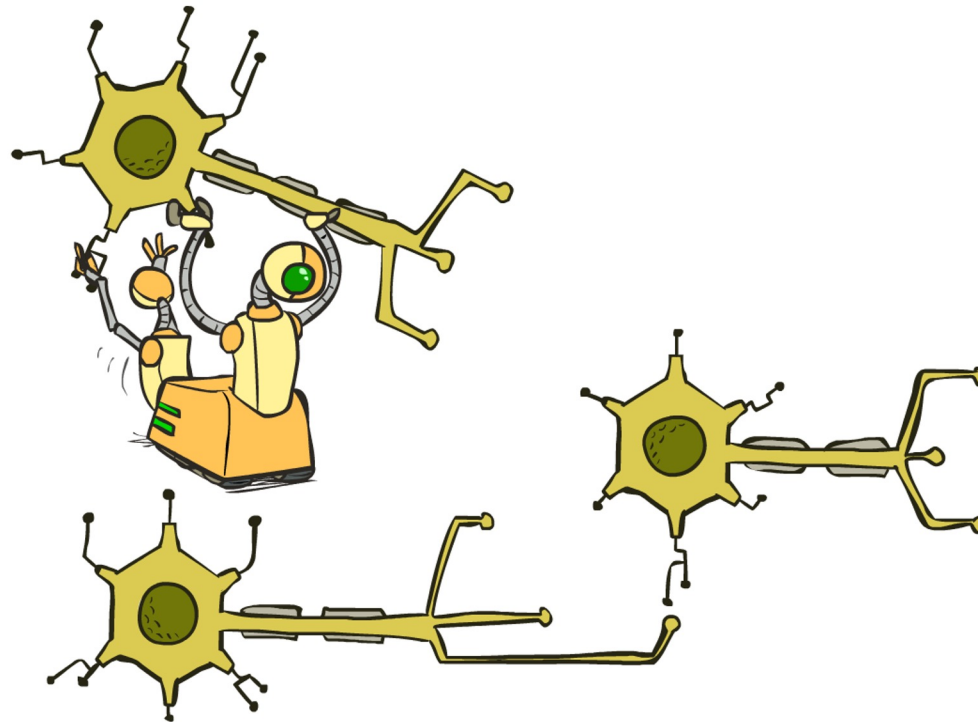# CS 188: Artificial Intelligence
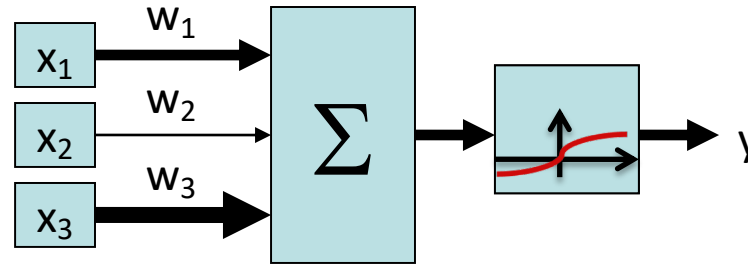
## Neural Networks
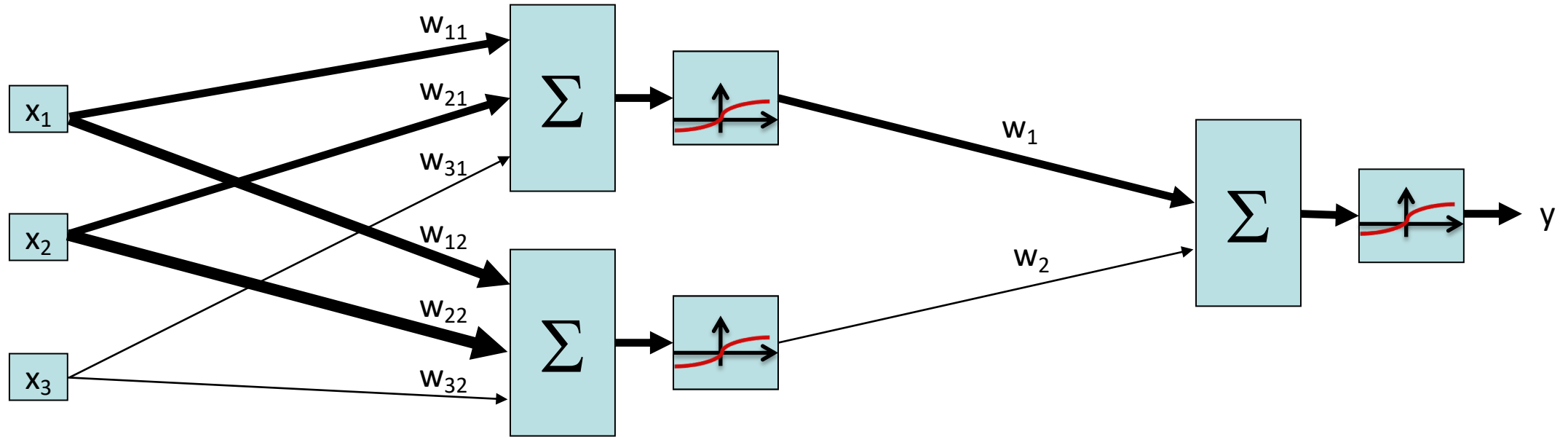
# Recall: Perceptron with Sigmoid Activation
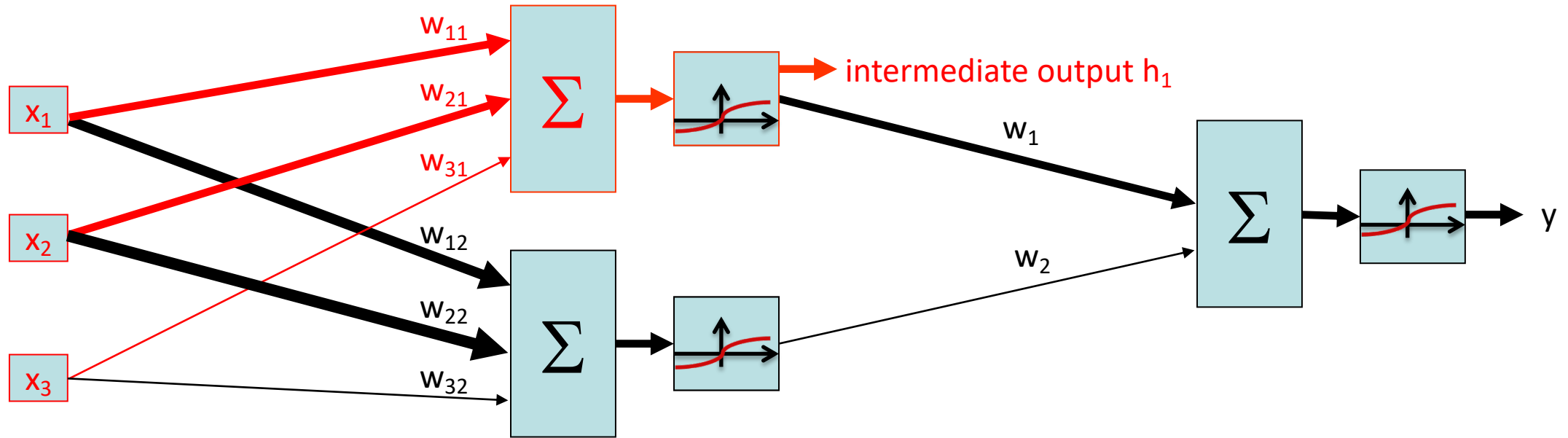


