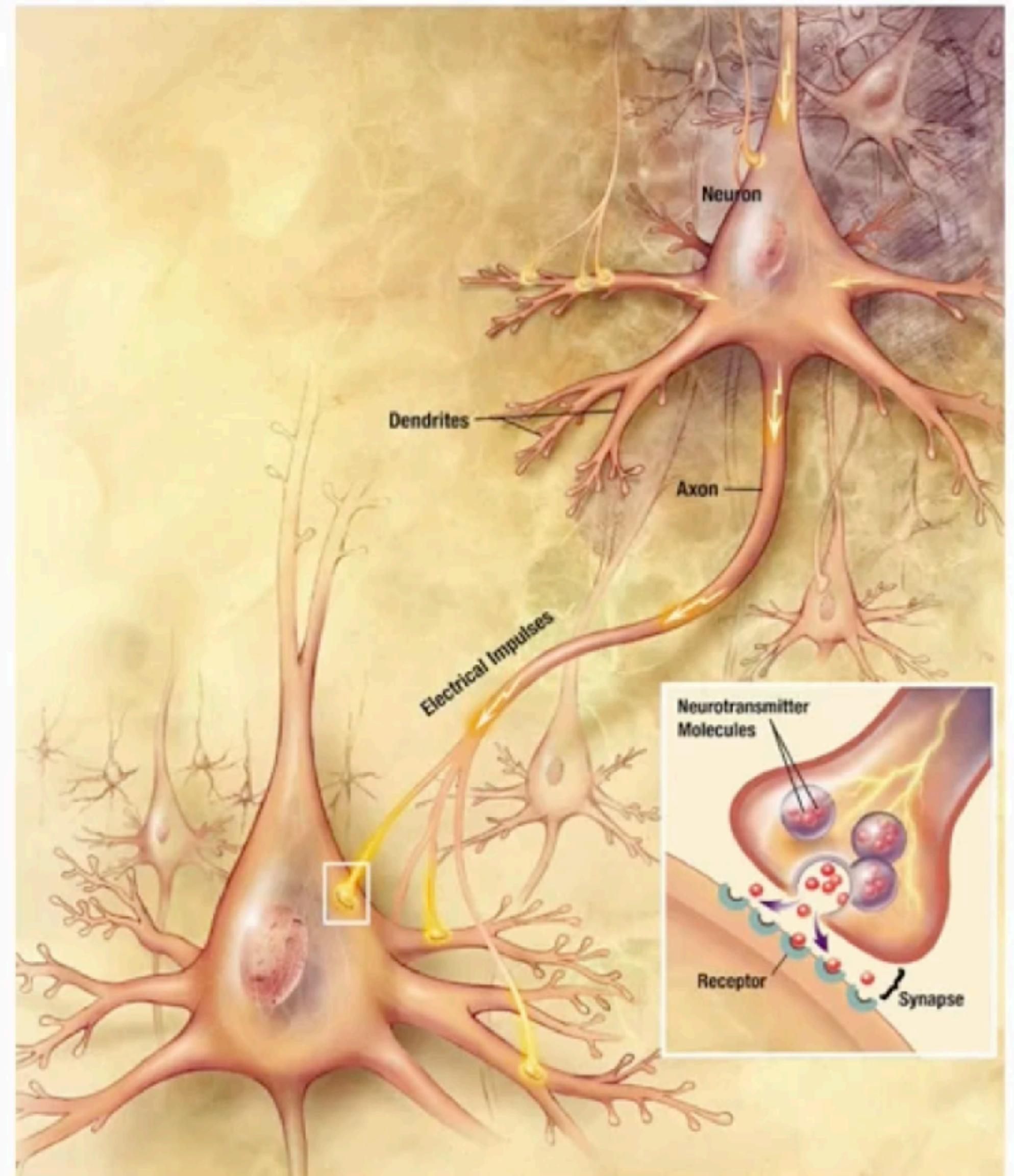


Biological Neuron

Origins: Algorithms that try to mimic the brain.

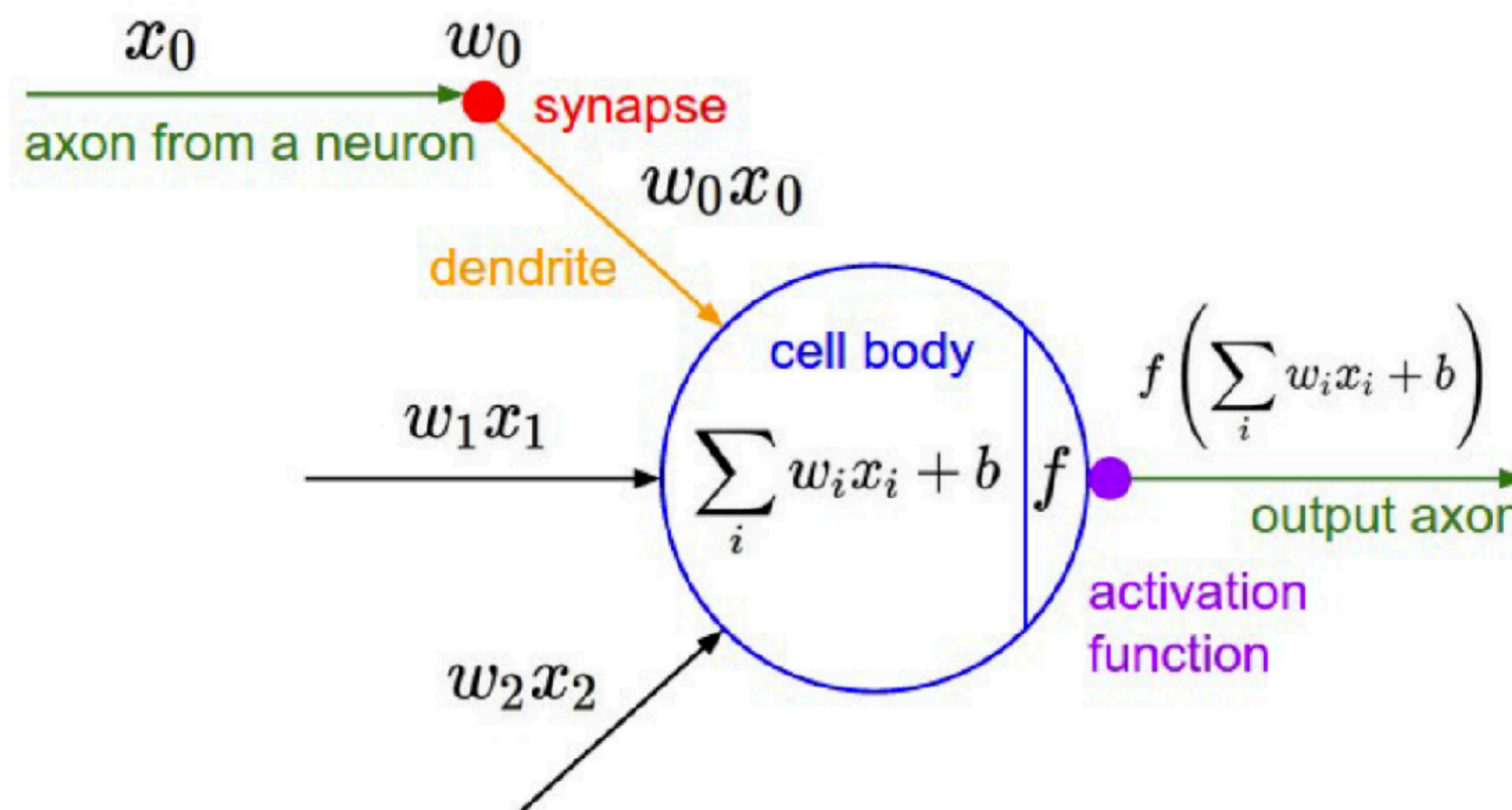
Used in the 1980's and early 1990's.
Fell out of favor in the late 1990's.

Resurgence from around 2005.



Artificial Neural Net

- Neural networks come with their own terminological baggage
- But if you understand how softmax models work, then you can easily understand the operation of a neuron!



