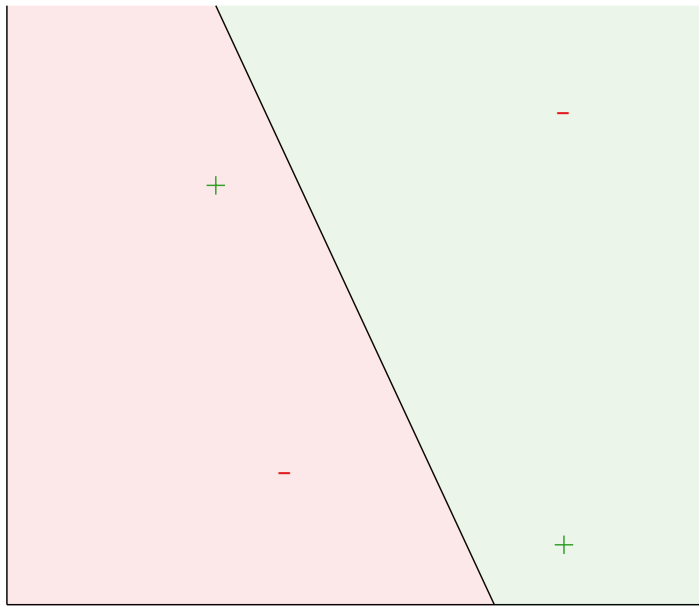


CSE 417T: Introduction to Machine Learning

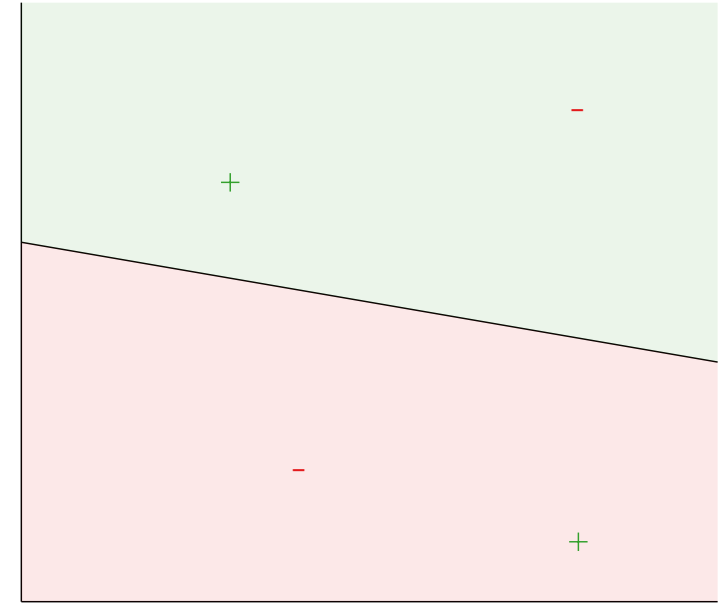
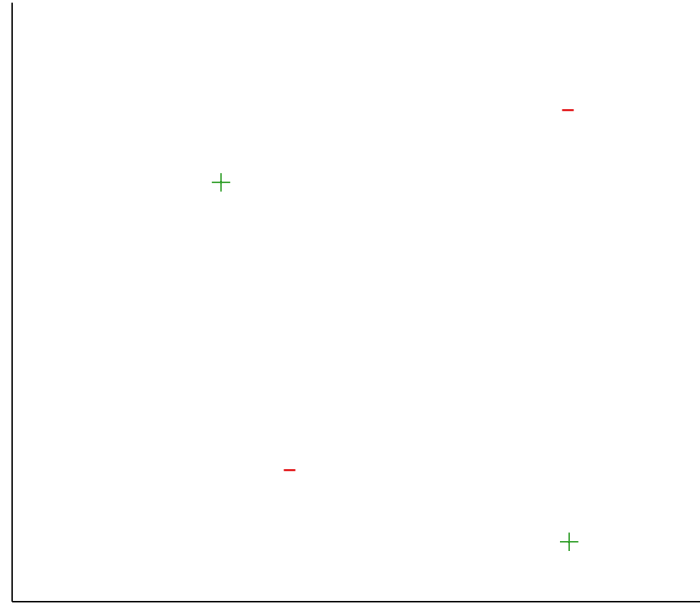
Lecture 25: Backpropagation