$$y = \phi(w_1 x_1 + w_2 x_2 + w_3 x_3)$$

nonlinear (probability)

$$= \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + w_3 x_3)}}$$

# Recall: 2-Layer, 2-Neuron Neural Network

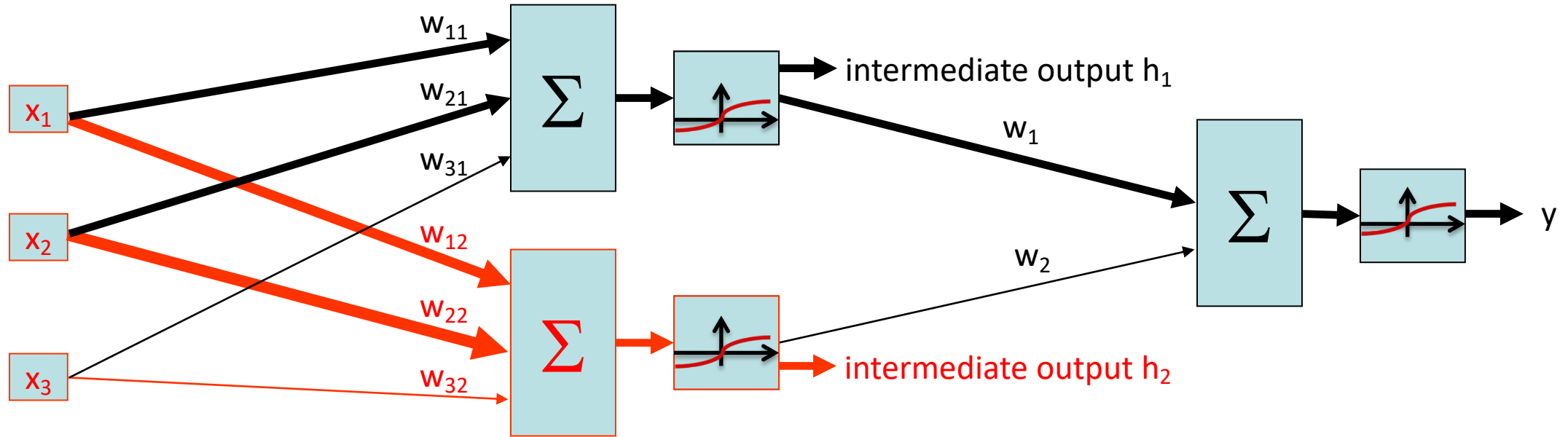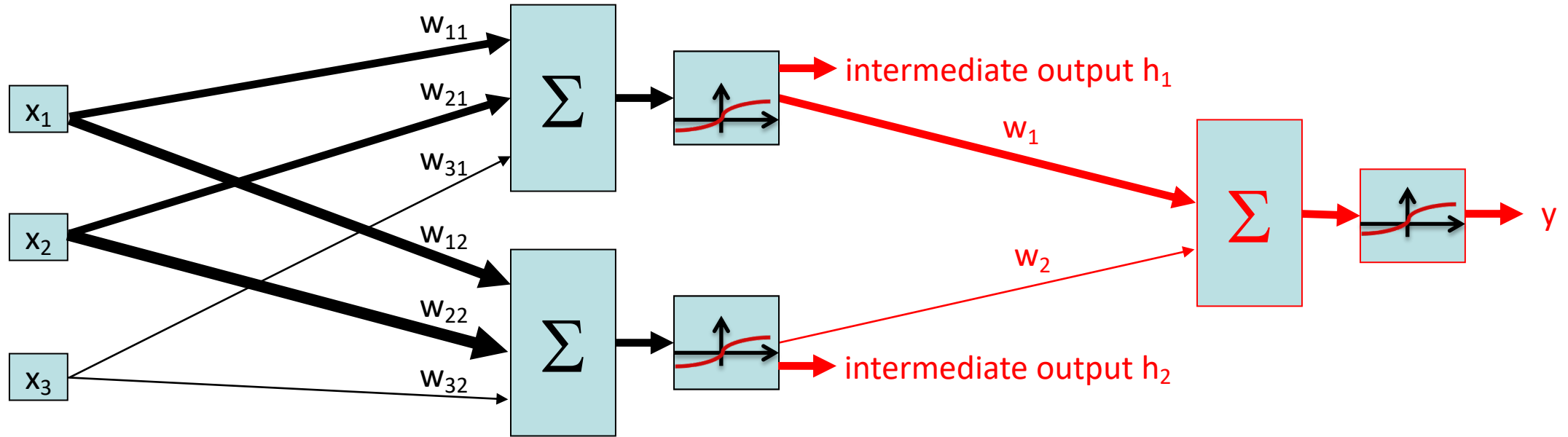# Recall: 2-Layer, 2-Neuron Neural Network



