
Neural Networks

Zhiyao Duan

Associate Professor of ECE and CS

University of Rochester

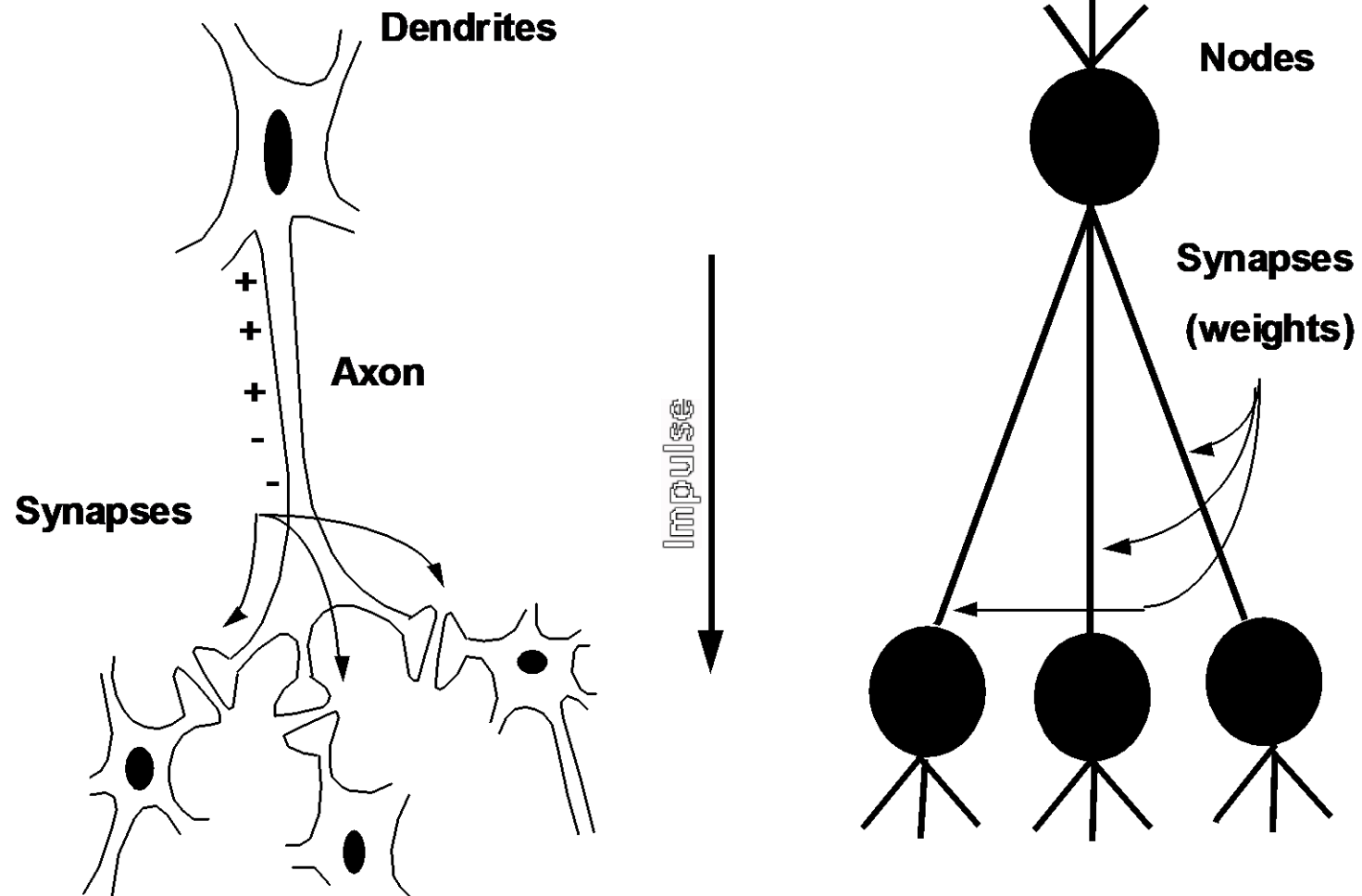
Some figures are copied from the following books

- **LWLS** - Andreas Lindholm, Niklas Wahlström, Fredrik Lindsten, Thomas B. Schön, *Machine Learning: A First Course for Engineers and Scientists*, Cambridge University Press, 2022.

Biological Motivation

- Human brain: a densely interconnected network
 - $\sim 10^{11}$ neurons
 - Each neuron connects to $\sim 10^4$ other neurons
 - Two states of neuron activity: excited vs. inhibited
 - Neuron switching speed: $\sim 1\text{kHz}$
 - CPU clock frequency: GHz
 - Yet many tasks (e.g., face recognition) can be completed within 0.1 s
- This suggests
 - Highly parallel processing
 - Distributed representations