Henry Chai

11/29/18

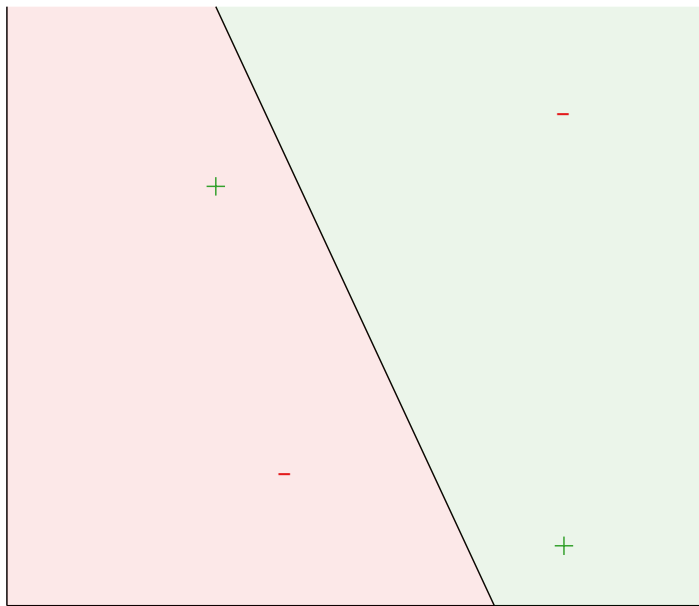


h_1

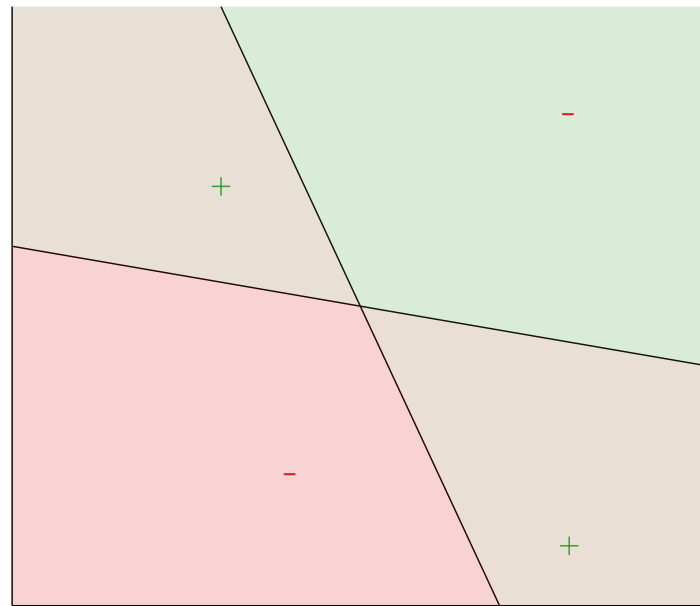


h_2

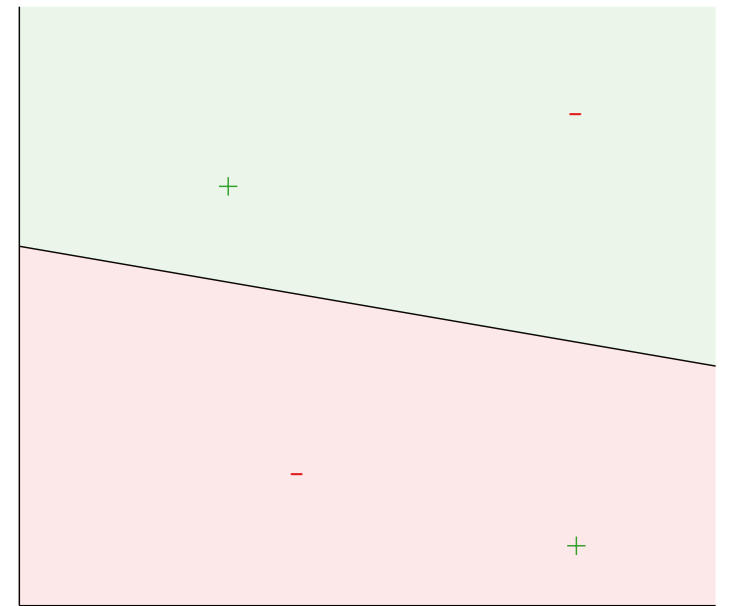
Combining Perceptrons



h_1



h_1



h_2

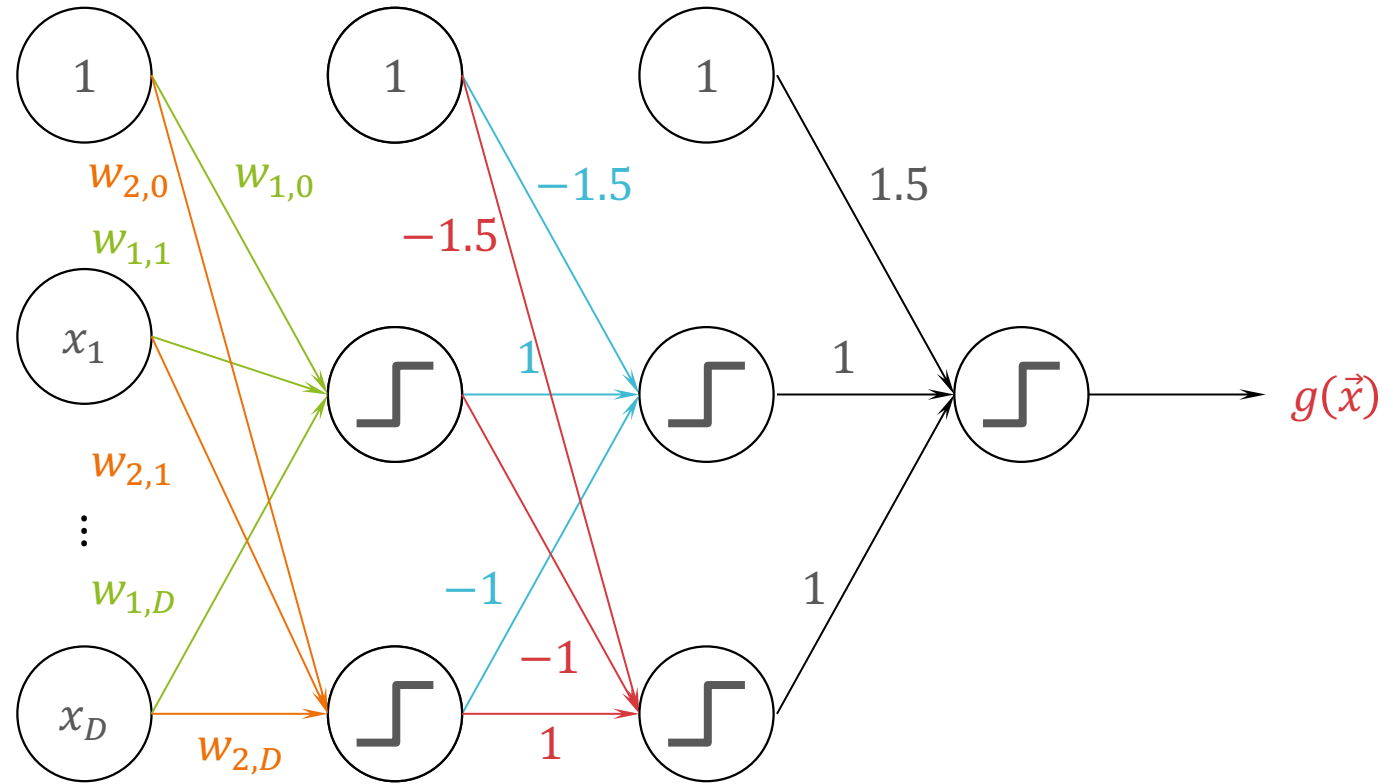
$$g(\vec{x}) = OR \left(AND(h_1(\vec{x}), \overline{h_2(\vec{x})}), AND(\overline{h_1(\vec{x})}, h_2(\vec{x})) \right)$$

Boolean Algebra

- Boolean variables are either $+1$ ("true") or -1 ("false")
- Basic Boolean operations
 - Negation: $\bar{z} = -1 * z$
 - And: $AND(z_1, z_2) = \text{sign}(z_1 + z_2 - 1.5)$
 - Or: $OR(z_1, z_2) = \text{sign}(z_1 + z_2 + 1.5)$