$$\text{intermediate output } h_1 = \phi(w_{11}x_1 + w_{21}x_2 + w_{31}x_3)$$

$$= \frac{1}{1 + e^{-(w_{11}x_1 + w_{21}x_2 + w_{31}x_3)}}$$

# Recall: 2-Layer, 2-Neuron Neural Network



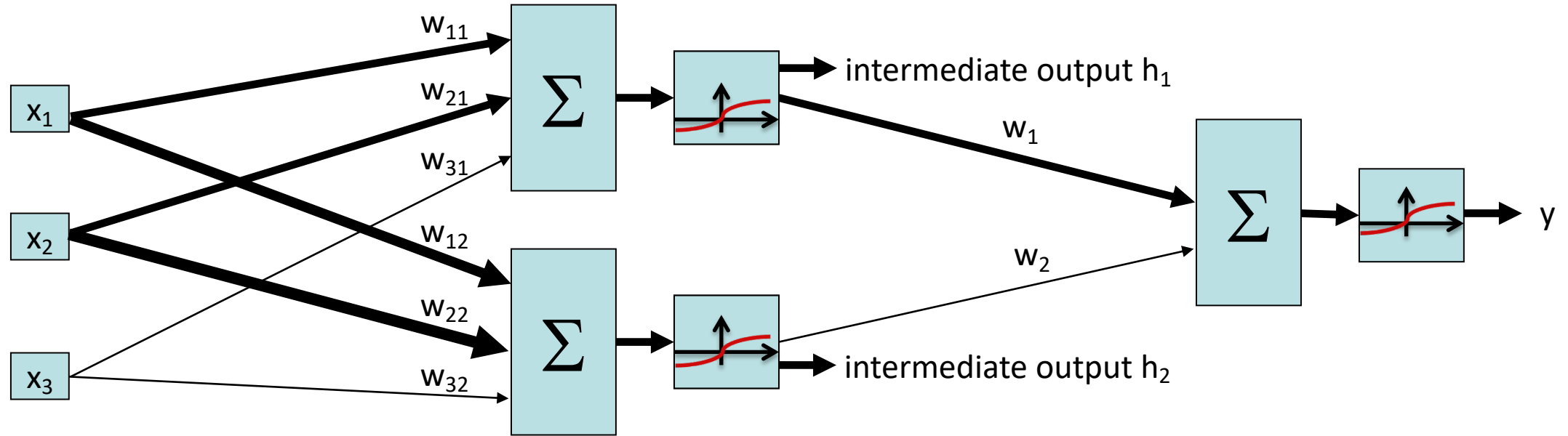$$\text{intermediate output } h_2 = \phi(w_{12}x_1 + w_{22}x_2 + w_{32}x_3)$$

$$= \frac{1}{1 + e^{-(w_{12}x_1 + w_{22}x_2 + w_{32}x_3)}}$$

# Recall: 2-Layer, 2-Neuron Neural Network



intermediate output h$_1$

intermediate output h$_2$
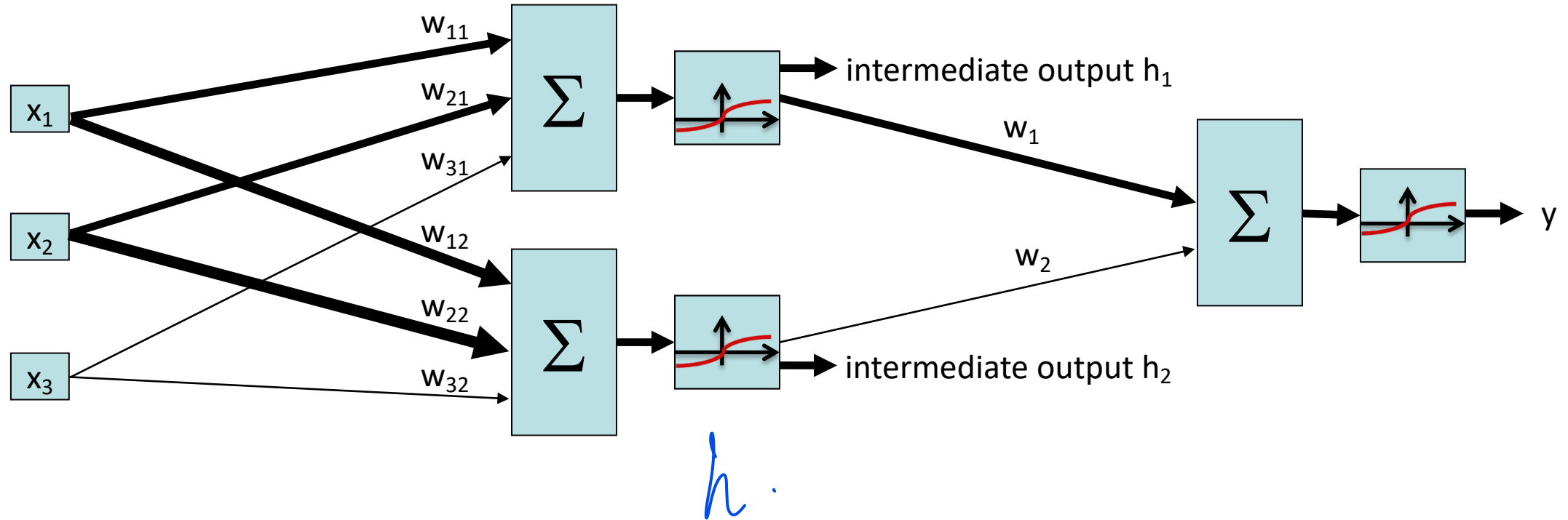
$$y = \phi(w_1 h_1 + w_2 h_2)$$
$$= \frac{1}{1 + e^{-(w_1 h_1 + w_2 h_2)}}$$