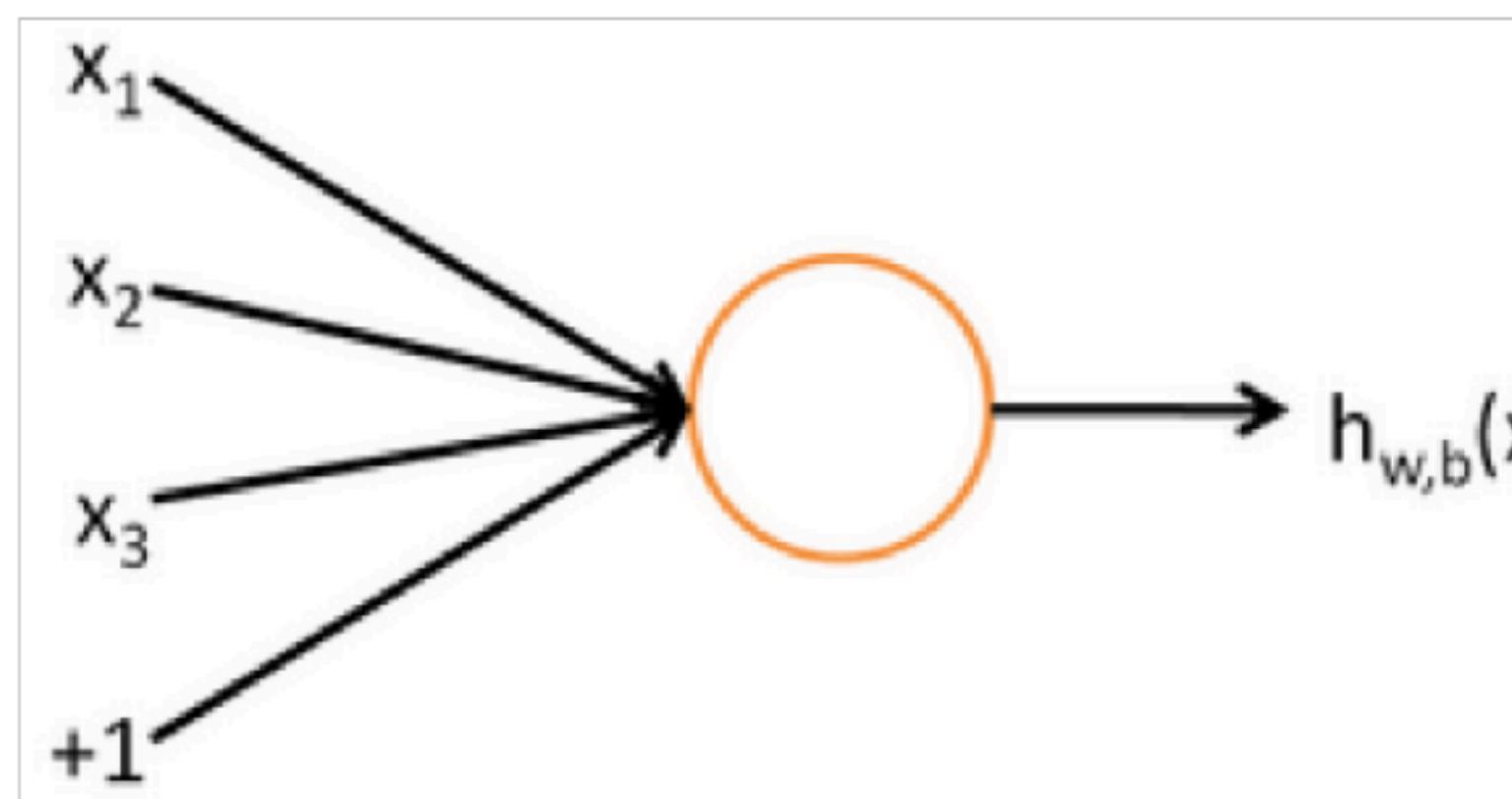
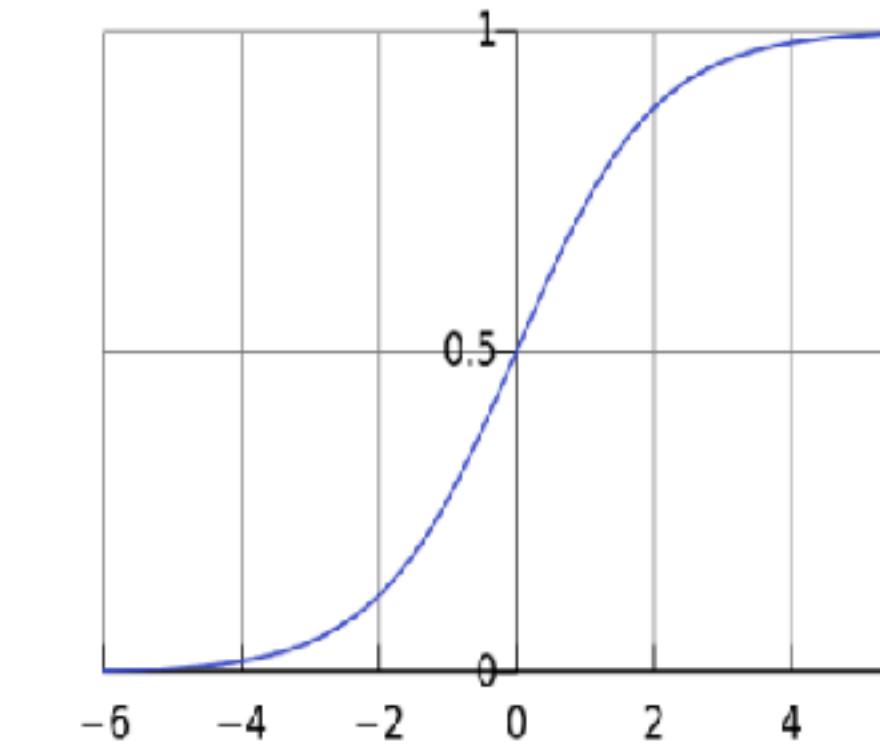
A neuron can be a binary logistic regression unit

f = nonlinear activation fct. (e.g. sigmoid), w = weights, b = bias, h = hidden, x = inputs

$$h_{w,b}(x) = f(w^T x + b)$$

b : We can have an “always on” feature, which gives a class prior, or separate it out, as a bias term