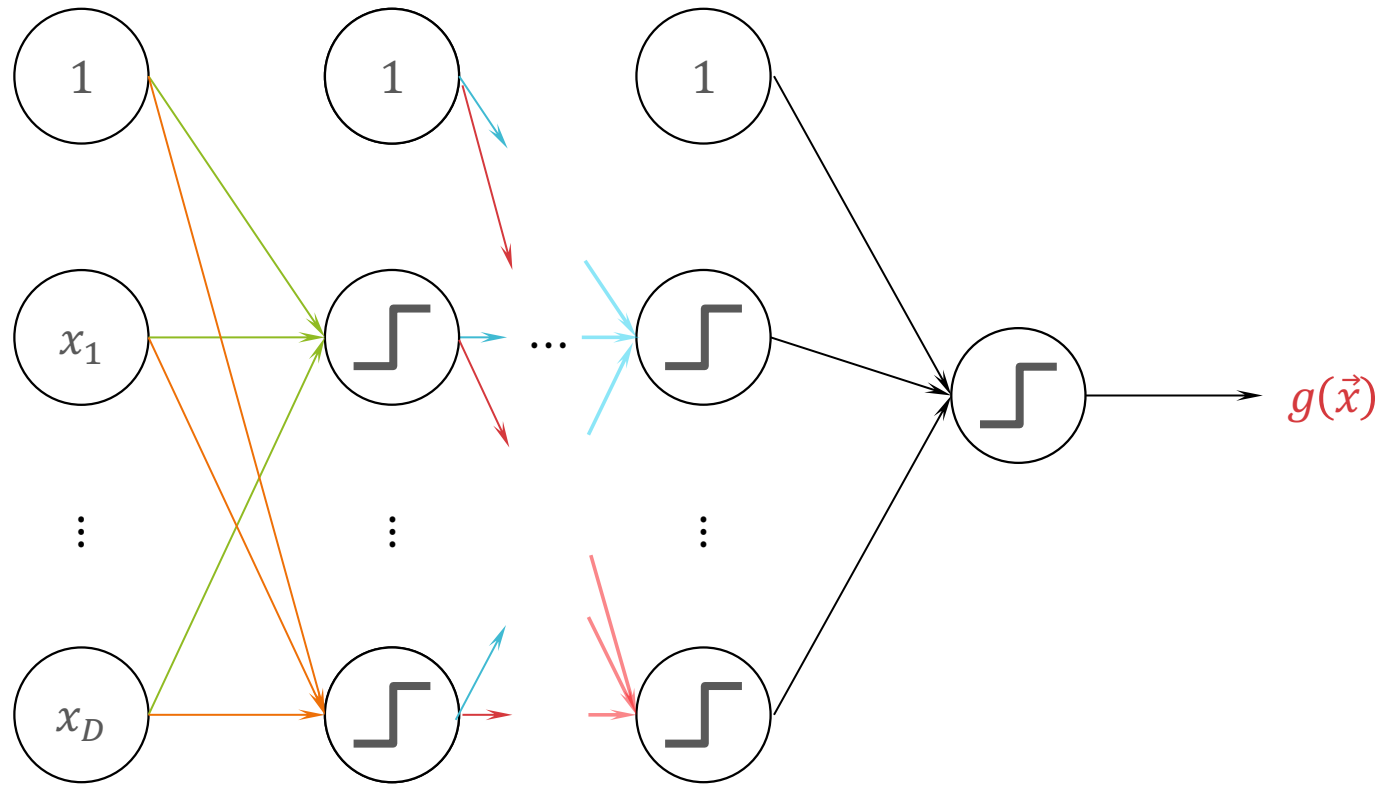
Building a Network

$$g(\vec{x}) = OR \left(AND(h_1(\vec{x}), \overline{h_2(\vec{x})}), AND(\overline{h_1(\vec{x})}, h_2(\vec{x})) \right)$$

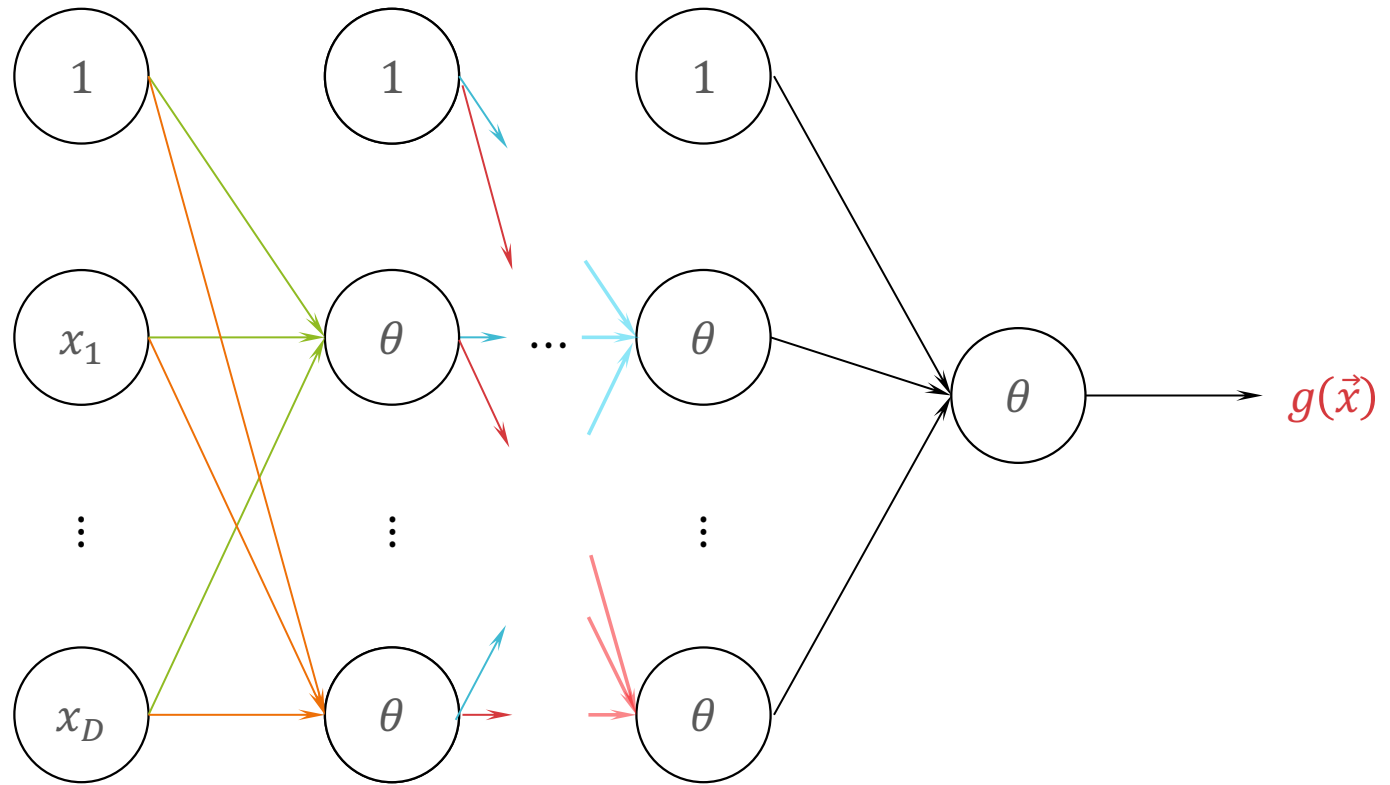


$$g(\vec{x}) = \text{sign} \left(\text{sign} \left(\text{sign}(\overline{w_1^T \vec{x}}) - \text{sign}(\overline{w_2^T \vec{x}}) - 1.5 \right) + \text{sign} \left(-\text{sign}(\overline{w_1^T \vec{x}}) + \text{sign}(\overline{w_2^T \vec{x}}) - 1.5 \right) + 1.5 \right)$$