# Recall: 2-Layer, 2-Neuron Neural Network



$$y = \phi(w_1 h_1 + w_2 h_2)$$
$$= \phi(w_1 \phi(w_{11} x_1 + w_{21} x_2 + w_{31} x_3) + w_2 \phi(w_{12} x_1 + w_{22} x_2 + w_{32} x_3))$$
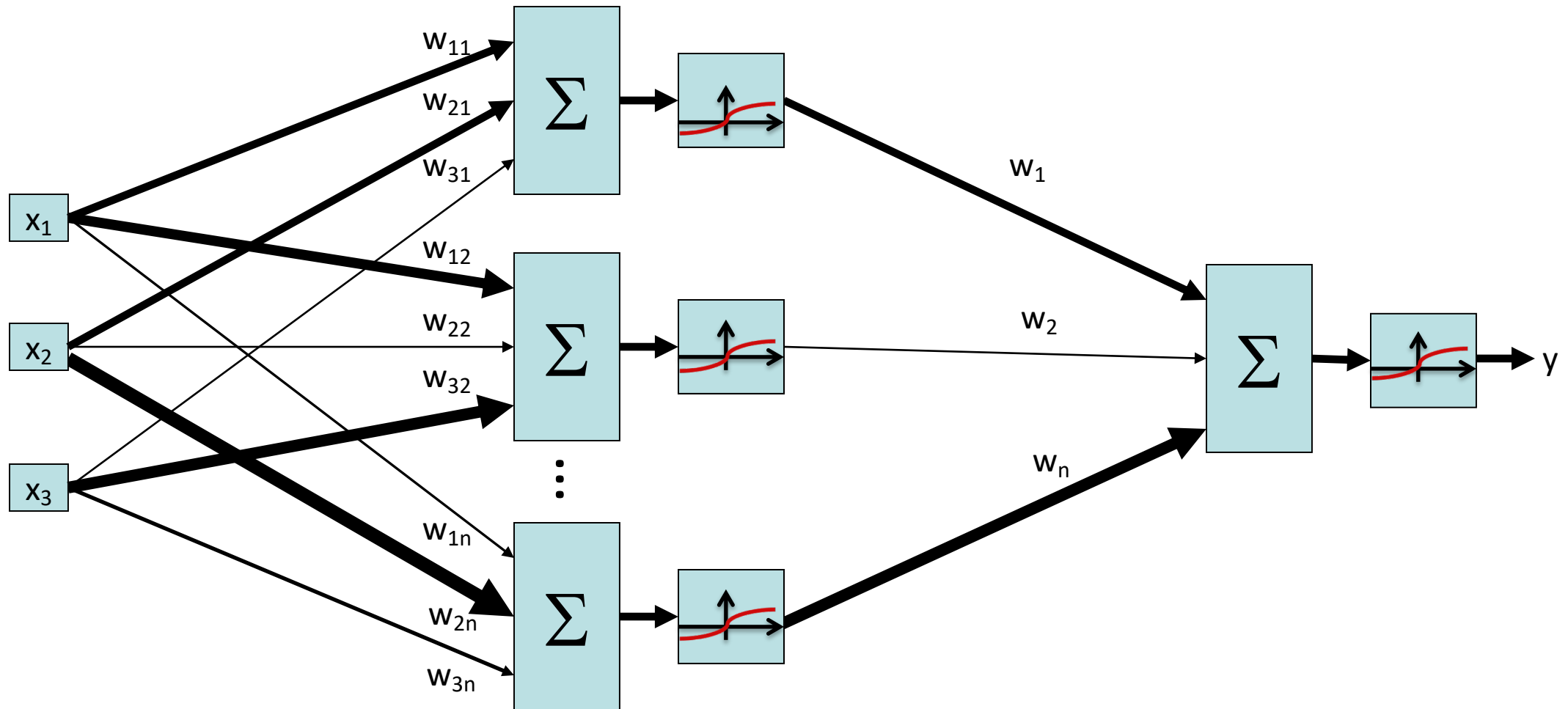
# Recall: 2-Layer, 2-Neuron Neural Network