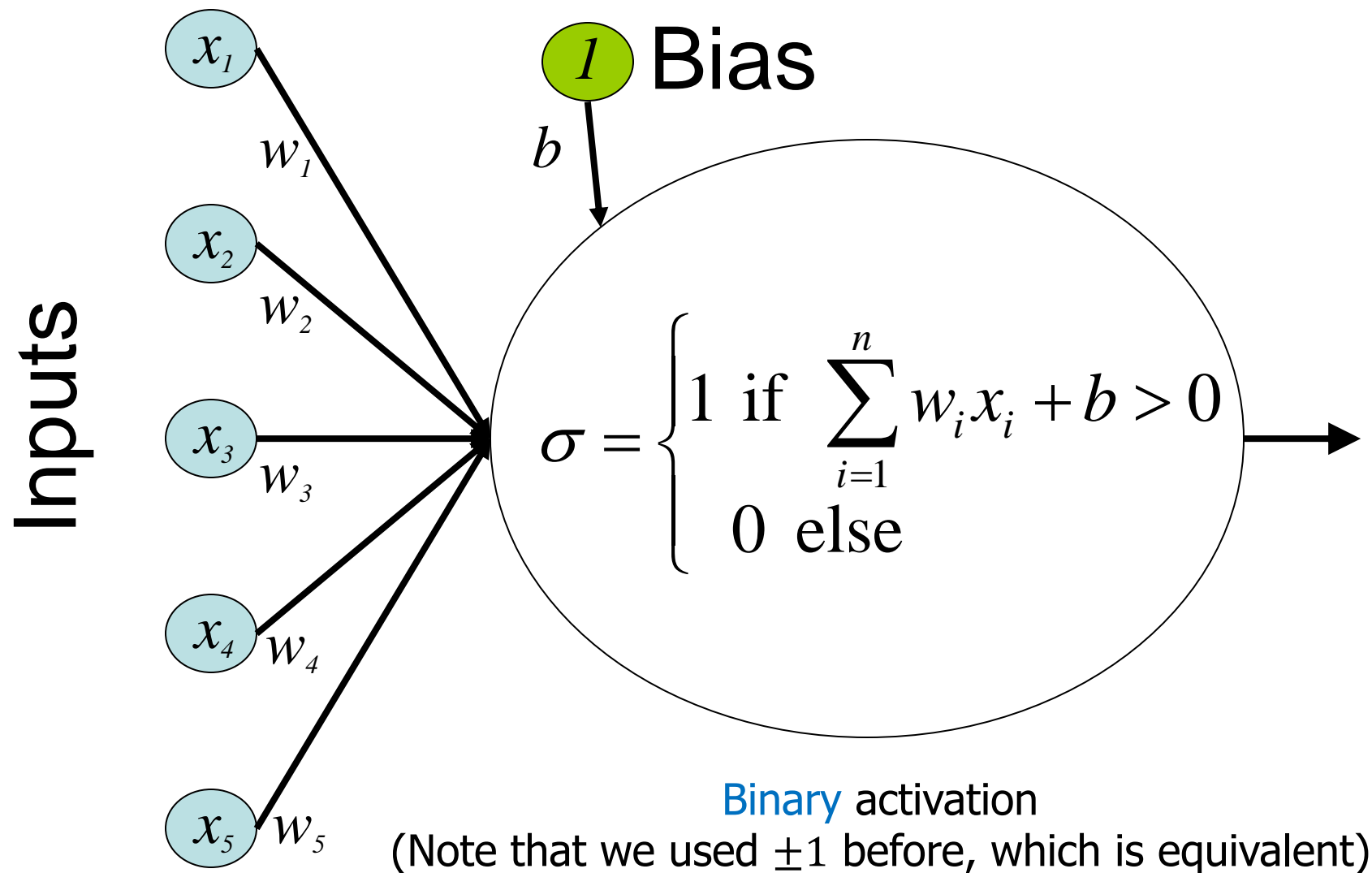
Biological Analogy



History of Neural Networks

- 1943 – first neural network computing model by McCulloch and Pitts
- 1958 – Perceptron by Rosenblatt
- 1960's – a big wave
- 1969 – Minsky & Papert's book "Perceptrons"
- 1970's – "winter" of neural networks
- 1975 – Backpropagation algorithm by Werbos
- 1980's – another big wave
- 1990's – overtaken by SVM proposed in 1993 by Vapnik
- 2006 – a fast learning algorithm for training deep belief networks by Hinton
- 2010's – another big wave
- 2018 – Turing Award to Hinton, Bengio & LeCun
- 2022 – ChatGPT!
- 2024 – Sora!
- Present – continue to transform various domains

Perceptron

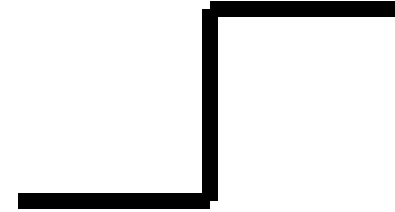


Nonlinear Activation Functions

- Step function

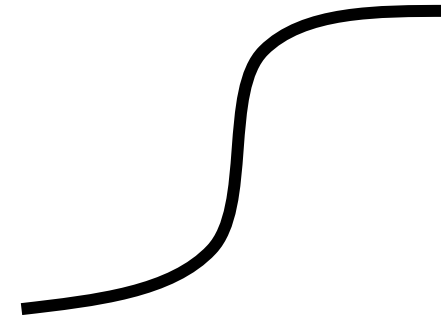
$$output = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

- Note: previously we used $\{-1,1\}$ for sign function for perceptron, which is equivalent