Multi-Layer Perceptron (MLP)

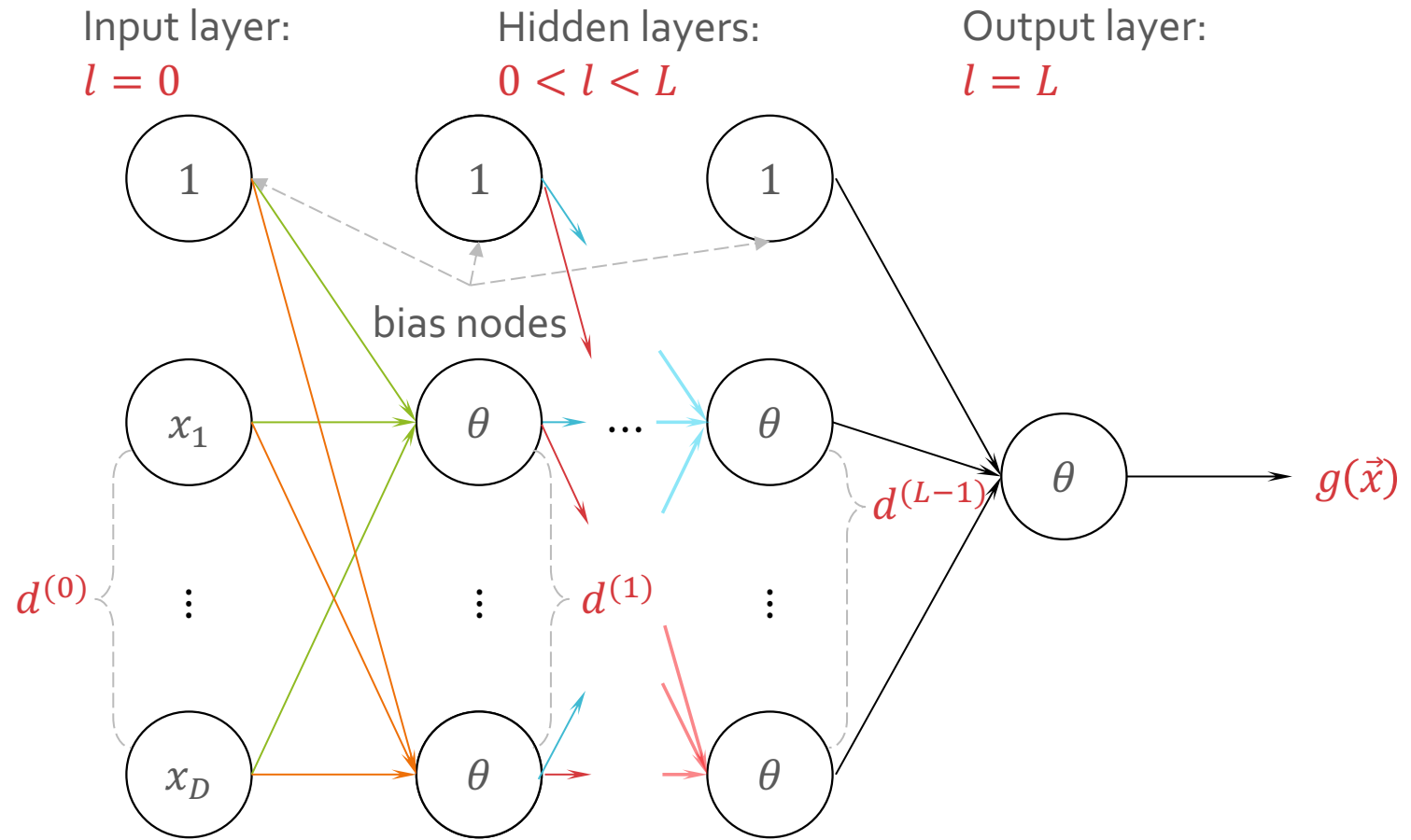


Feed-Forward Neural Network (NN)



What happens if $\theta(\cdot)$ is a linear function e.g. $\theta(x) = Cx$ for some constant C ?

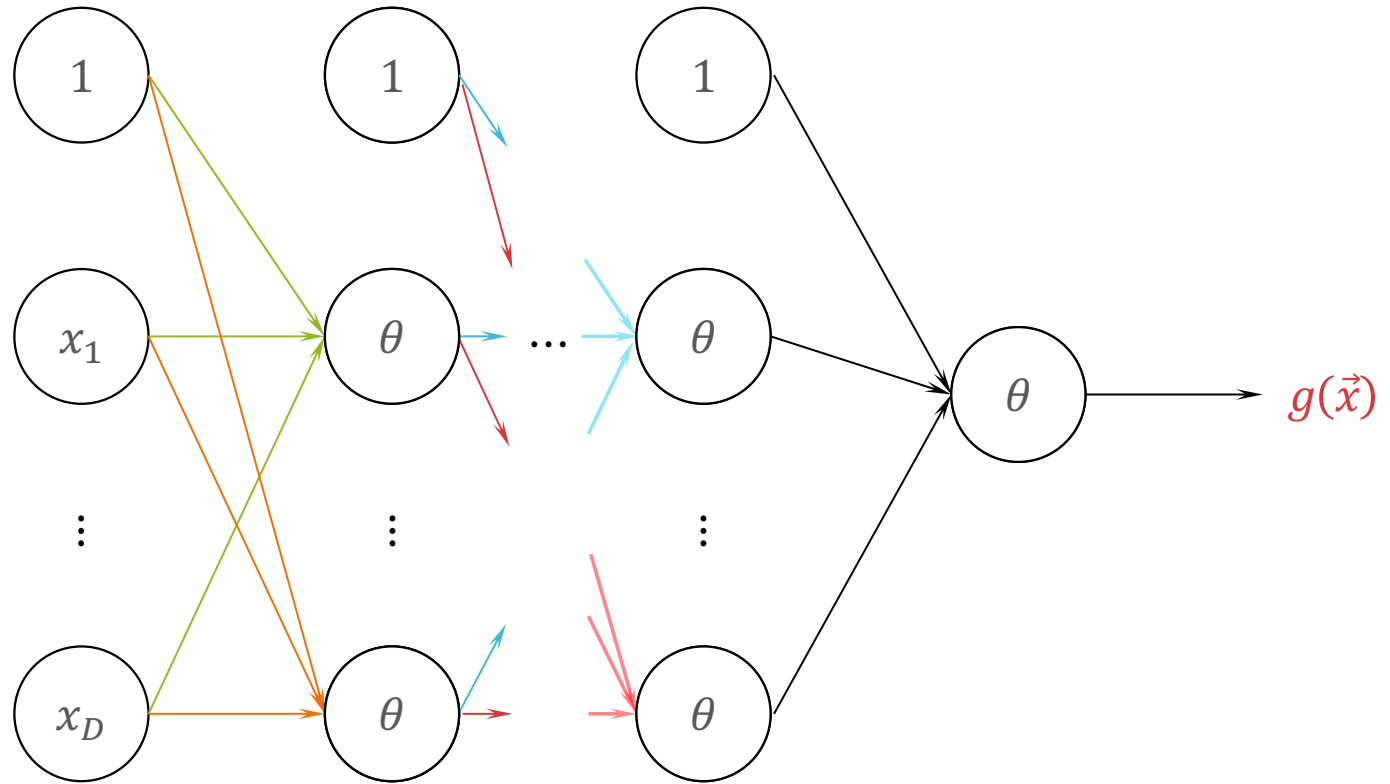
Feed-Forward Neural Network (NN)