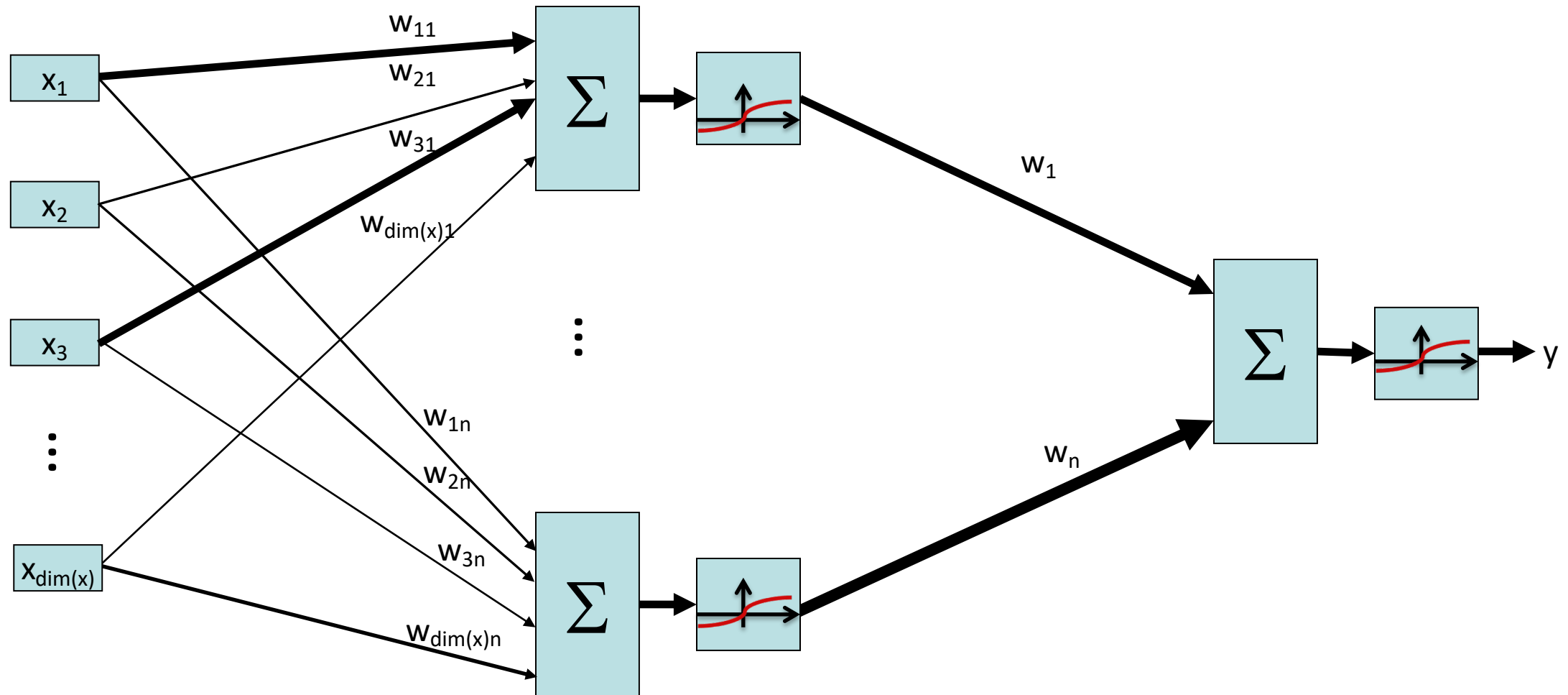

$$\phi(x \times W_{\text{layer 1}}) = h \qquad \phi(h \times W_{\text{layer 2}}) = y$$
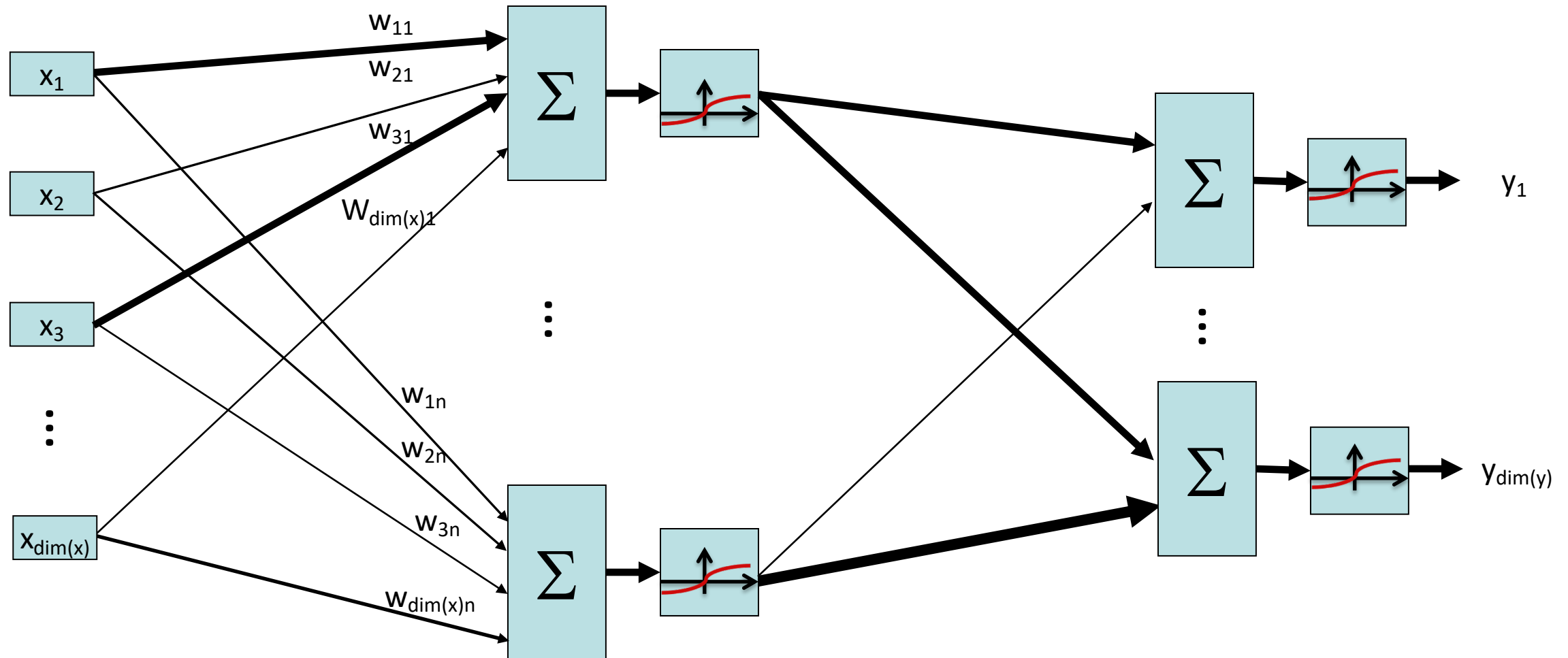
# Recall: generalize number of hidden neurons



The hidden layer doesn't necessarily need to have 3 neurons; it could have any arbitrary number *n* neurons.