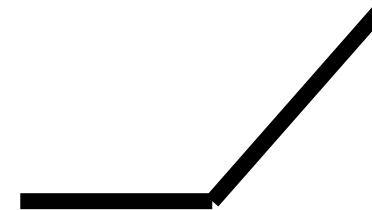
- Sigmoid function

$$output = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$



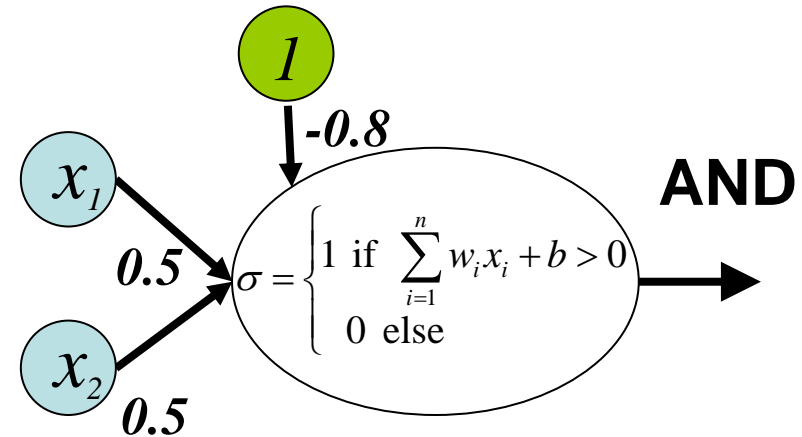
- Rectified Linear Unit (ReLU)

$$output = \max\{0, \mathbf{w}^T \mathbf{x} + b\}$$

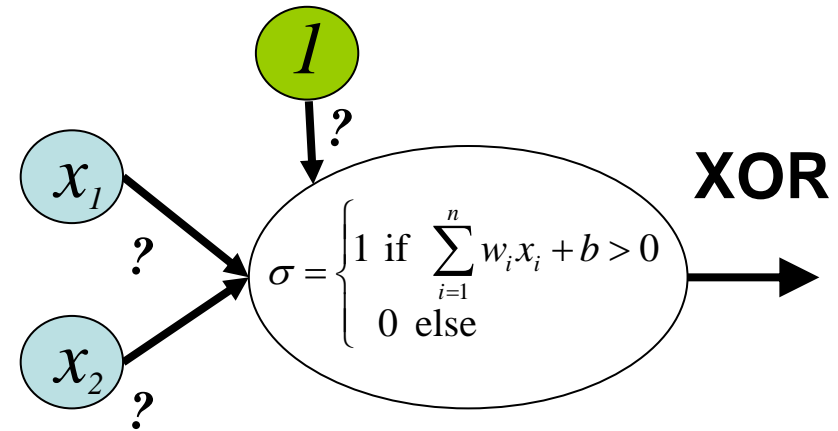


Limitations of 1-layer Nets

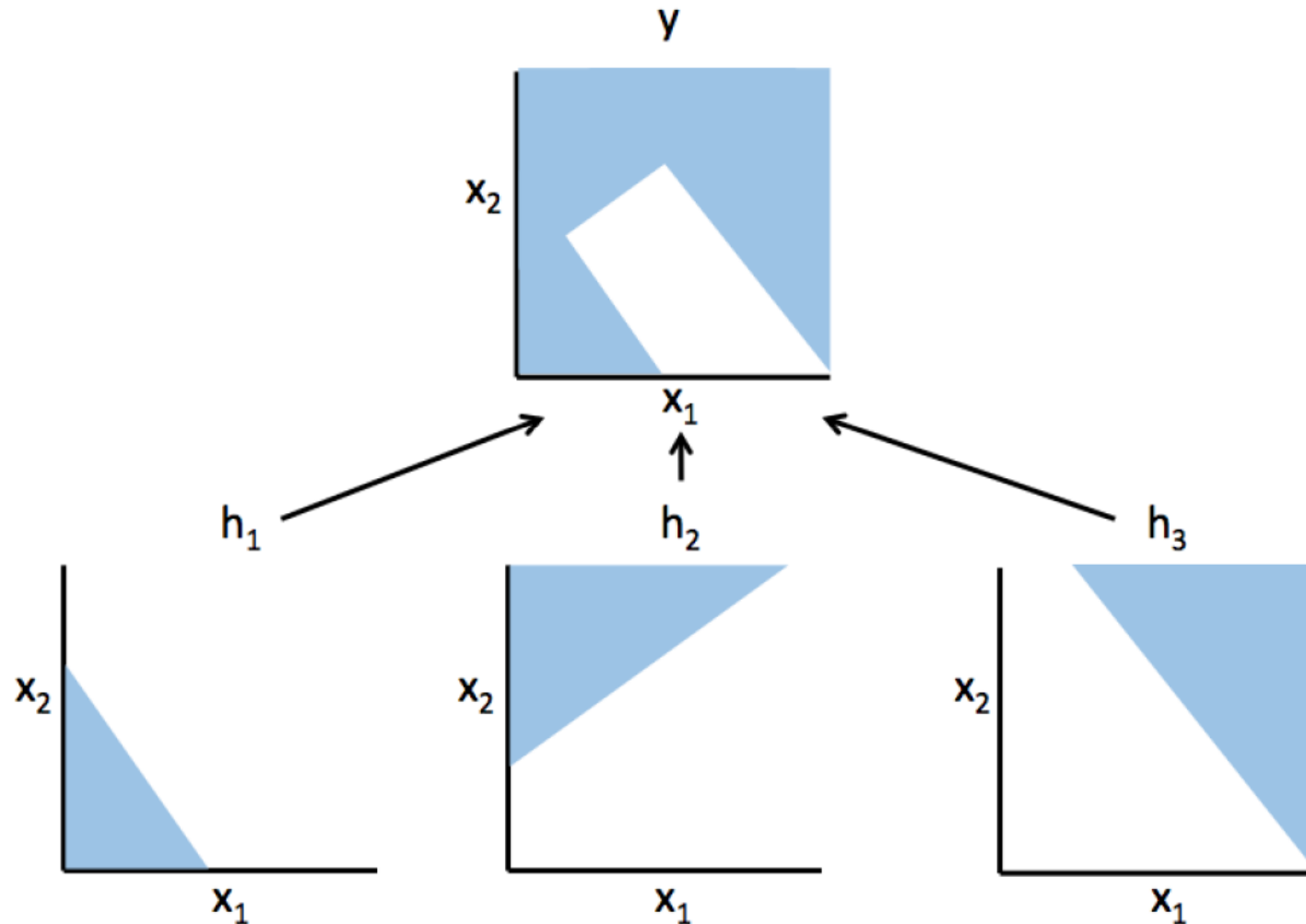
- Only express linearly separable cases
 - For example, they are good as logic operators “AND”, “NOT”, and “OR”



