

CS60050

MACHINE LEARNING

Neural Networks - Backpropagation

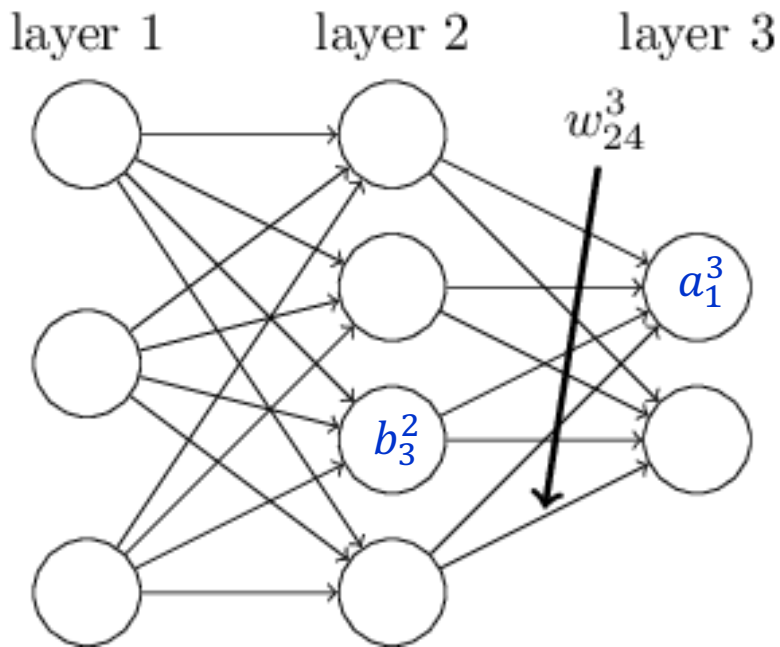
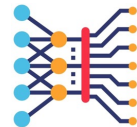
Somak Aditya
Assistant Professor

Sudeshna Sarkar

Department of CSE, IIT Kharagpur
Oct 12, 2023



Notations



w_{jk}^l is the weight from the k^{th} neuron in the $(l-1)^{\text{th}}$ layer to the j^{th} neuron in the l^{th} layer

b_j^l : Bias for j^{th} neuron in l^{th} layer.

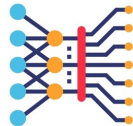
a_j^l : Activation of the j^{th} neuron in the l^{th} layer.

$$a_j^l = g(\sum_k w_{jk}^l a_k^{l-1} + b_j^l)$$

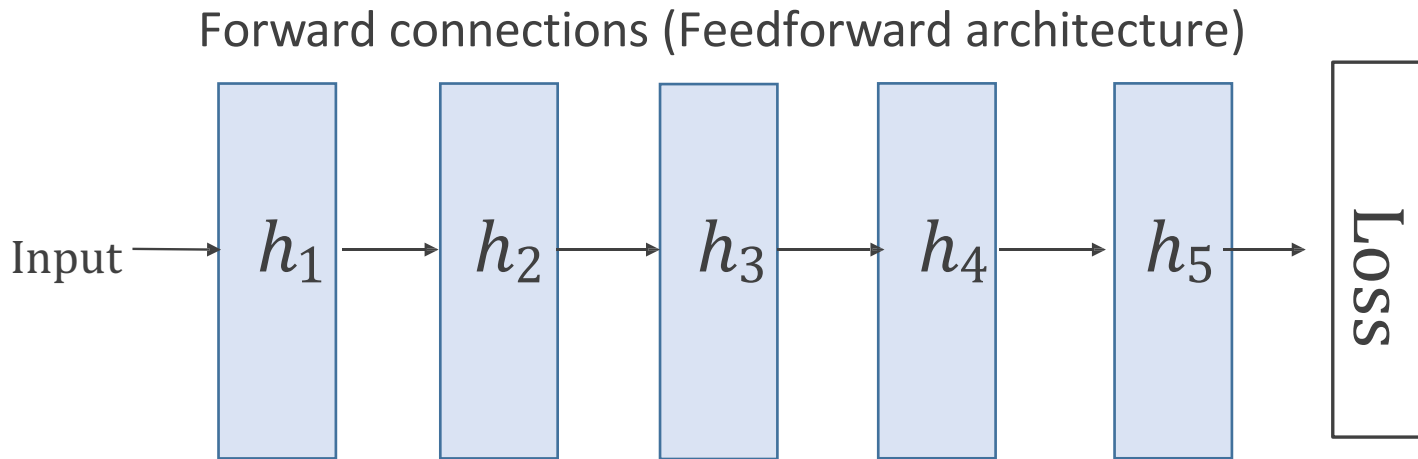
→
Vectorize

$$\begin{aligned} a^l &= g(w^l a^{l-1} + b^l) \\ z^l &= w^l a^{l-1} + b^l \\ a^l &= g(z^l) \end{aligned}$$

Neural networks in blocks

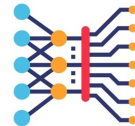


- We can visualize $a_L = h_L \circ h_{L-1} \circ \dots \circ h_1(x)$ as a cascade of blocks.