# Recall: generalize number of input features
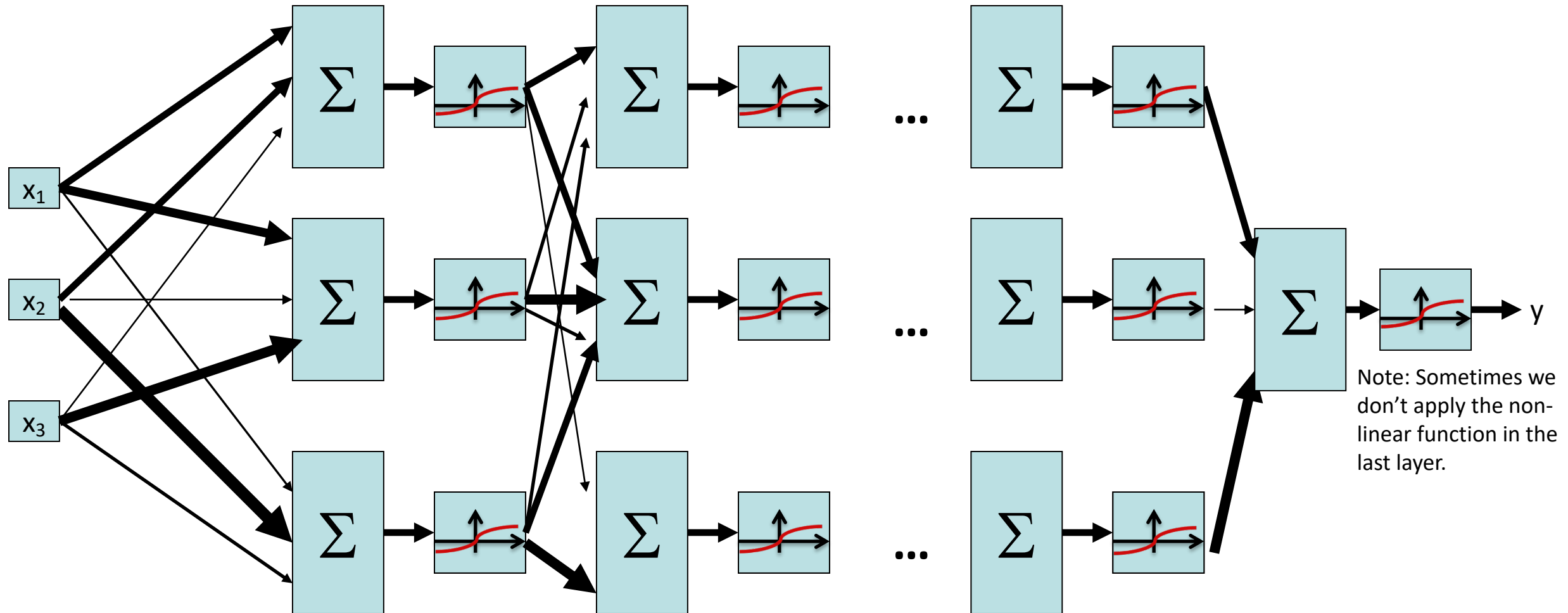


The input feature vector doesn't necessarily need to have 3 features; it could have some arbitrary number *dim(x)* of features.
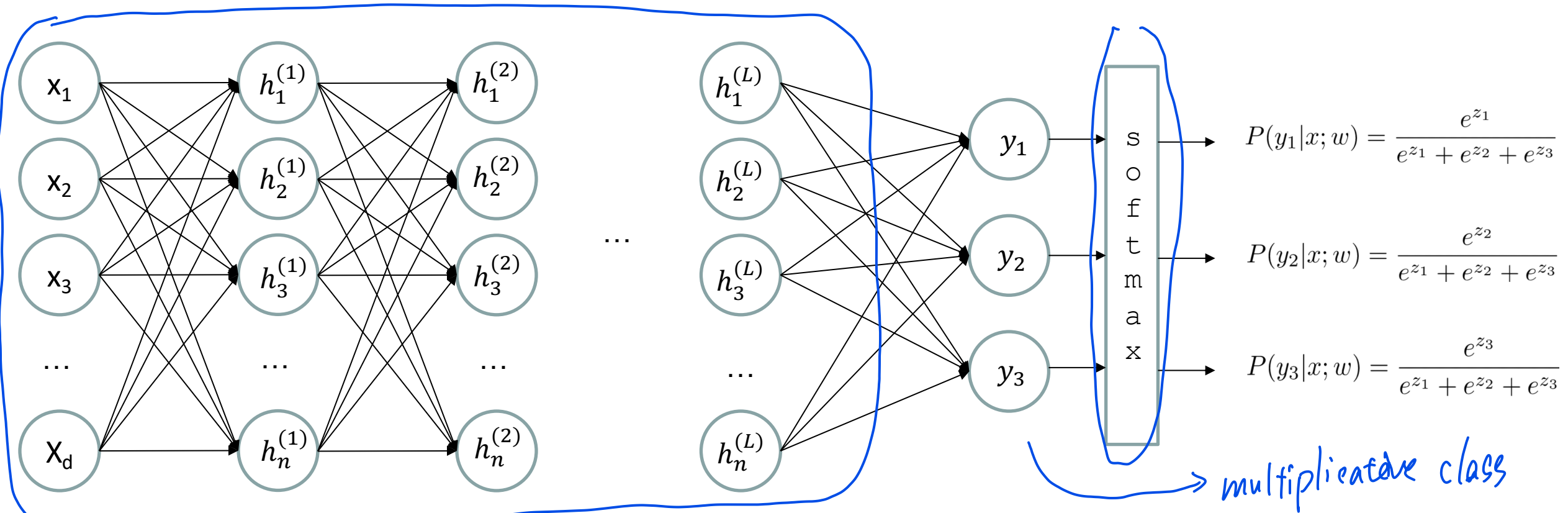
# Recall: generalize number of outputs



The output doesn't necessarily need to be just one number; it could be some arbitrary *dim(y)* length vector.

Note: Sometimes we don't apply the non-linear function in the last layer.