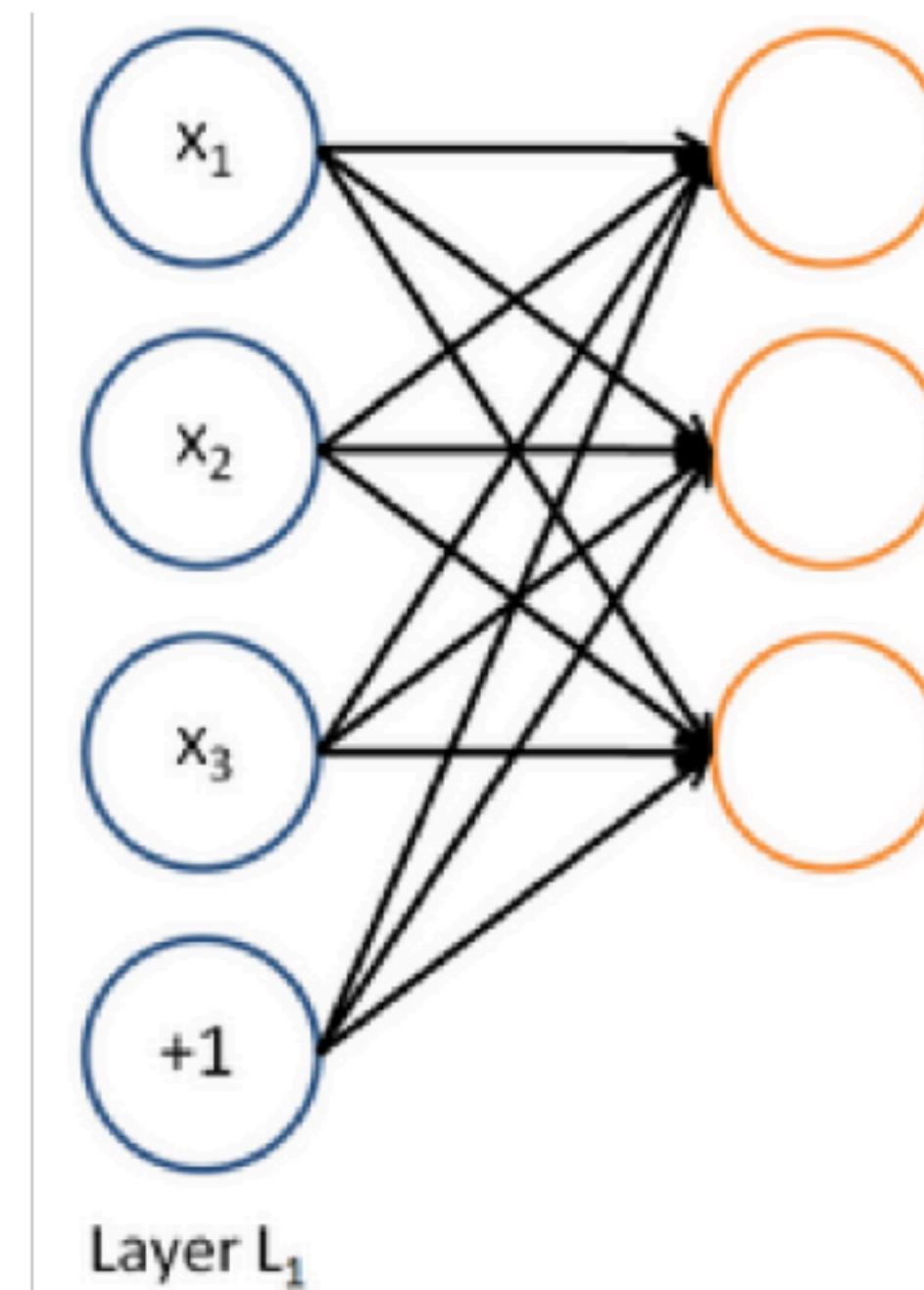
$$f(z) = \frac{1}{1 + e^{-z}}$$



w , b are the parameters of this neuron
i.e., this logistic regression model

Running several logistic regression

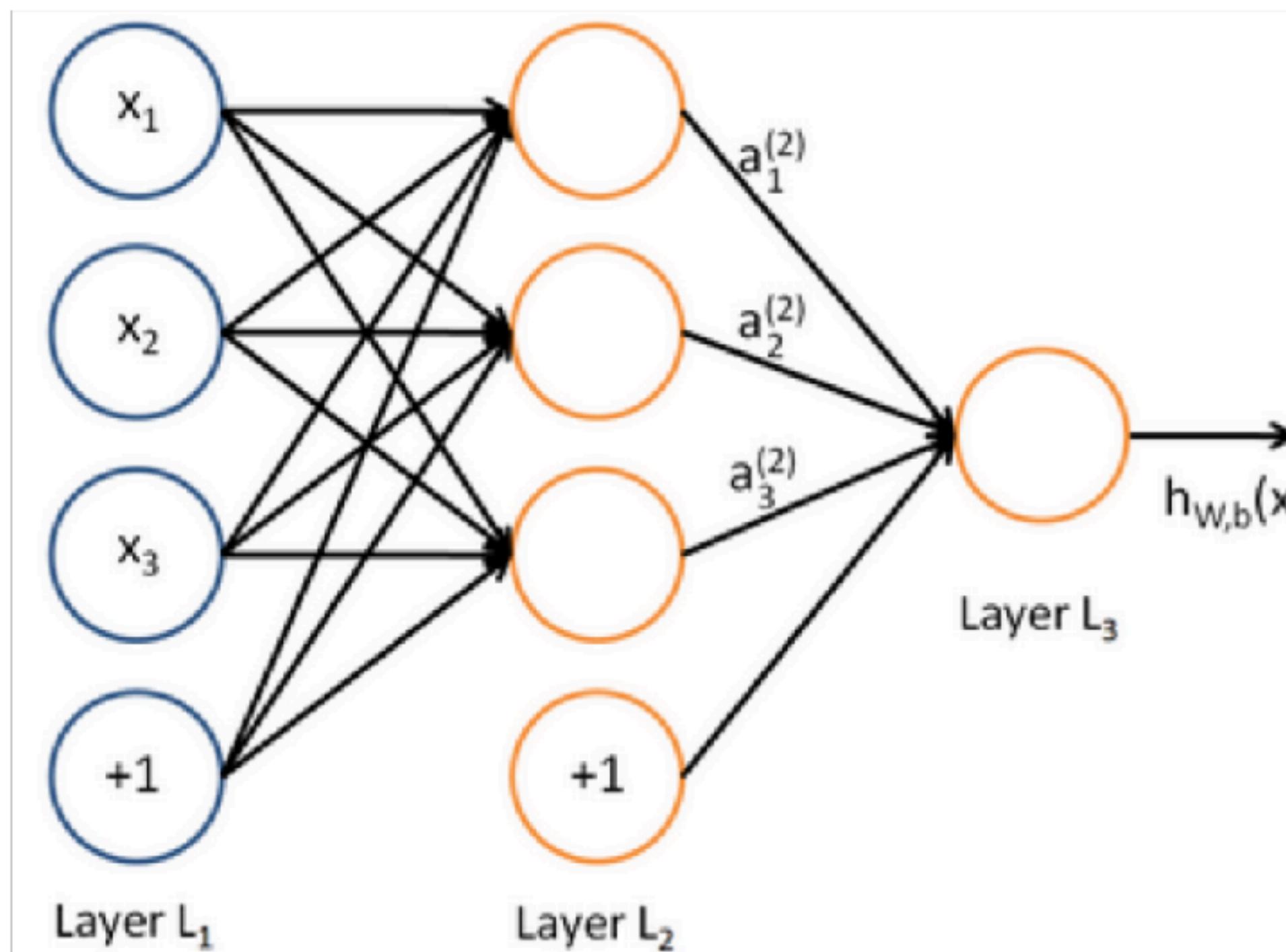
If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...



But we don't have to decide ahead of time what variables these logistic regressions are trying to predict!

Running several logistic regression

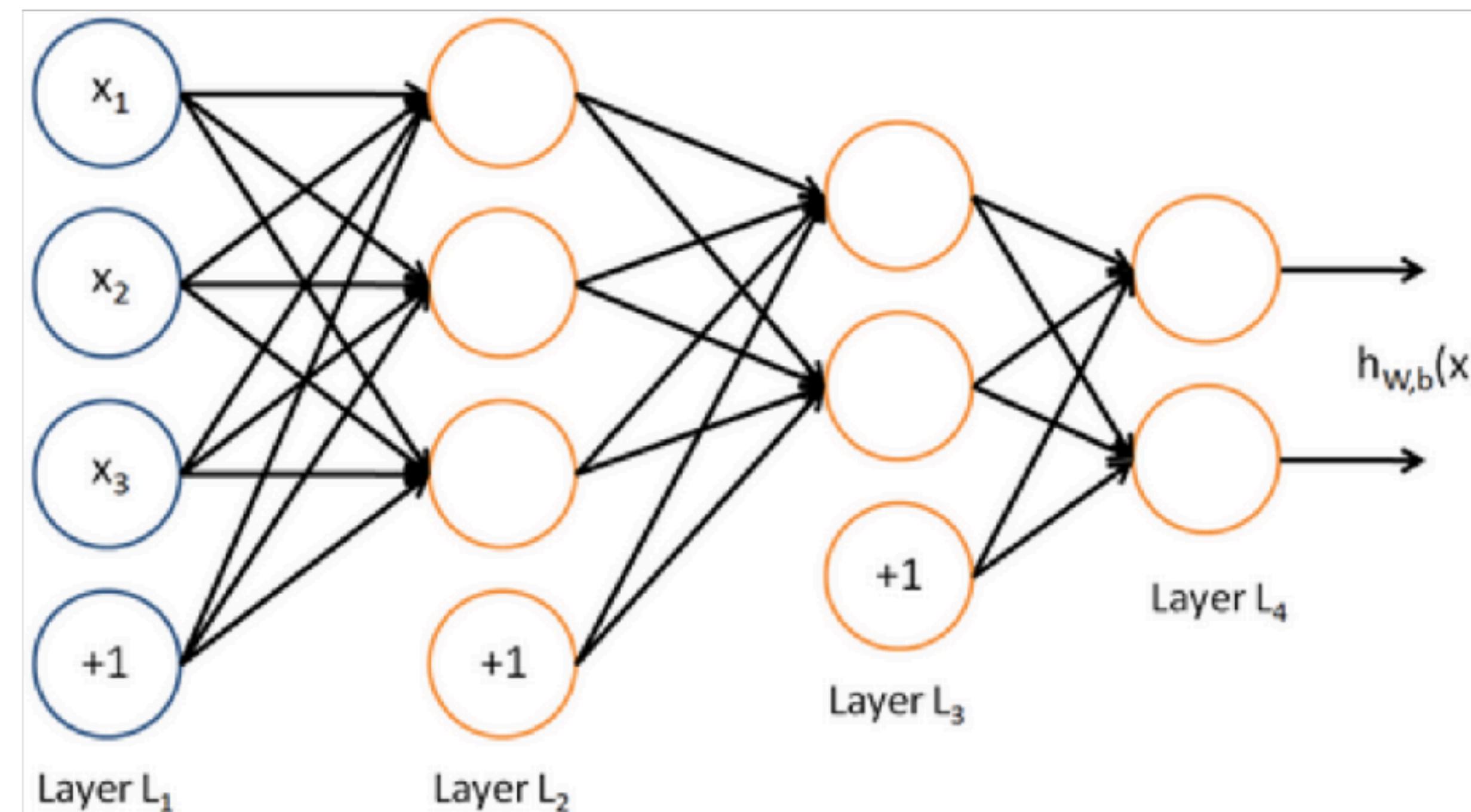
... which we can feed into another logistic regression function