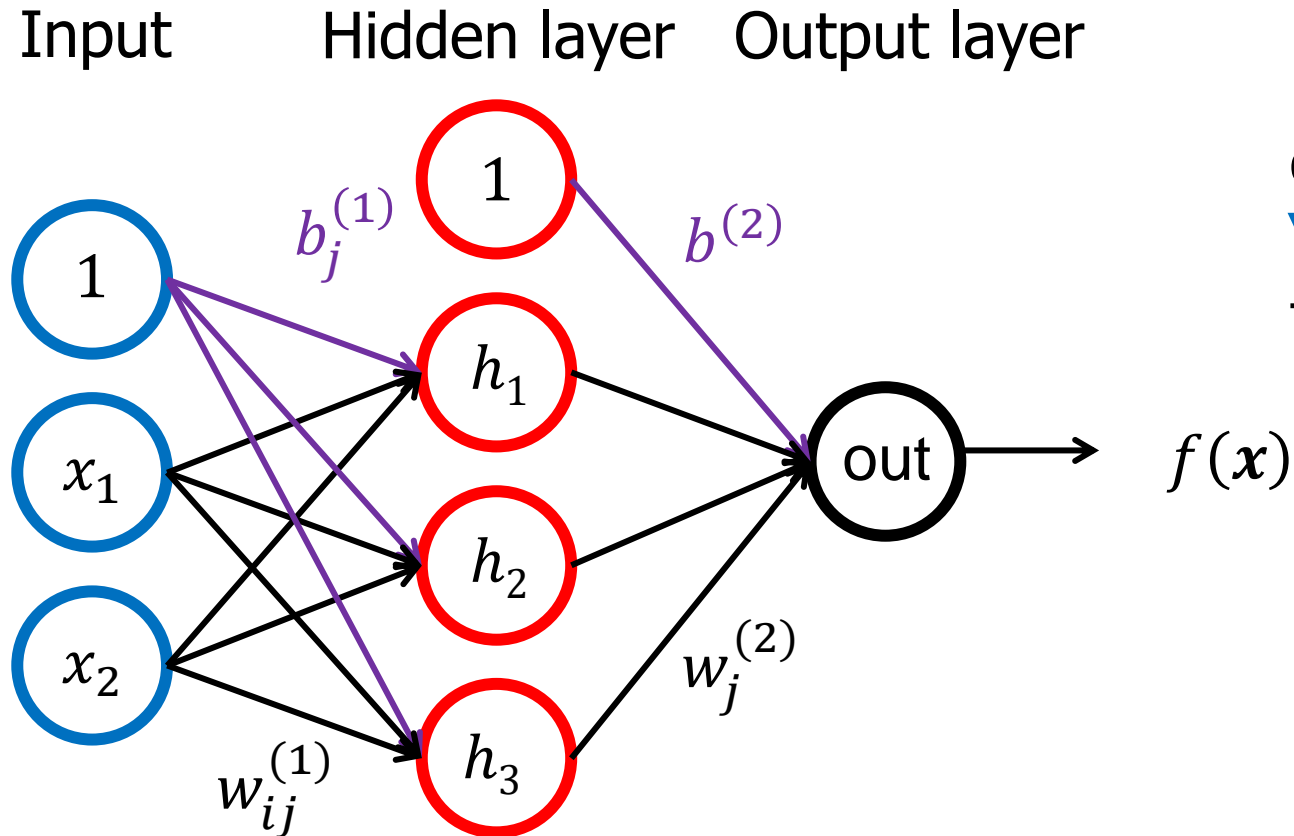
- Cannot represent “XOR”, which is not linearly separable



But, we can combine them!