Layer l has dimension $d^{(l)}$ \rightarrow Layer l has $d^{(l)} + 1$ nodes, counting the bias node

Architecture

The architecture of a NN is the vector $\vec{d} = [d^{(0)}, d^{(1)}, \dots, d^{(L)}]$

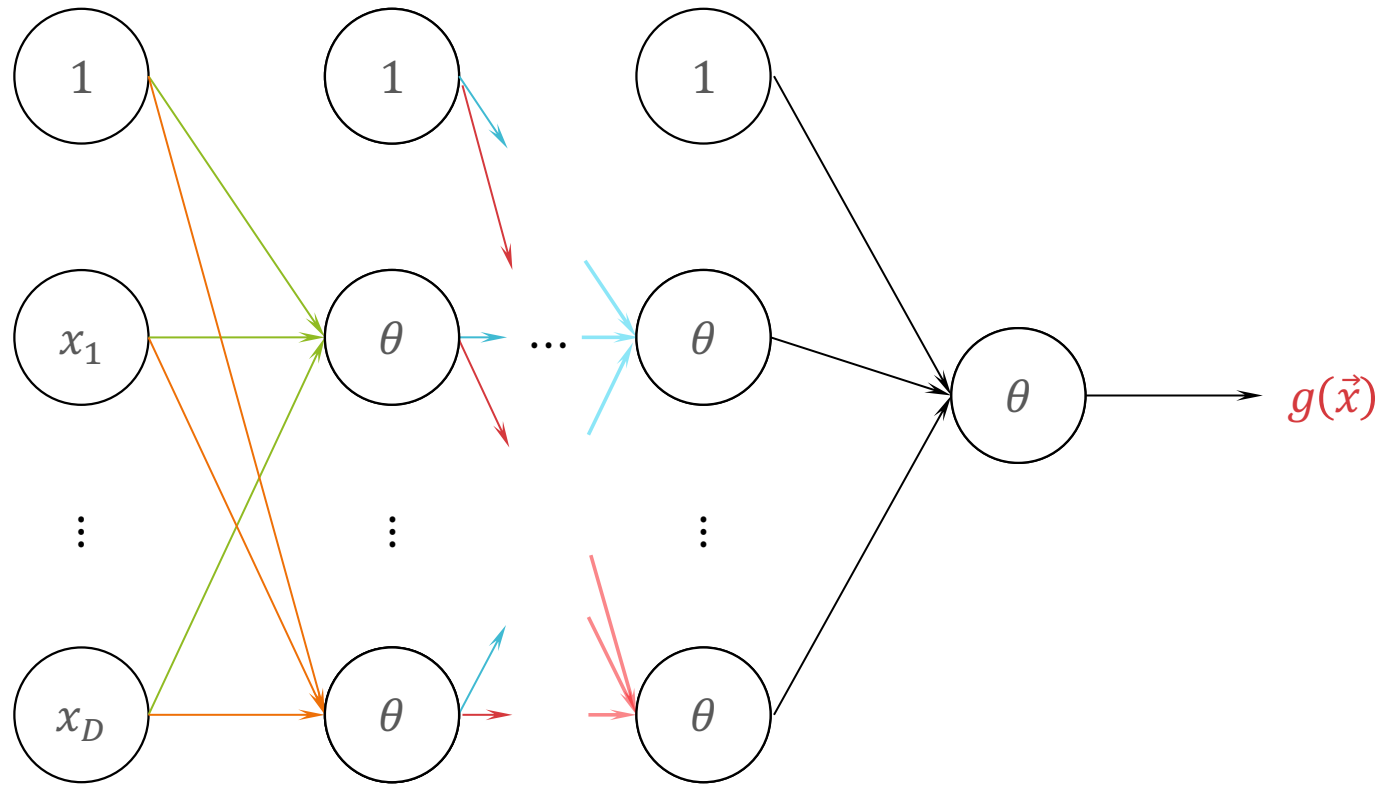


Every architecture corresponds to a hypothesis set: $\mathcal{H}_{NN}(\vec{d})$

A hypothesis $h \in \mathcal{H}_{NN}(\vec{d})$ is specified by setting all the weights between nodes

Weights

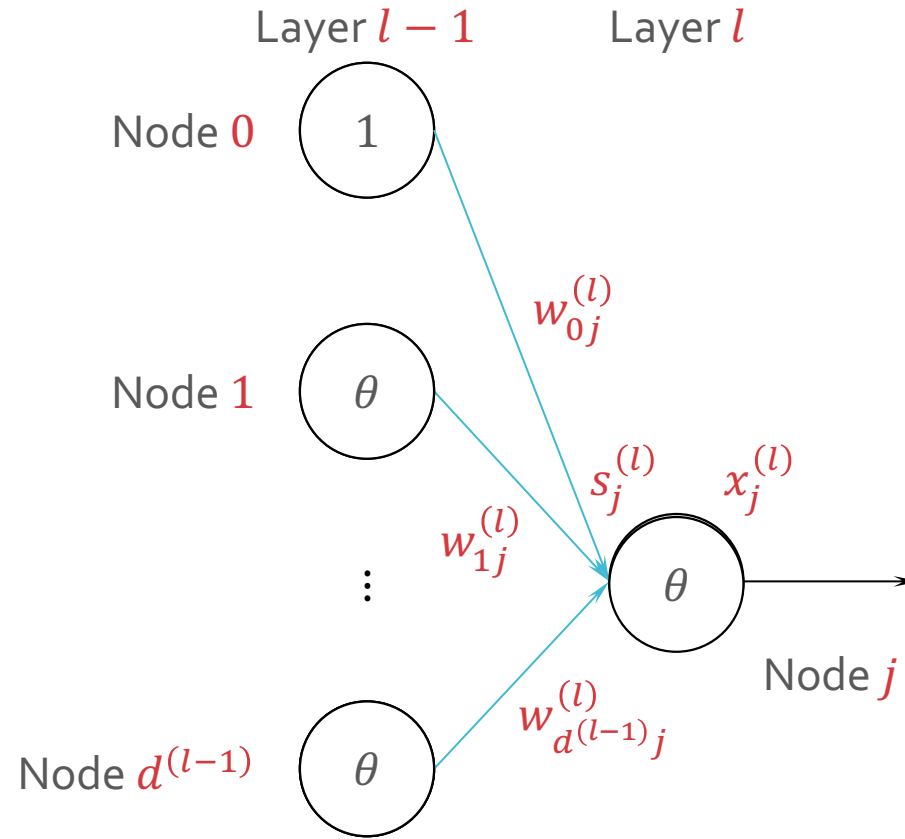
The weights between layer $l - 1$ and layer l are a matrix: $W^{(l)} \in \mathbb{R}^{(d^{(l-1)}+1) \times d^{(l)}}$