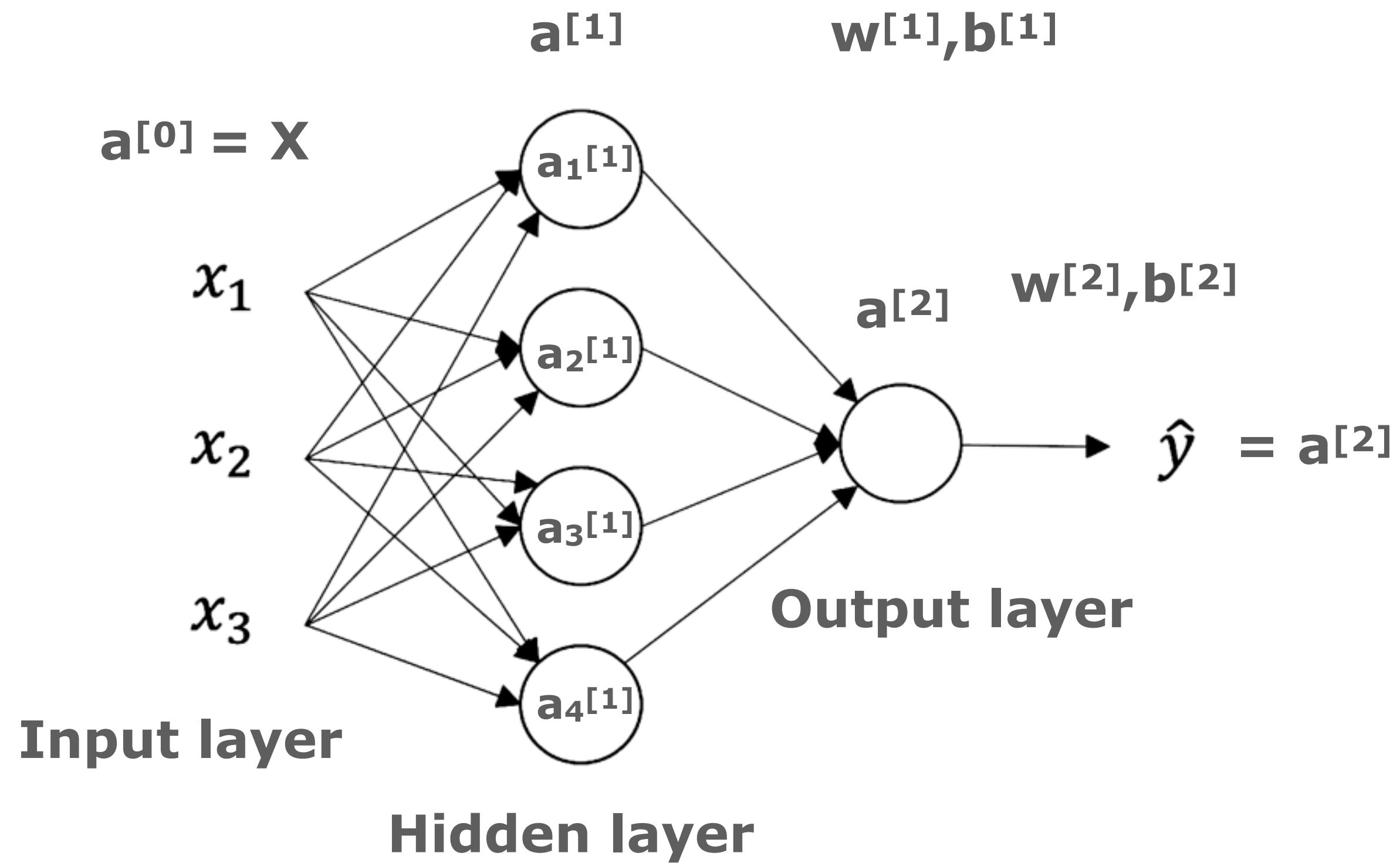
It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.

Running several logistic regression

Before we know it, we have a multilayer neural network....