2-layer Nets



$$f(\mathbf{x}) = \sigma \left(\sum_j w_j^{(2)} h_j + b^{(2)} \right) = \sigma \left(\sum_j w_j^{(2)} \sigma \left(\sum_i w_{ij}^{(1)} x_i + b_j^{(1)} \right) + b^{(2)} \right)$$

Matrix Notation

$$f(\mathbf{x}) = \sigma \left(\sum_j w_j^{(2)} \sigma \left(\sum_i w_{ij}^{(1)} x_i + b_j^{(1)} \right) + b^{(2)} \right)$$

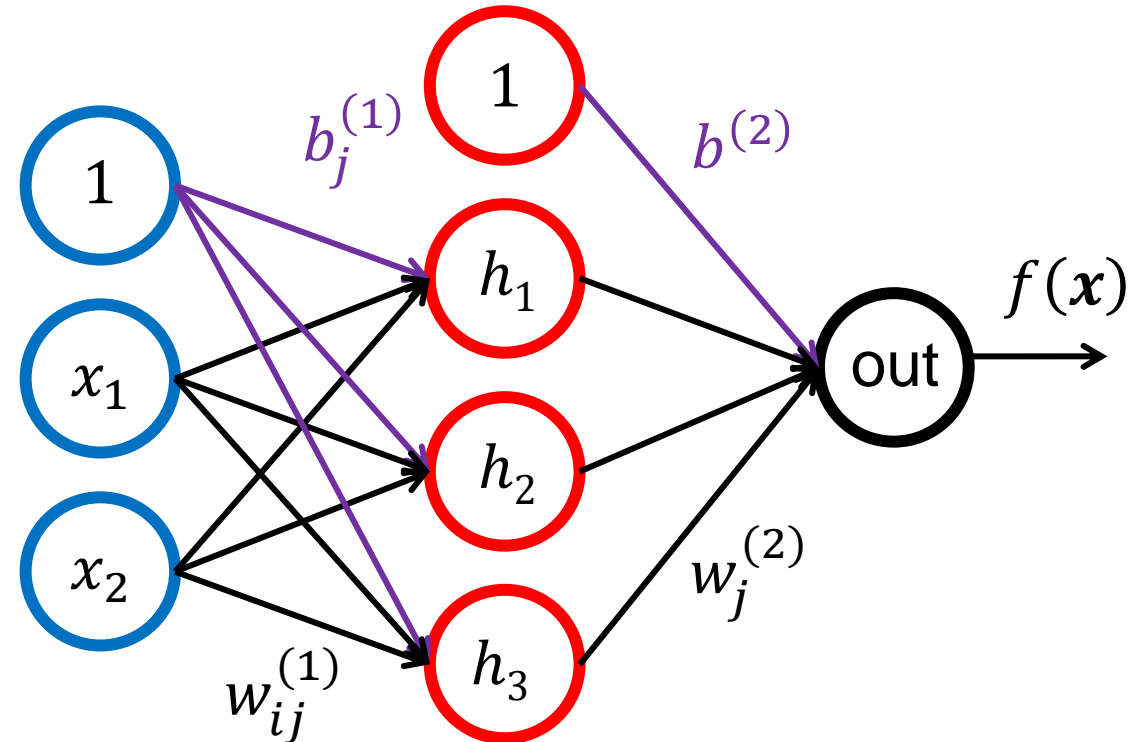
Input Hidden layer Output layer

$$f(\mathbf{x}) = \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1) + b_2)$$