$w_{ij}^{(l)}$ is the weight between node i in layer $l - 1$ and node j in layer l

Signal and Outputs

Every node has an incoming signal and outgoing output



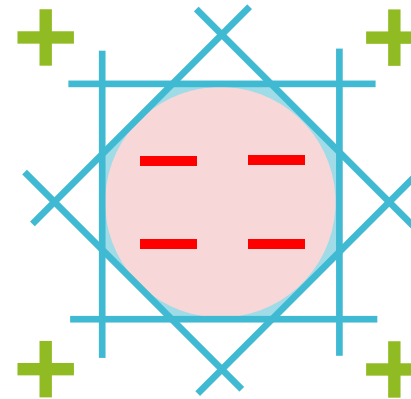
$$\overrightarrow{x^{(l)}} = \begin{bmatrix} 1 \\ \theta \left(\overrightarrow{s^{(l)}} \right) \end{bmatrix} \text{ and } \overrightarrow{s^{(l)}} = W^{(l)T} \overrightarrow{x^{(l-1)}}$$

Forward Propagation

- Input: weights $W^{(1)}, \dots, W^{(L)}$ and a query point \vec{x}
- Initialize $\overrightarrow{x^{(0)}} = \begin{bmatrix} 1 \\ \vec{x} \end{bmatrix}$
- For $l = 1, \dots, L$
 - $\overrightarrow{s^{(l)}} = W^{(l)T} \overrightarrow{x^{(l-1)}}$
 - $\overrightarrow{x^{(l)}} = \begin{bmatrix} 1 \\ \theta(\overrightarrow{s^{(l)}}) \end{bmatrix}$
- Output: $\overrightarrow{x^{(1)}}, \dots, \overrightarrow{x^{(L)}}$

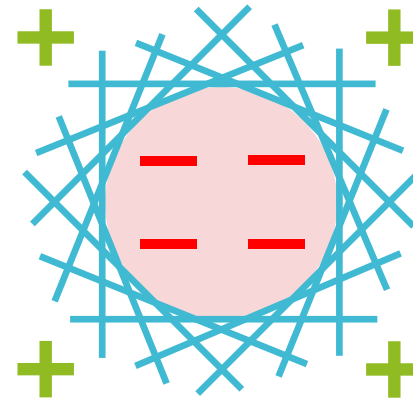
MLPs as Universal Approximators

- Theorem: any function that can be decomposed into perceptrons can be modelled exactly using a 3-layer MLP
- Any smooth decision boundary can be approximated to an arbitrary precision using a finite number of perceptrons



MLPs as Universal Approximators

- Theorem: any function that can be decomposed into perceptrons can be modelled exactly using a 3-layer MLP
- Any smooth decision boundary can be approximated to an arbitrary precision using a finite number of perceptrons



- Theorem: Any smooth decision boundary can be approximated to an arbitrary precision using a 3-layer MLP

NNs as Universal Approximators

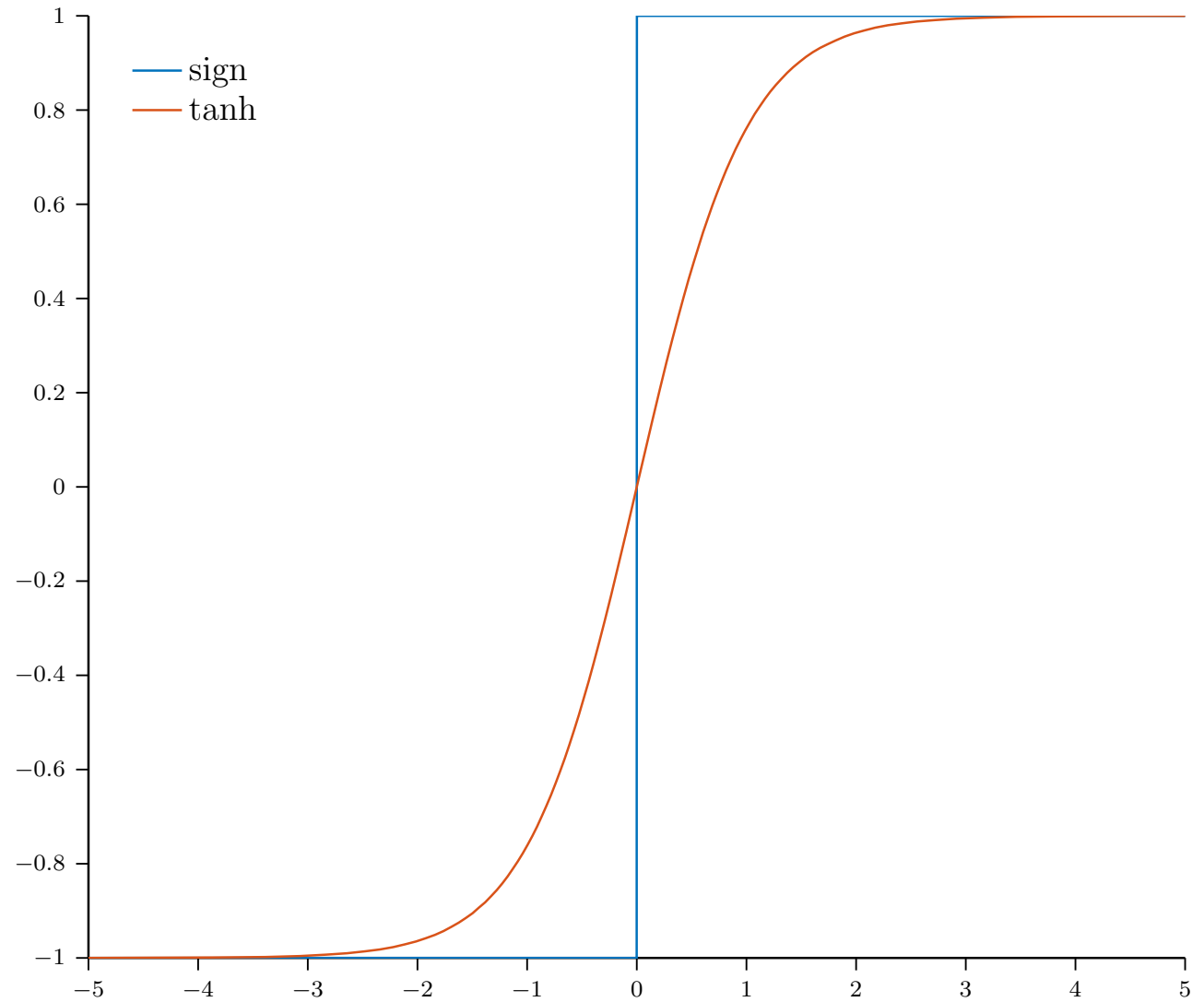
- Theorem: Any smooth decision boundary can be approximated to an arbitrary precision using a 2-layer feed-forward NN if the activation function, θ , is continuous, bounded and non-constant.