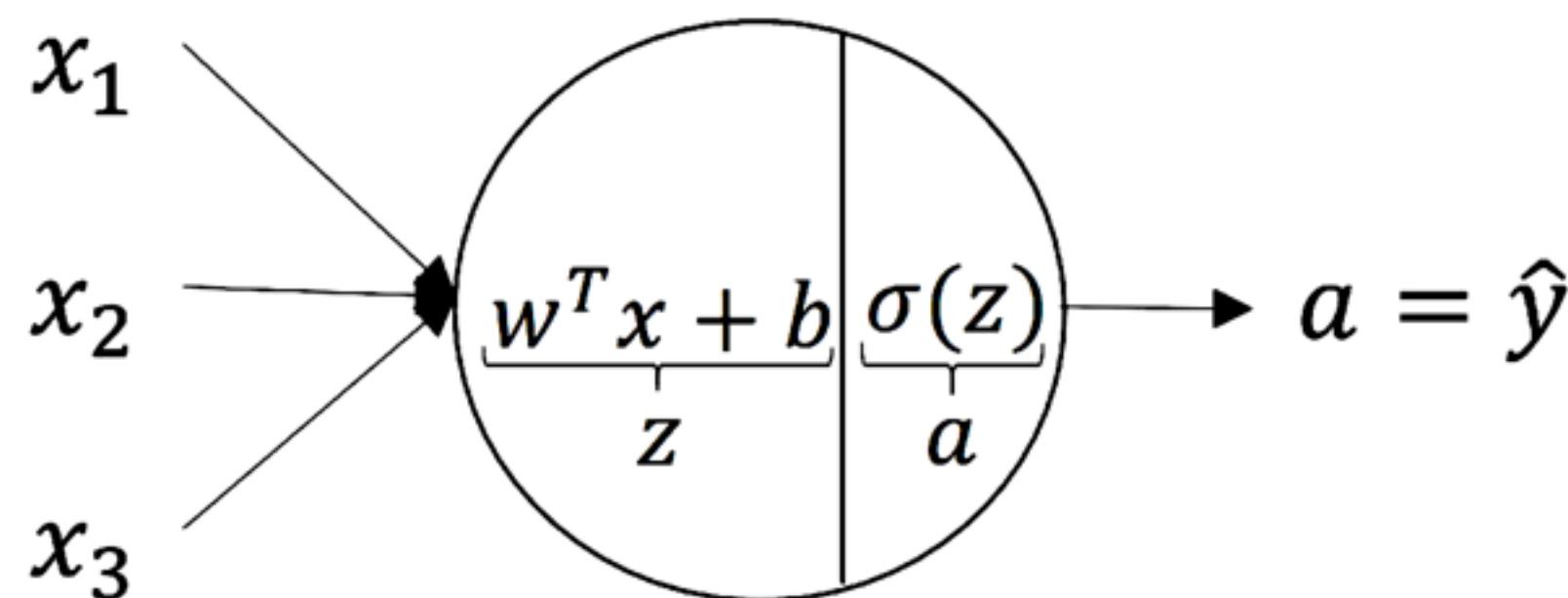


NN Representation



2 layer NN

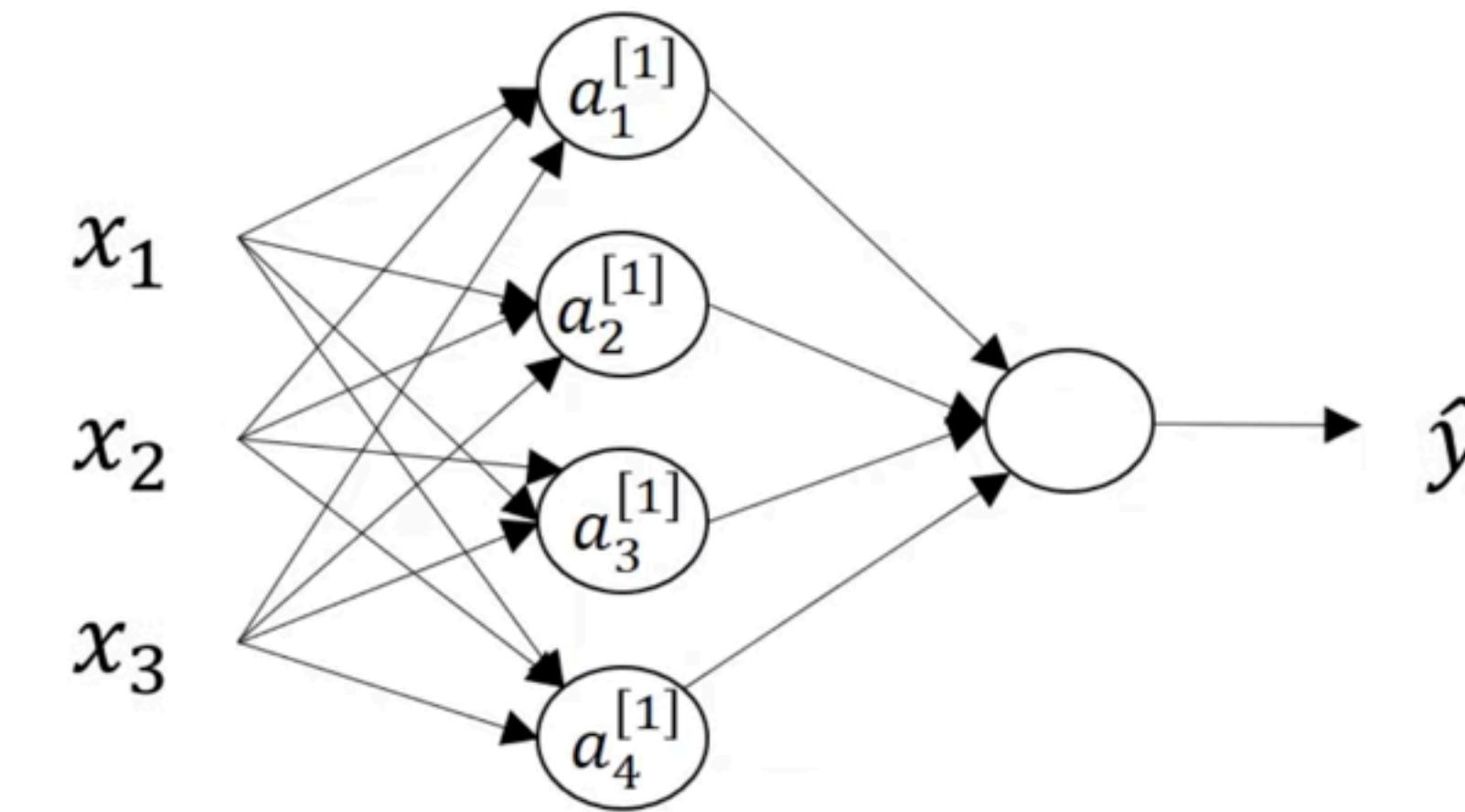
Computing NN output: Forward Pass



$$z = w^T x + b$$

$$a = \sigma(z)$$

$$\begin{bmatrix} w_1^{(1)T} \\ w_2^{(1)T} \\ w_3^{(1)T} \\ w_4^{(1)T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} = \begin{bmatrix} \rightarrow w_1^{(1)T} x + b_1^{(1)} \\ \rightarrow w_2^{(1)T} x + b_2^{(1)} \\ \rightarrow w_3^{(1)T} x + b_3^{(1)} \\ \rightarrow w_4^{(1)T} x + b_4^{(1)} \end{bmatrix}$$



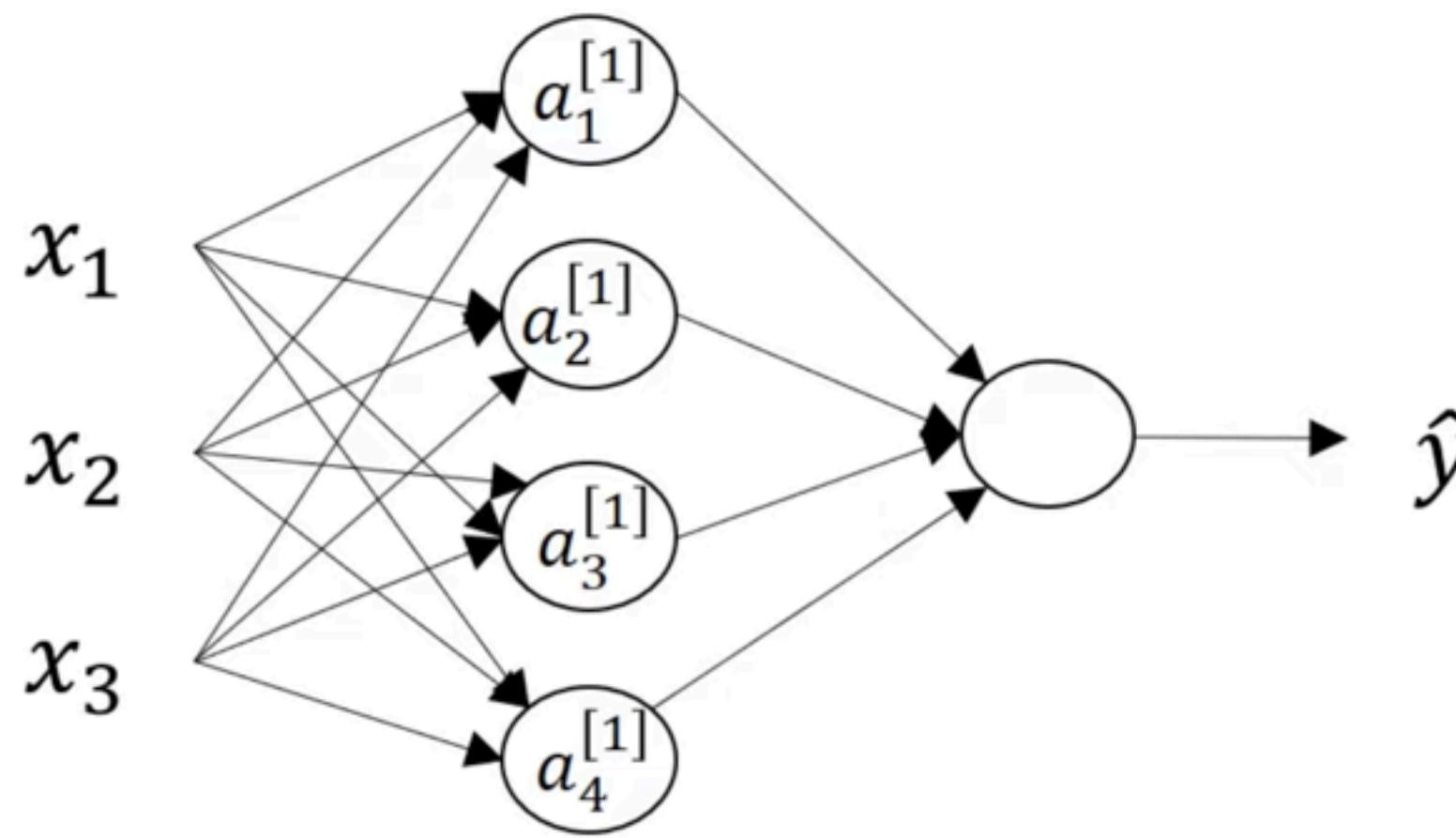
$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