where

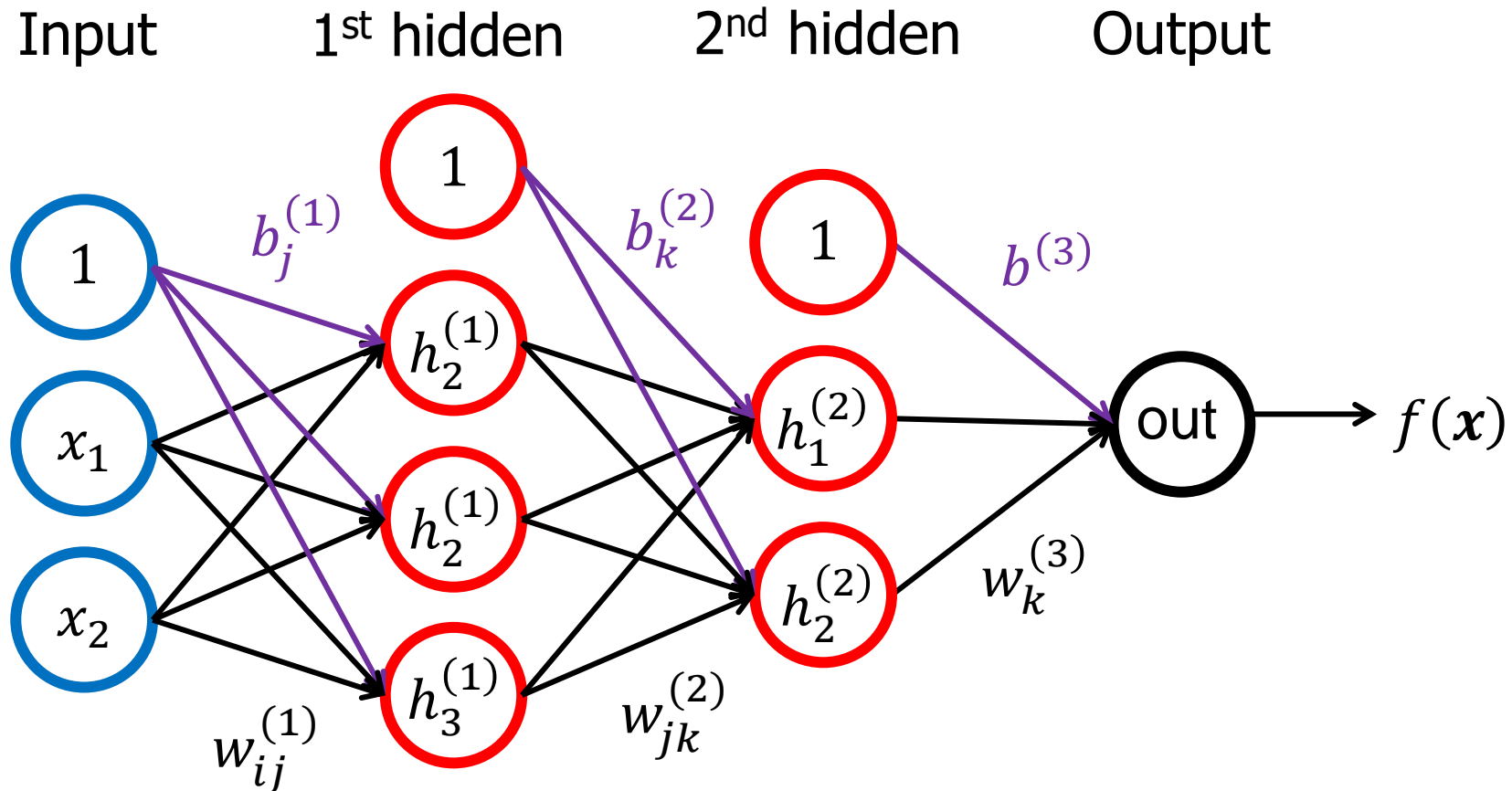
$$\mathbf{W}_1 = [w_{ij}^{(1)}]_{d \times l_1}, \mathbf{b}_1 = [b_j^{(1)}]_{l_1 \times 1}$$

$$\mathbf{W}_2 = [w_{jk}^{(2)}]_{l_1 \times l_2}, b_2 = b^{(2)}$$



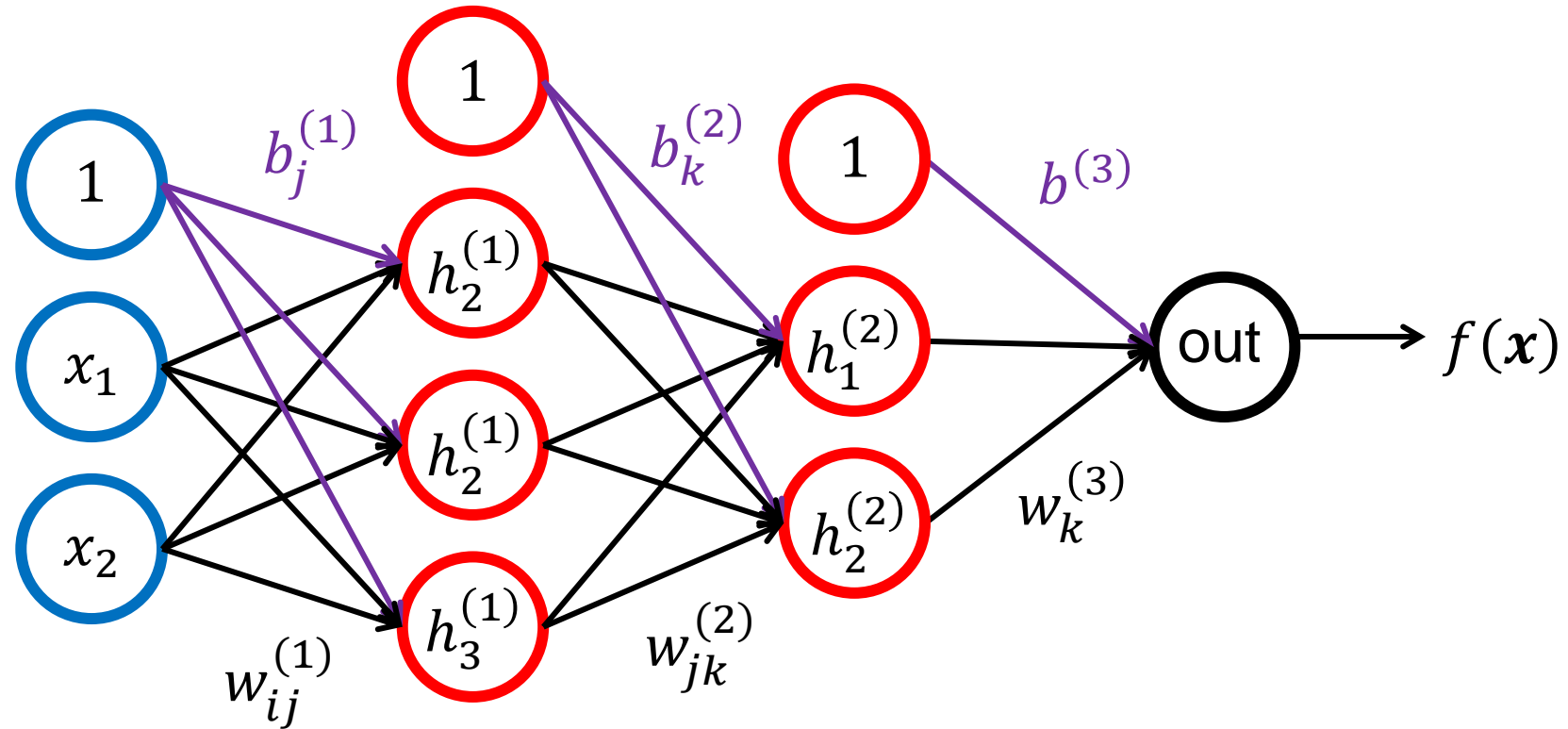
- What does $\mathbf{W}_1^T \mathbf{x}$ compute?
 - Inner products between columns of \mathbf{W}_1 and \mathbf{x}
 - Columns of \mathbf{W}_1 are "receptors" or "filters"
 - $\mathbf{W}_1^T \mathbf{x}$ are their **responses** to input