$\theta(\cdot)$

- Hyperbolic tangent:

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)} = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

- $\frac{\partial \tanh(z)}{\partial z} = 1 - \tanh(z)^2$



Error of a Neural Network

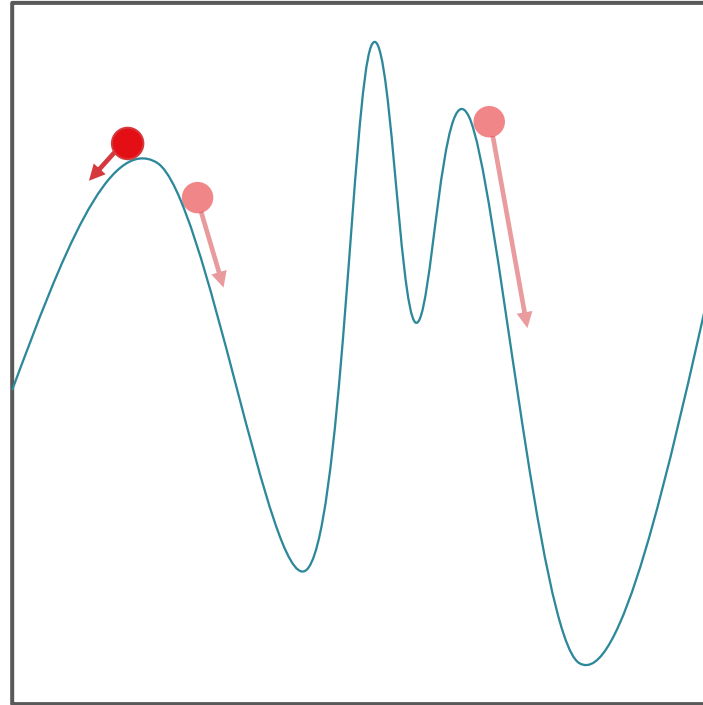
$$E_{in}(W^{(1)}, \dots, W^{(L)}) = \frac{1}{n} \sum_{i=1}^n (h(\vec{x}_i | W^{(1)}, \dots, W^{(L)}) - y_i)^2$$



awkward silence

Recall: Gradient Descent

- Iterative method for minimizing functions
- Requires the gradient to exist everywhere



Gradient Descent for Neural Networks

- Input: $\mathcal{D} = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}, \eta_0$
- Initialize all weights $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$ to small, random numbers and set $t = 0$
- While some termination condition is not satisfied
 - For $l = 1, \dots, L$
 - Compute $G^{(l)} = \nabla_{W^{(l)}} E_{in} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$
 - Update $W^{(l)}$: $W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)}$
 - Increment t : $t = t + 1$
- Output: $W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}$

Computing Gradients

$$E_{in} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) = \frac{1}{n} \sum_{i=1}^n e \left(h \left(\vec{x}_i \mid W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right), y_i \right)$$

$$\nabla_{W^{(l)}} E_{in} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$$

$$= \begin{bmatrix} \frac{\partial E_{in}}{\partial w_{01}^{(l)}} & \frac{\partial E_{in}}{\partial w_{02}^{(l)}} & \dots & \frac{\partial E_{in}}{\partial w_{0d^{(l)}}^{(l)}} \\ \frac{\partial E_{in}}{\partial w_{11}^{(l)}} & \frac{\partial E_{in}}{\partial w_{12}^{(l)}} & \dots & \frac{\partial E_{in}}{\partial w_{1d^{(l)}}^{(l)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial E_{in}}{\partial w_{d^{(l-1)}1}^{(l)}} & \frac{\partial E_{in}}{\partial w_{d^{(l-1)}2}^{(l)}} & \dots & \frac{\partial E_{in}}{\partial w_{d^{(l-1)}d^{(l)}}^{(l)}} \end{bmatrix}$$

Computing Gradients

$$E_{in} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) = \frac{1}{n} \sum_{i=1}^n e \left(h \left(\vec{x}_i \mid W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right), y_i \right)$$

$$\frac{\partial E_{in} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)}{\partial w_{ab}^{(l)}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial e \left(h \left(\vec{x}_i \mid W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right), y_i \right)}{\partial w_{ab}^{(l)}}$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial e \left(\vec{x}^{(L)}, y_i \right)}{\partial w_{ab}^{(l)}}$$

Computing Gradients

Computing $\nabla_{W^{(l)}} \text{Ein} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

$$\frac{\partial e \left(\vec{x}^{(L)}, y_i \right)}{\partial w_{ab}^{(l)}}$$

Insight: $w_{ab}^{(l)}$ only affects $e \left(\vec{x}^{(L)}, y_i \right)$ via $s_b^{(l)}$

Use the chain rule: $\frac{\partial e \left(\vec{x}^{(L)}, y_i \right)}{\partial w_{ab}^{(l)}} = \frac{\partial e \left(\vec{x}^{(L)}, y_i \right)}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{ab}^{(l)}} \right)$

Computing Gradients

To compute $\frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{ab}^{(l)}} \right)$, recall:

$$s_b^{(l)} = \sum_{a=0}^{d^{(l-1)}} w_{ab}^{(l)} x_a^{(l-1)} \rightarrow \frac{\partial s_b^{(l)}}{\partial w_{ab}^{(l)}} = x_a^{(l-1)}$$

Can compute all outputs $\vec{x}^{(l)} \forall l \in \{0, \dots, L\}$ using forward propagation

$$\text{Let } \delta_b^{(l)} = \frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial s_b^{(l)}}$$