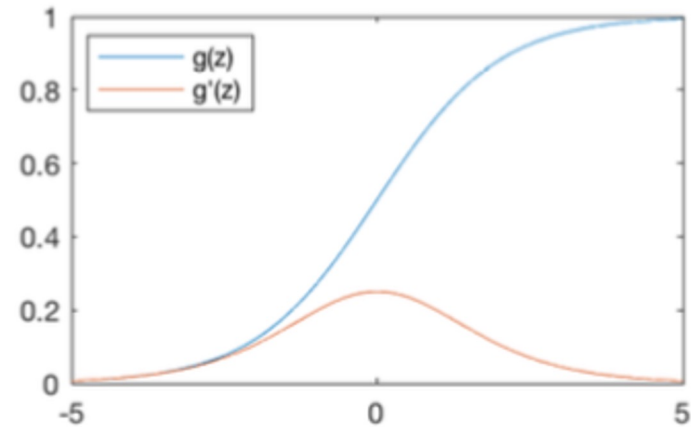# Deep Neural Network for 3-way classification



$$h_i^{(\text{layer } l)} = \phi\left(\sum_j w_{ji}^{(\text{layer } l)} \cdot h_j^{(\text{layer } l-1)}\right)$$

$\phi$ = nonlinear activation function

$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

→ multiplicative class

- Neural network with L layers
- $h^{(l)}$: activations at layer l
- $w^{(l)}$: weights taking activations from layer l-1 to layer l
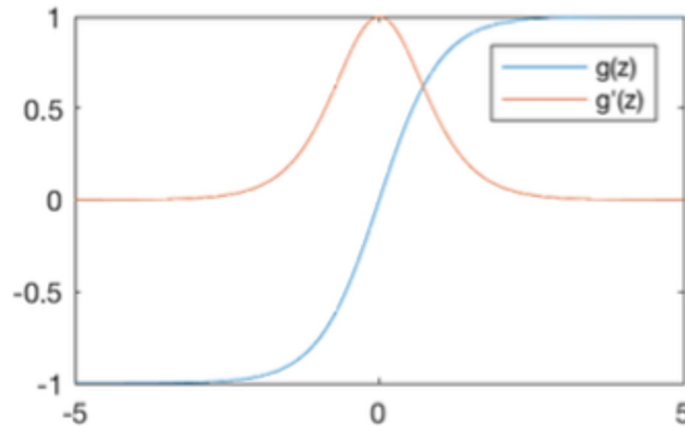
# Recall: Common Activation Functions

## Sigmoid Function



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

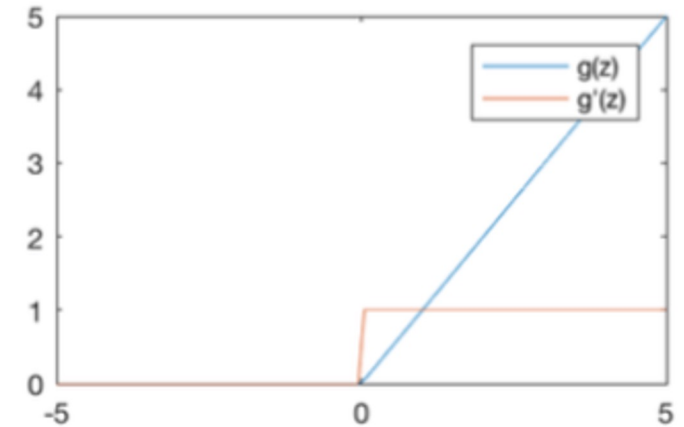$$g'(z) = g(z)(1 - g(z))$$

## Hyperbolic Tangent



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

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

## Rectified Linear Unit (ReLU)



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

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Recall: Sizes of neural networks



$\phi(\quad x \quad \times \quad W_{layer\ 1} \quad) = \quad h$

$\phi(\quad h \quad \times \quad W_{layer\ 2} \quad) = \quad y$

We have a neural network with the matrices drawn.

1. How many layers are in the network?
   2

2. How many input dimensions dim(x)?
   3

3. How many hidden neurons n?
   2

4. How many output dimensions dim(y)?
   1

5. What is the batch size?
   4

Training the deep neural network is just like logistic regression:

$$\max_{w} \quad ll(w) = \max_{w} \quad \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

just w tends to be a much, much larger vector

-> just run gradient ascent

+ stop when log likelihood of hold-out data starts to decrease

# Batch Gradient Ascent on the Log Likelihood Objective

$$\max_w \; ll(w) = \max_w \; \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

$$\underbrace{\phantom{\sum_i \log P(y^{(i)}|x^{(i)}; w)}}_{g(w)}$$

```
init w
for iter = 1, 2, …
```

$$w \leftarrow w + \underset{10^{-3}}{\alpha} * \sum_i \nabla \log P(y^{(i)}|x^{(i)}; w)$$

$$\nabla g(w)$$

# How about computing all the derivatives?

Derivatives tables:

$$\frac{d}{dx}(a) = 0$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(au) = a\frac{du}{dx}$$

$$\frac{d}{dx}(u+v-w) = \frac{du}{dx} + \frac{dv}{dx} - \frac{dw}{dx}$$

$$\frac{d}{dx}(uv) = u\frac{dv}{dx} + v\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{u}{v}\right) = \frac{1}{v}\frac{du}{dx} - \frac{u}{v^2}\frac{dv}{dx}$$

$$\frac{d}{dx}(u^n) = nu^{n-1}\frac{du}{dx}$$

$$\frac{d}{dx}(\sqrt{u}) = \frac{1}{2\sqrt{u}}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u}\right) = -\frac{1}{u^2}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u^n}\right) = -\frac{n}{u^{n+1}}\frac{du}{dx}$$

$$\frac{d}{dx}[f(u)] = \frac{d}{du}[f(u)]\frac{du}{dx}$$

$$\frac{d}{dx}[\ln u] = \frac{d}{dx}[\log_e u] = \frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}[\log_a u] = \log_a e\frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}e^u = e^u\frac{du}{dx}$$

$$\frac{d}{dx}a^u = a^u \ln a\frac{du}{dx}$$

$$\frac{d}{dx}(u^v) = vu^{v-1}\frac{du}{dx} + \ln u \; u^v\frac{dv}{dx}$$

$$\frac{d}{dx}\sin u = \cos u\frac{du}{dx}$$

$$\frac{d}{dx}\cos u = -\sin u\frac{du}{dx}$$

$$\frac{d}{dx}\tan u = \sec^2 u\frac{du}{dx}$$

$$\frac{d}{dx}\cot u = -\csc^2 u\frac{du}{dx}$$

$$\frac{d}{dx}\sec u = \sec u \tan u\frac{du}{dx}$$

$$\frac{d}{dx}\csc u = -\csc u \cot u\frac{du}{dx}$$

# How about computing all the derivatives?

- But neural net f is never one of those?