NN Forward Pass



$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

Given input x:

$$z_{(4,1)}^{[1]} = W_{(4,3)}^{[1]} x + b_{(4,1)}^{[1]}$$

$$a_{(4,1)}^{[1]} = \sigma(z_{(4,1)}^{[1]})$$

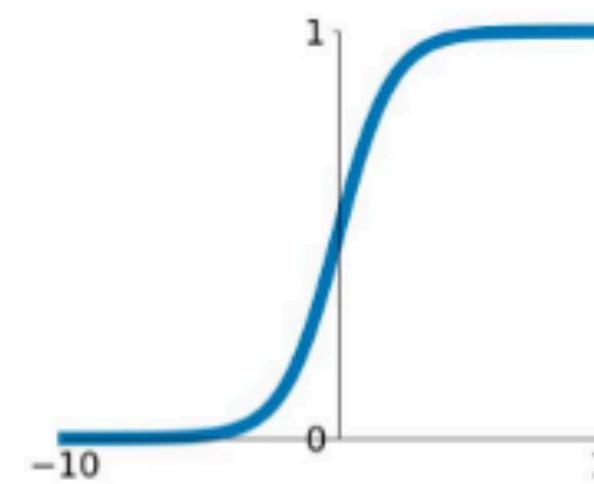
$$z_{(1,1)}^{[2]} = W_{(1,4)}^{[2]} a_{(4,1)}^{[1]} + b_{(1,1)}^{[2]}$$

$$a_{(1,1)}^{[2]} = \sigma(z_{(1,1)}^{[2]})$$

Activation Functions

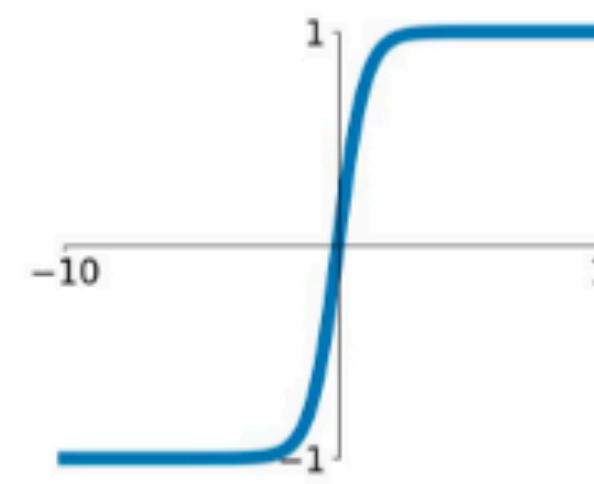
Sigmoid

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



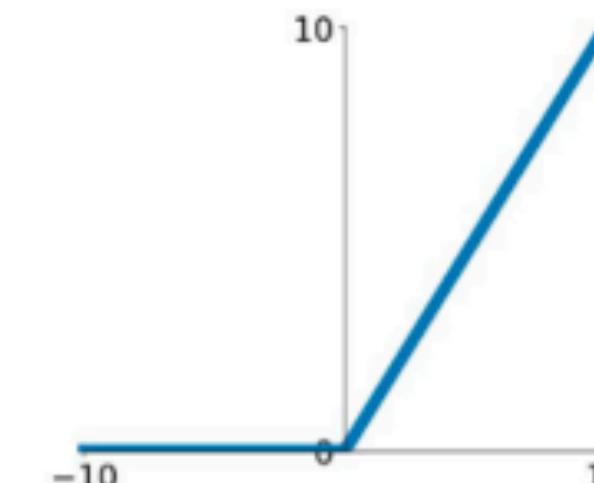
tanh

$$\tanh(x)$$



ReLU

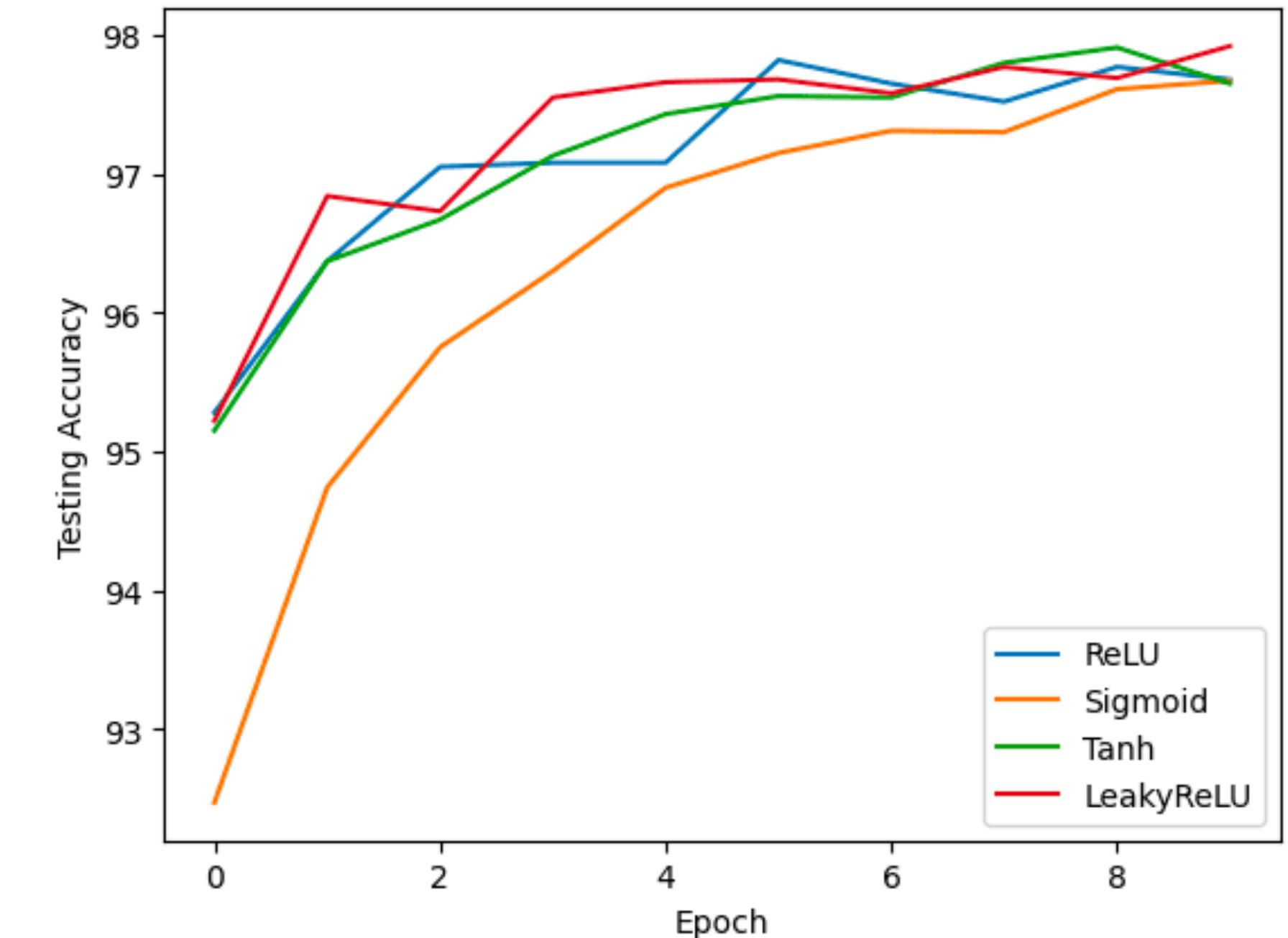
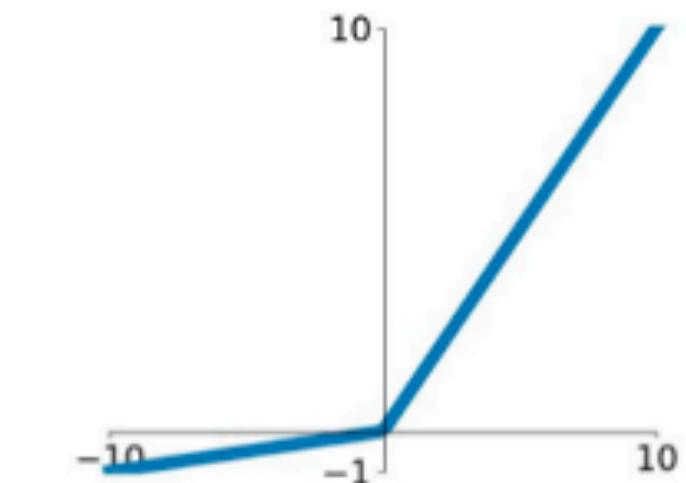
$$\max(0, x)$$