3-layer Nets



$$f(\mathbf{x}) = \sigma \left(\sum_k w_k^{(3)} h_k^{(2)} + b^{(3)} \right) = \sigma \left(\sum_k w_k^{(3)} \sigma \left(\sum_j w_{jk}^{(2)} h_j^{(1)} + b_k^{(2)} \right) + b^{(3)} \right) = \sigma \left(\sum_k w_k^{(3)} \sigma \left(\sum_j w_{jk}^{(2)} \sigma \left(\sum_i w_{ij}^{(1)} x_i + b_j^{(1)} \right) + b_k^{(2)} \right) + b^{(3)} \right)$$

Matrix Notation



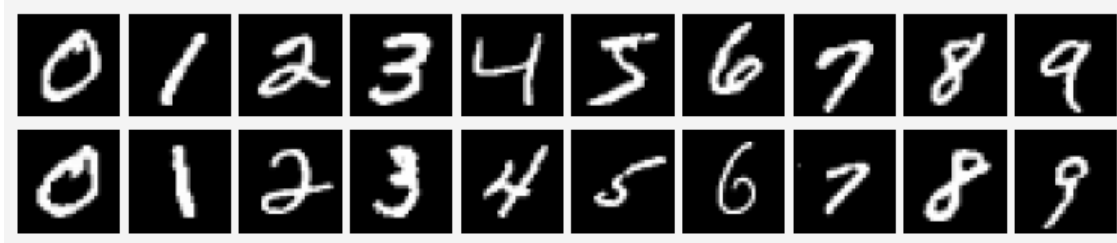
$$f(\mathbf{x}) = \sigma \left(\sum_k w_k^{(3)} \sigma \left(\sum_j w_{jk}^{(2)} \sigma \left(\sum_i w_{ij}^{(1)} x_i + b_j^{(1)} \right) + b_k^{(2)} \right) + b^{(3)} \right)$$

$$f(\mathbf{x}) = \sigma(\mathbf{W}_3^T \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) + b_3)$$

Richer Representations with More Layers

- 1-layer nets (e.g., perceptron) only model linear hyperplanes
- 2-layer nets can approximate any continuous function, given **enough** hidden nodes
- ≥ 3 -layer nets can do so with fewer nodes and weights
- Nonlinear activation is key!
 - Multiple layers of linear activations is still linear!

Example Application



(Fig. 6.5 in LWLS, from MNIST dataset)
70,000 grayscale images (28*28) from 10 classes

- One-layer MLP (i.e., logistic regression)
 - Input: $28*28=784$ -d vectors
 - Output layer size: 10 nodes
 - #parameters: $784*10+10 = 7,850$
- Two-layer MLP
 - Input: $28*28=784$ -d vectors
 - Hidden layer size: 200 nodes
 - Output layer size: 10 nodes
 - #parameters for hidden layer: $784*200+200$
 - #parameters for output layer: $200*10+10$
 - #Total parameters = 159,010

Properties of NNs

- Large capacity: able to learn **complex** relations between input and output
- Support various data formats: continuous, discrete, categorical (needs to be encoded into numeric)
- Robust to some level of noise in training data
- Inference (i.e., making predictions on test examples) is fast

- Data hungry
- Training is slow
- Lack of mathematical analysis and difficult to interpret

How to learn the weights?

- Given training data - input and label pairs $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$
- Update network weights to minimize the difference (error) between $f(\mathbf{x}^{(i)})$ and $y^{(i)}$
 - Calculate derivative of error w.r.t. weights
 - Gradient descent to update weights
 - **Backpropagation** algorithm: recursive computation of these gradients
- See derivation on white board

Backpropagation Recap

- Assume we use **sigmoid activation** and the **squared error loss**
 - We can also use other activations, e.g., ReLU
 - We can also use other losses, e.g., cross entropy
- Then the loss on the **entire** training set is

$$E(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2 = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta}))^2$$

where $\boldsymbol{\theta}$ denotes network parameters, i.e., network weights

- We compute gradient $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$ (called the **true gradient**, versus **stochastic gradient** computed on a subset of data), and then update $\boldsymbol{\theta}$ along the negative gradient direction iteratively
- The computation of $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$ is recursive, **backward** from the last layer to the first layer, leveraging the layer-wise structure of the network
- The computation also requires node outputs at each layer, which are computed in a **forward** pass