$\vec{\delta}^{(l)}$ is called the sensitivity vector for layer l

Computing Gradients

Insight: $s_b^{(l)}$ only affects $e(\vec{x}^{(L)}, y_i)$ via $x_b^{(l)}$

Use the chain rule:
$$\delta_b^{(l)} = \frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial x_b^{(l)}} \left(\frac{\partial x_b^{(l)}}{\partial s_b^{(l)}} \right)$$

Recall: $x_b^{(l)} = \theta(s_b^{(l)}) \rightarrow \frac{\partial x_b^{(l)}}{\partial s_b^{(l)}} = \frac{\partial \theta(s_b^{(l)})}{\partial s_b^{(l)}}$

$$= 1 - (x_b^{(l)})^2$$

when $\theta(\cdot) = \tanh(\cdot)$

Computing Gradients

Insight: $x_b^{(l)}$ affects $e(\vec{x}^{(L)}, y_i)$ via $s_1^{(l+1)}, \dots, s_{d^{(l+1)}}^{(l+1)}$

Use the chain rule:
$$\frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial x_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial s_c^{(l+1)}} \left(\frac{\partial s_c^{(l+1)}}{\partial x_b^{(l)}} \right)$$
$$= \sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \left(w_{bc}^{(l+1)} \right)$$

Computing Gradients

$$\begin{aligned}\delta_b^{(l)} &= \frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial x_b^{(l)}} \left(\frac{\partial x_b^{(l)}}{\partial s_b^{(l)}} \right) \\ &= \left(\sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \left(w_{bc}^{(l+1)} \right) \right) \left(1 - \left(x_b^{(l)} \right)^2 \right)\end{aligned}$$

$$\vec{\delta}^{(l)} = W^{(l+1)} \vec{\delta}^{(l+1)} \otimes \left(1 - \vec{x}^{(l)} \otimes \vec{x}^{(l)} \right)$$

where \otimes is the element-wise product operation

$$\frac{\partial e(\vec{x}^{(L)}, y_i)}{\partial w_{ab}^{(l)}} = \delta_b^{(l)} \left(\frac{\partial s_b^{(l)}}{\partial w_{ab}^{(l)}} \right) = \delta_b^{(l)} \left(x_a^{(l-1)} \right)$$

$$\nabla_{W^{(l)}} e(\vec{x}^{(L)}, y_i) = \vec{x}^{(l-1)} \left(\vec{\delta}^{(l)} \right)^T$$

Computing Gradients

Can recursively compute $\overrightarrow{\delta^{(l)}}$ using $\overrightarrow{\delta^{(l+1)}}$; need to compute the base case: $\overrightarrow{\delta^{(L)}}$

Assume the output layer is a single node and the error function is the squared error:

$$\overrightarrow{\delta^{(L)}} = \delta_1^{(L)}, \overrightarrow{x^{(L)}} = x_1^{(L)} \text{ and } e(x_1^{(L)}, y_i) = (x_1^{(L)} - y_i)^2$$

$$\delta_1^{(L)} = \frac{\partial e(x_1^{(L)}, y_i)}{\partial s_1^{(L)}} = \frac{\delta}{\partial s_1^{(L)}} (x_1^{(L)} - y_i)^2$$

$$= 2(x_1^{(L)} - y_i) \frac{\delta x_1^{(L)}}{\partial s_1^{(L)}} = 2(x_1^{(L)} - y_i)^2 (1 - (x_1^{(L)})^2)$$

when $\theta(\cdot) = \tanh(\cdot)$

Back- propagation

- Input: weights $W^{(1)}, \dots, W^{(L)}$ and a query point \vec{x}
- Run forward propagation to get $\vec{x}^{(1)}, \dots, \vec{x}^{(L)}$
- Initialize $\delta_1^{(L)} = 2 \left(x_1^{(L)} - y_i \right) \left(1 - \left(x_1^{(L)} \right)^2 \right)$
- For $l = L - 1, \dots, 1$
 - Compute $\vec{\delta}^{(l)} = W^{(l+1)} \vec{\delta}^{(l+1)} \otimes \left(1 - \vec{x}^{(l)} \otimes \vec{x}^{(l)} \right)$
- Output: $\vec{\delta}^{(1)}, \dots, \vec{\delta}^{(L)}$

Computing Gradients

- Input: $W^{(1)}, \dots, W^{(L)}$ and $\mathcal{D} = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$
- Initialize $E_{in} = 0$ and $G^{(l)} = 0 * W^{(l)}$ for $l = 1, \dots, L$
- For $i = 1, \dots, n$
 - Run forward propagation to get $\vec{x}^{(1)}, \dots, \vec{x}^{(L)}$
 - Run backpropagation to get $\vec{\delta}^{(1)}, \dots, \vec{\delta}^{(L)}$
 - Increment E_{in} : $E_{in} = E_{in} + \frac{1}{n} (\vec{x}^{(L)} - y_i)^2$
 - For $l = 1, \dots, L$
 - Compute $G_i^{(l)} = \vec{x}^{(l-1)} (\vec{\delta}^{(l)})^T$
 - Increment $G^{(l)}$: $G^{(l)} = G^{(l)} + \frac{1}{n} G_i^{(l)}$
- Output: $G^{(1)}, \dots, G^{(L)}$, the gradients of E_{in} w.r.t $W^{(1)}, \dots, W^{(L)}$

Complexity

- Both forward and backpropagation contain matrix multiplications involving $W^{(1)}, \dots, W^{(L)}$ \rightarrow both take time $O(|W^{(1)}| + \dots + |W^{(L)}|)$...
- Computing $G^{(1)}, \dots, G^{(L)}$ requires running forward and backpropagation for each training point $(\vec{x}, y) \in \mathcal{D}$...
- Each iteration of gradient descent for a neural network takes time $O(n(|W^{(1)}| + \dots + |W^{(L)}|))$
- Use stochastic gradient descent instead!
- Also use parallelization and GPUs / TPUs!

Stochastic Gradient Descent for Neural Networks

- Input: $\mathcal{D} = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}, \eta_0$
- Initialize all weights $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$ to small, random numbers and set $t = 0$
- While some termination condition is not satisfied
 - For $l = 1, \dots, L$
 - Randomly select a point $(\vec{x}_i, y_i) \in \mathcal{D}$
 - Compute $G^{(l)'} = \nabla_{W^{(l)}} e \left(h \left(\vec{x}_i \mid W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right), y_i \right)$
 - Update $W^{(l)}$: $W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)'}$
 - Increment t : $t = t + 1$
- Output: $W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}$

Initialization and Termination

- Initialization:
 - Randomness is good for non-convex optimization
 - Initialize weights by sampling from $N(0, \sigma^2)$
- Termination:
 - For complicated surfaces, the gradient's magnitude is not a good metric for proximity to a minimum
 - A simple solution: combine multiple termination criteria e.g. stop if enough iterations have passed and the improvement in error is small