$$\tanh: a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Leaky ReLU

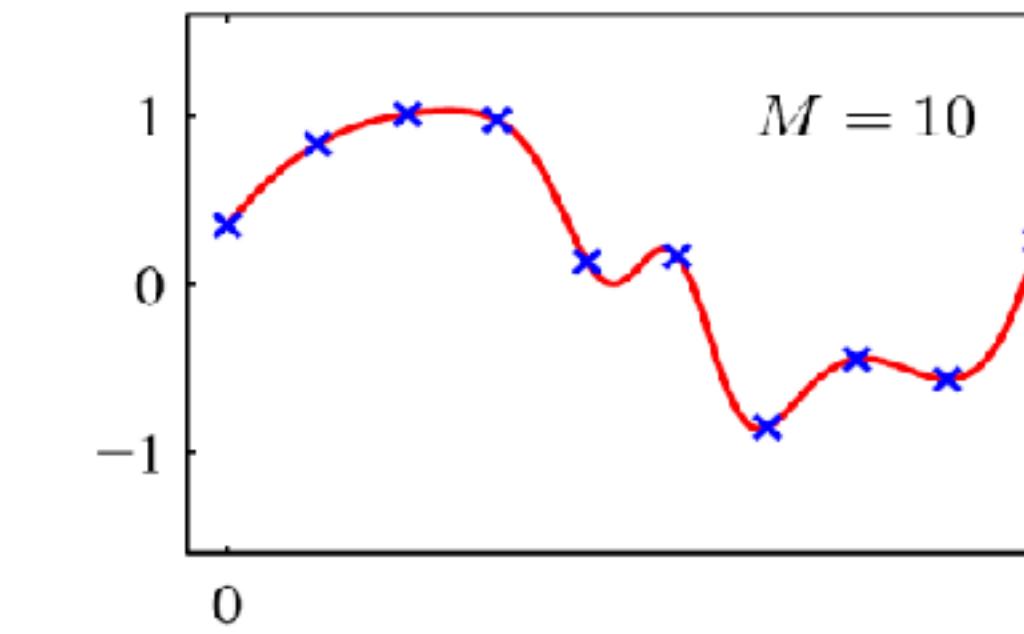
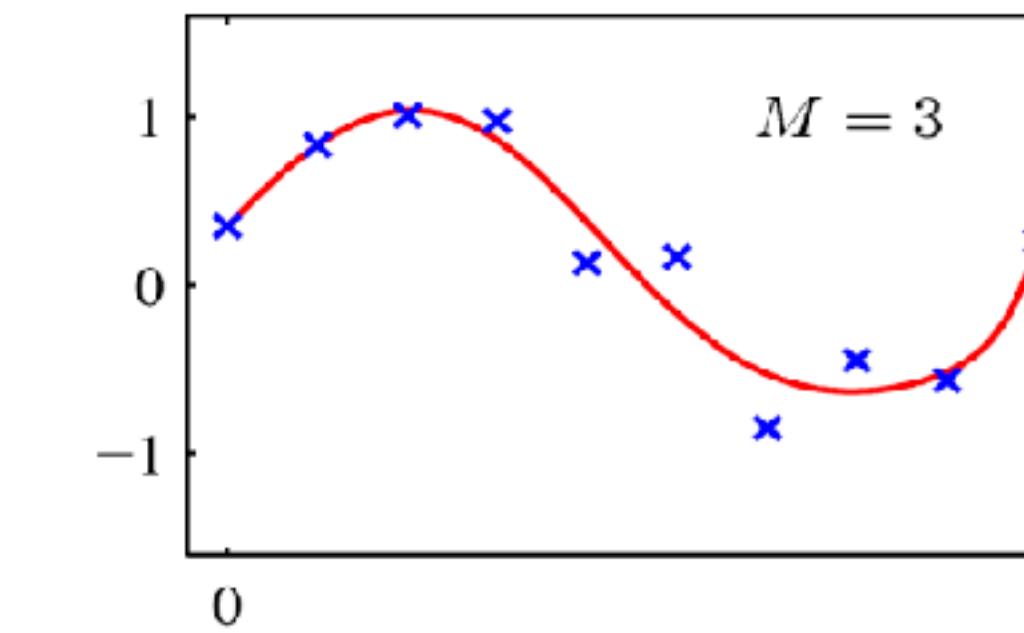
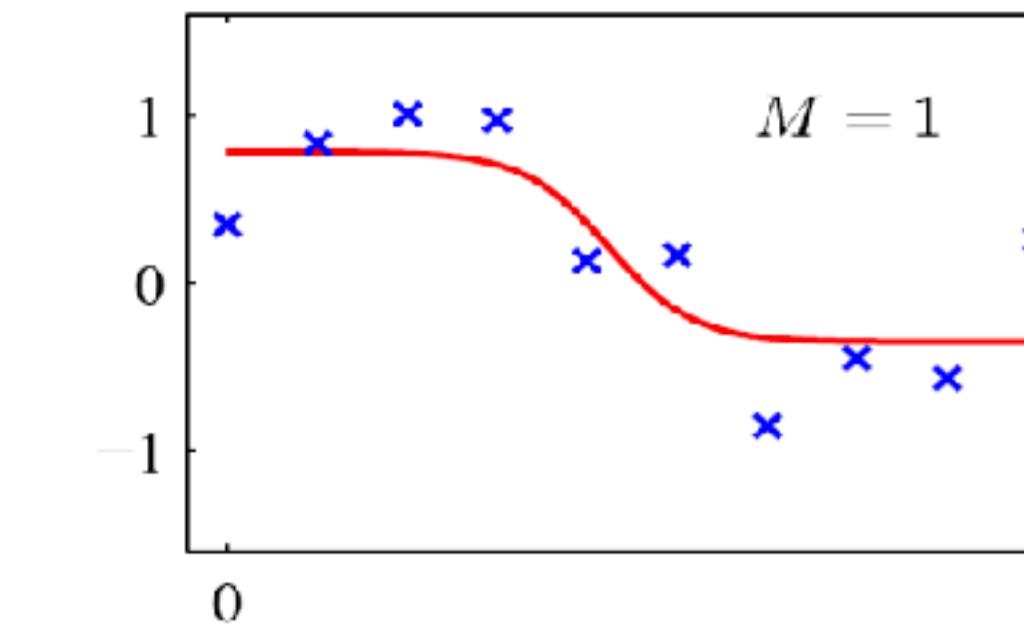
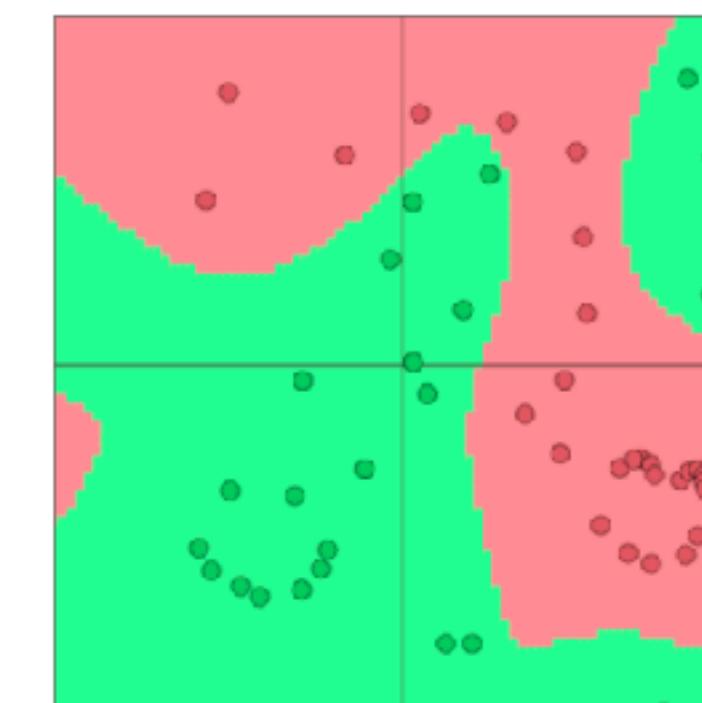
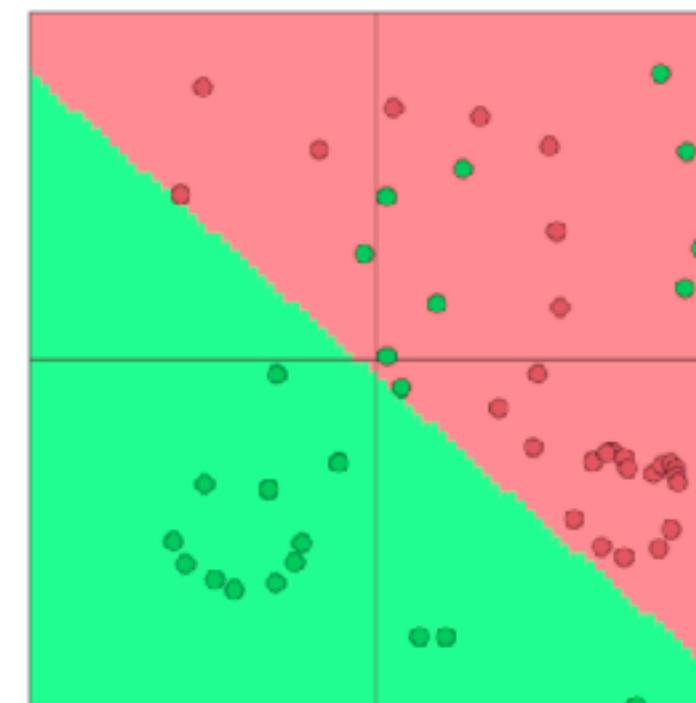
$$\max(0.1x, x)$$



Non-linearities: Why they're needed

Example: function approximation,
e.g., regression or classification

- Without non-linearities, deep neural networks can't do anything more than a linear transform
- Extra layers could just be compiled down into a single linear transform: $W_1 W_2 x = Wx$
- With more layers, they can approximate more complex functions!