Forward Pass In Matrix Notation


- Start from **input** $\mathbf{X}_{N \times d} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]^T$ corresponding to all N points
- Compute first hidden layer **net input** \mathbf{Z}_1
$$[\mathbf{Z}_1]_{N \times l_1} = [\mathbf{X}\mathbf{W}_1]_{N \times l_1} + [\text{ repmat}(\mathbf{b}_1^T)]_{N \times l_1}$$
- Compute first hidden layer **output** \mathbf{H}_1
$$[\mathbf{H}_1]_{N \times l_1} = \sigma(\mathbf{Z}_1)$$
- Compute second hidden layer **net input** \mathbf{Z}_2
$$[\mathbf{Z}_2]_{N \times l_2} = [\mathbf{H}_1\mathbf{W}_2]_{N \times l_2} + [\text{ repmat}(\mathbf{b}_2^T)]_{N \times l_2}$$
- Compute second hidden layer **output** \mathbf{H}_2
$$[\mathbf{H}_2]_{N \times l_2} = \sigma(\mathbf{Z}_2)$$
-
- Compute final **output** $\hat{\mathbf{y}}$, a vector corresponding to all N points

Backward Pass in Matrix Notation


- Mean squared error computed on all data: $E(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2 = \frac{1}{2N} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})$
- Compute gradients w.r.t. weights in the output layer (the M -th layer)

$$\begin{aligned} \left[\frac{\partial E}{\partial \hat{\mathbf{y}}} \right]_{N \times 1} &= \frac{1}{N} (\hat{\mathbf{y}} - \mathbf{y}) \\ [\sigma'(\mathbf{z}_M)]_{N \times 1} &= \hat{\mathbf{y}} \odot (1 - \hat{\mathbf{y}}) \end{aligned}$$

$$\left[\frac{\partial E}{\partial \mathbf{w}_M} \right]_{l_{M-1} \times 1} = \left[\frac{\partial \mathbf{z}_M}{\partial \mathbf{w}_M} \right]_{l_{M-1} \times N} \cdot \left[\left[\frac{\partial E}{\partial \hat{\mathbf{y}}} \right]_{N \times 1} \odot [\sigma'(\mathbf{z}_M)]_{N \times 1} \right]$$

\mathbf{H}_{M-1}^T 

$$\frac{\partial E}{\partial b_M} = \left[\frac{\partial \mathbf{z}_M}{\partial b_M} \right]_{1 \times N} \cdot \left[\left[\frac{\partial E}{\partial \hat{\mathbf{y}}} \right]_{N \times 1} \odot [\sigma'(\mathbf{z}_M)]_{N \times 1} \right]$$

$\mathbf{1}^T$ 

Backward Pass in Matrix Notation

- Compute gradients w.r.t. weights in the $(m - 1)$ -th layer **recursively**

$$\left[\frac{\partial E}{\partial \mathbf{H}_{m-1}} \right]_{N \times l_{m-1}} = \left[\frac{\partial E}{\partial \mathbf{H}_m} \right]_{N \times l_m} \odot [\sigma'(\mathbf{Z}_m)]_{N \times l_m} \cdot [\mathbf{W}_m^T]_{l_m \times l_{m-1}}$$

$$\left[\frac{\partial E}{\partial \mathbf{W}_{m-1}} \right]_{l_{m-2} \times l_{m-1}} = [\mathbf{H}_{m-2}^T]_{l_{m-2} \times N} \cdot \left[\left[\frac{\partial E}{\partial \mathbf{H}_{m-1}} \right]_{N \times l_{m-1}} \odot [\sigma'(\mathbf{Z}_{m-1})]_{N \times l_{m-1}} \right]$$

$$\left[\frac{\partial E}{\partial \mathbf{b}_{m-1}} \right]_{l_{m-1} \times 1} = \left[\left[\frac{\partial E}{\partial \mathbf{H}_{m-1}} \right]_{N \times l_{m-1}} \odot [\sigma'(\mathbf{Z}_{m-1})]_{N \times l_{m-1}} \right]^T \cdot \mathbf{1}_{N \times 1}$$

Problems of BP for Deep Networks

- Vanishing gradient problem
 - Gradients **vanish** when they are propagated back to early layers, hence their weights are hard to adjust
 - Sigmoid activation → ReLU activation
- Many local minima
 - Which will trap gradient decent methods
 - In practice, local minima are pretty good

Summary

- (Artificial) neural networks are inspired by biological neural networks
 - Parallel processing + distributed representation
- Feedforward neural networks use a **layer-wise structure**
 - Full connection between adjacent layers
 - Linear mapping + nonlinear activation
- Representation power
 - 1-layer NNs are just perceptron or logistic regression
 - 2-layer NNs can represent (almost) any continuous function, with sufficient hidden nodes
 - ≥ 3 -layer NNs can do so with much fewer nodes
- **Gradient descent** to update network weights using training data
- **Backpropagation** algorithm to recursively compute gradients
 - **Vanishing gradient** issues for sigmoid activation