  - No problem: CHAIN RULE:

  If
  $$f(x) = g(h(x))$$

  Then
  $$f'(x) = g'(h(x))h'(x)$$

  **Derivatives can be computed by following well-defined procedures**

# Automatic Differentiation

Automatic differentiation software

    e.g. TensorFlow, PyTorch, Jax

    Only need to program the function g(x,y,w)

    Can automatically compute all derivatives w.r.t. all entries in w

    This is typically done by caching info during forward computation pass of f, and then doing a backward pass = "backpropagation"

    Autodiff / Backpropagation can often be done at computational cost comparable to the forward pass

Need to know this exists

How this is done? Details outside of scope of CS188, but we'll show a basic example

# Backpropagation*

- Gradient of $g(w_1, w_2, w_3) = w_1^4 w_2 + 5w_3$ at $w_1 = 2$, $w_2 = 3$, $w_3 = 2$
- Think of $g$ as a composition of many functions
  - Then, we can use the chain rule to compute the gradient
- $g = b + c$

$$\frac{\partial g}{\partial b} = 1, \frac{\partial g}{\partial c} = 1$$

- $b = a \times w_2$

$$\frac{\partial g}{\partial a} = \frac{\partial g}{\partial b}\frac{\partial b}{\partial a} = 1 \cdot w_2 = 3 \qquad \frac{\partial g}{\partial w_2} = \frac{\partial g}{\partial b}\frac{\partial b}{\partial w_2} = 1 \cdot a = 16$$

- $a = w_1^4$

$$\frac{\partial g}{\partial w_1} = \frac{\partial g}{\partial a}\frac{\partial a}{\partial w_1} = 3 \cdot 4w_1^3 = 96$$

- $c = 5w_3$

$$\frac{\partial g}{\partial w_3} = \frac{\partial g}{\partial c}\frac{\partial c}{\partial w_3} = 1 \cdot 5 = 5$$

# Properties of Neural Networks

# Neural Networks Properties

- Theorem (Universal Function Approximators).  A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

# Universal Function Approximation Theorem*

**Hornik theorem 1:** Whenever the activation function is *bounded and nonconstant*, then, for any finite measure $\mu$, standard multilayer feedforward networks can approximate any function in $L^p(\mu)$ (the space of all functions on $R^k$ such that $\int_{R^k} |f(x)|^p d\mu(x) < \infty$) arbitrarily well, provided that sufficiently many hidden units are available.

**Hornik theorem 2:** Whenever the activation function is *continuous, bounded and nonconstant*, then, for arbitrary compact subsets $X \subseteq R^k$, standard multilayer feedforward networks can approximate any continuous function on $X$ arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

- <u>In words:</u> Given any continuous function f(x), if a 2-layer neural network has enough hidden units, then there is a choice of weights that allow it to closely approximate f(x).

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

# Universal Function Approximation Theorem*

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation
Functions Can Approximate Any Function"

# Neural Networks Properties

- **Theorem (Universal Function Approximators).** A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

- Practical considerations
  - Can be seen as learning the features

  - Large number of neurons
    - Danger for overfitting
    - (hence early stopping!)

# Preventing Overfitting in Neural Networks

Early stopping:



Weight regularization

# Weight Regularization

What can go wrong when we maximize log-likelihood?

Example: logistic regression with only one datapoint: f(x)=1, y=+1

$$\max_{w} \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

- $P(y = +1|x; w) = \frac{1}{1+e^{-w \cdot f(x)}}$

$$\max_{w} \frac{1}{1 + e^{-w}}$$

$$\log \frac{1}{1+e^{-w}} = \log 1 - \log(1+e^{-w})$$

$$= 0 - \log(1+e^{-w})$$

$w$ can grow very large and lead to overfitting and learning instability

# Weight Regularization

What can go wrong when we maximize log-likelihood?

$$\max_{w} \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

$w$ can grow very large

Solution: add an objective term to penalize weight magnitude

$$\max_{w} \sum_{i} \log P(y^{(i)}|x^{(i)}; w) - \frac{\lambda}{2} \sum_{j} w_j^2$$

$\lambda$ is a hyperparameter (typically 0.1 to 0.0001 or smaller)

# Preventing Overfitting in Neural Networks

Early stopping:



Weight regularization: $\max_w \sum_i \log P(y^{(i)} | x^{(i)}; w) - \frac{\lambda}{2} \sum_j w_j^2$

Dropout

# Consistency vs. Simplicity

- Example: curve fitting (regression, function approximation)



- Consistency vs. simplicity
- Ockham's razor

# Consistency vs. Simplicity

- Usually algorithms prefer consistency by default (why?)

- Several ways to operationalize "simplicity"
  - Reduce the hypothesis/model space
    - Assume more: e.g. independence assumptions, as in naïve Bayes
    - Fewer features or neurons
    - Other limits on model structure
  - Regularization
    - Laplace Smoothing: cautious use of small counts
    - Small weight vectors in neural networks (stay close to zero-mean prior)
    - Hypothesis space stays big, but harder to get to the outskirts

# Fun Neural Net Demo Site

Demo-site:

http://playground.tensorflow.org/

# Summary of Key Ideas

Optimize probability of label given input $\qquad \max_w \; ll(w) = \max_w \; \sum_i \log P(y^{(i)}|x^{(i)}; w)$

## Continuous optimization

Gradient ascent:

Compute steepest uphill direction = gradient (= just vector of partial derivatives)

Take step in the gradient direction

Repeat (until held-out data accuracy starts to drop = "early stopping")

## Deep neural nets

Last layer = still logistic regression

Now also many more layers before this last layer

= computing the features

the features are learned rather than hand-designed

Universal function approximation theorem

`If`      neural net is large enough

`Then`   neural net can represent any continuous mapping from input to output with arbitrary accuracy

But remember: need to avoid overfitting / memorizing the training data ? early stopping!

Automatic differentiation gives the derivatives efficiently (how? = outside of scope of 188)

# Next: How well does deep learning work?