Gradient Computation

$$\Rightarrow J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

$$\min_{\Theta} J(\Theta)$$

Need code to compute:

- $J(\Theta)$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

Forward Pass: for a DNN

Given one training example (x, y):

Forward propagation:

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)}) \text{ Or } b^{(2)}$$

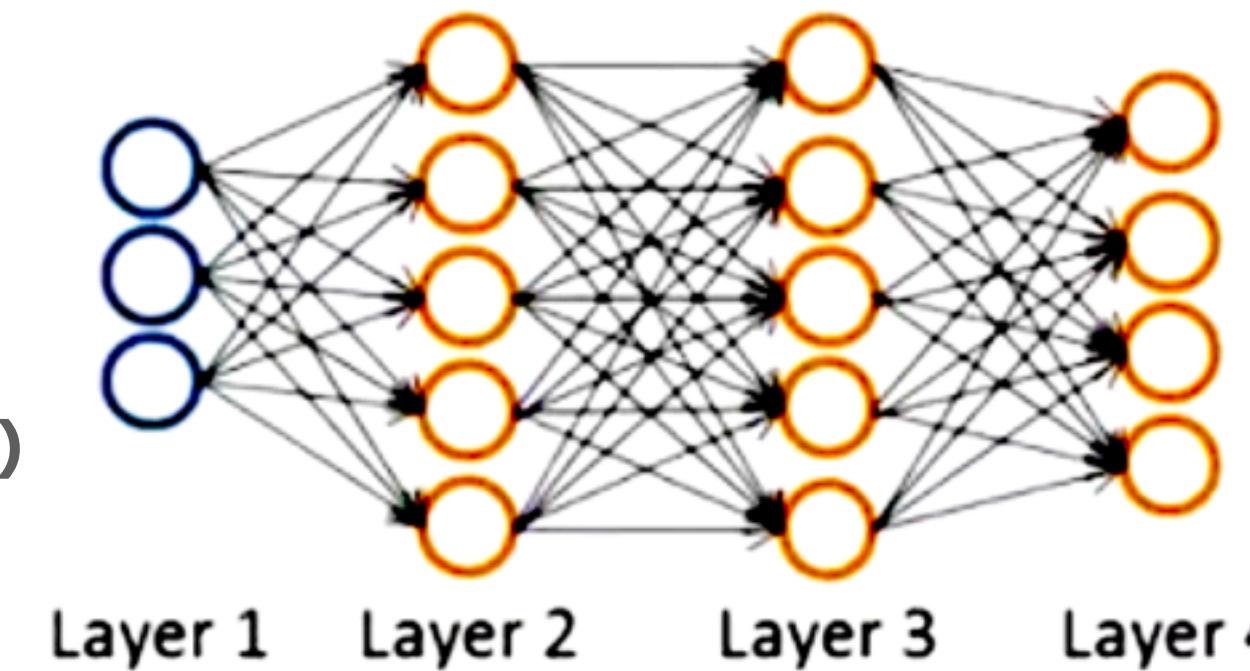
$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)}) \text{ Or } b^{(3)}$$

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$

Or \hat{y}



Backpropagation Algorithm

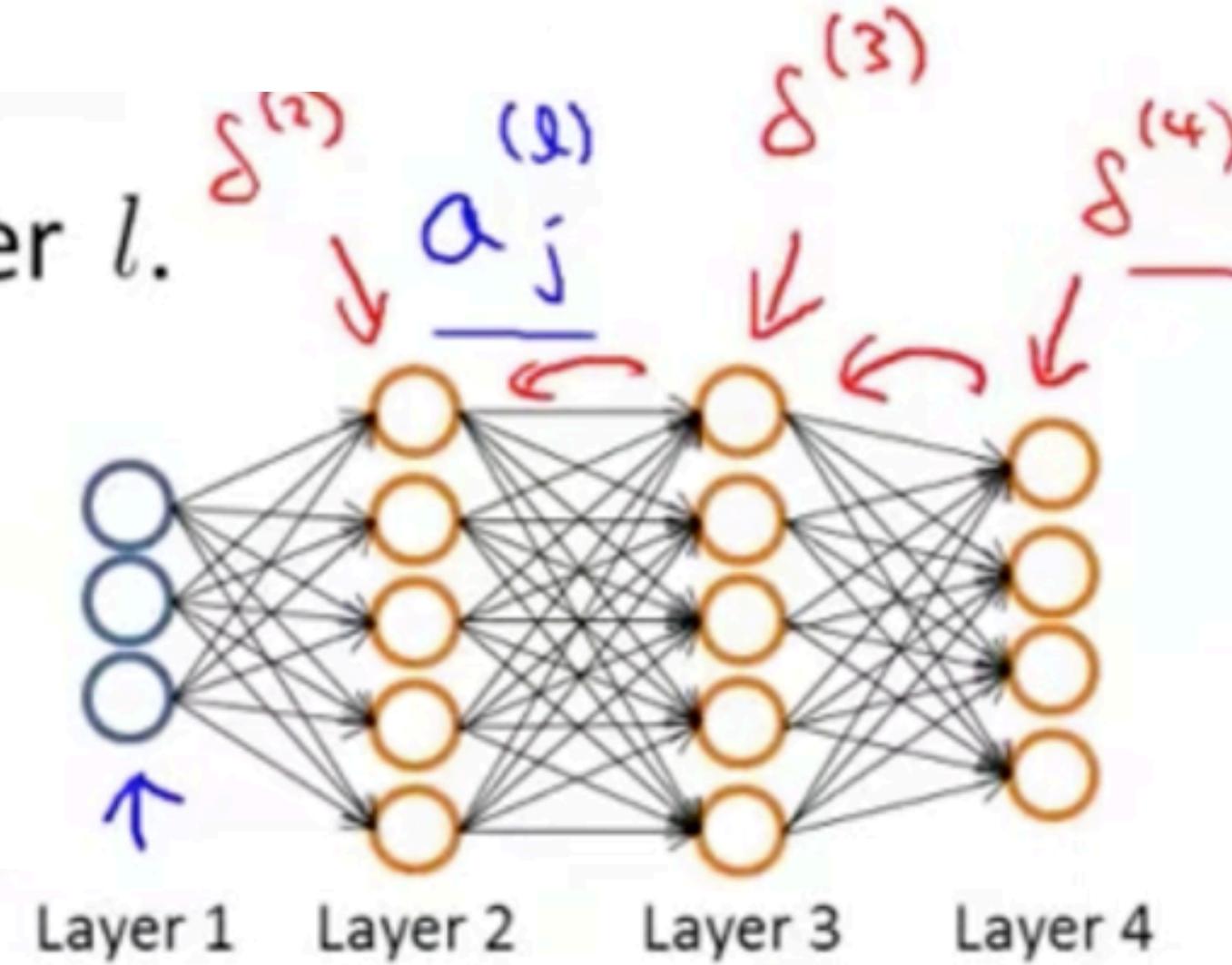
Intuition: $\delta_j^{(l)}$ = “error” of node j in layer l .

For each output unit (layer $L = 4$)

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \underline{\delta}^{(4)} = \underline{a}^{(4)} - \underline{y}$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)}) \rightarrow a^{(3)} \cdot * (1 - a^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)})$$



$$\frac{\partial}{\partial \Theta^{(k)}} J(\Theta) = a_j^{(k)} \delta_i^{(k+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = 0)$$

Backpropagation Algorithm

Training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j).

For $i = 1$ to m

Set $a^{(1)} = x^{(i)}$

Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$

Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \text{ if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \text{ if } j = 0$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

Why deep learning?