The activation functions must be **1st-order differentiable (almost) everywhere**

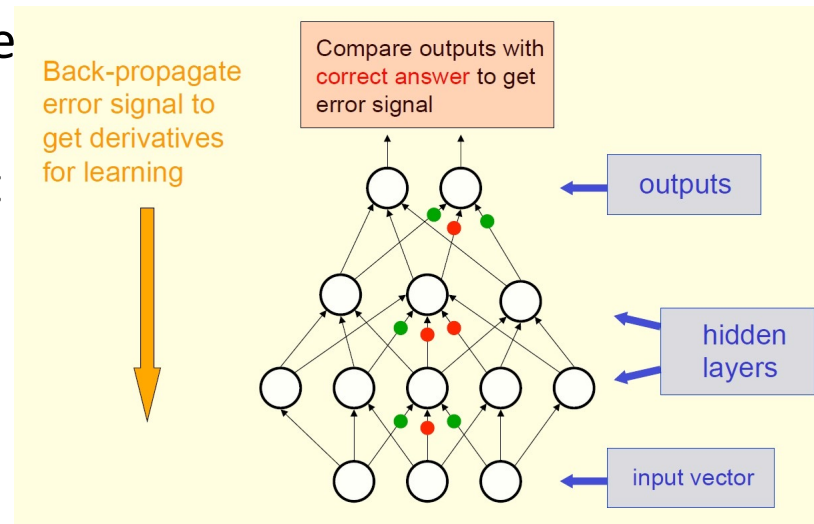
Backpropagation

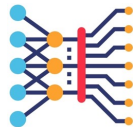


- Feedforward Propagation: Accept input $x^{(i)}$, pass through intermediate stages and obtain output $\hat{y}^{(i)}$
- During Training: Compute scalar cost $J(\theta)$

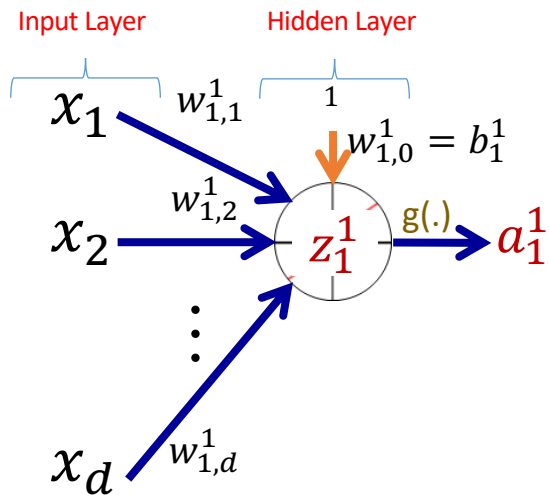
$$J(\theta) = \sum_i L(NN(x^{(i)}; \theta), y^{(i)})$$

- Backpropagation allows information to flow backwards from cost to compute the gradient





Multilayer Neural Network

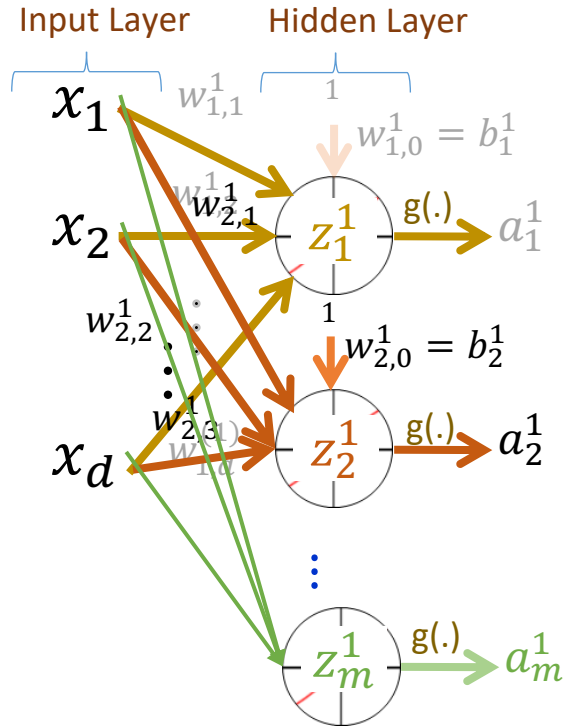


$$z_1^1 = b_1^1 + \sum_{i=1}^d w_{1,i}^1 x_i = [\mathbf{w}_1^1 b_1^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$a_1^1 = g(z_1^1)$$

$$[x_1 \ x_2 \ \dots \ x_d]^T$$

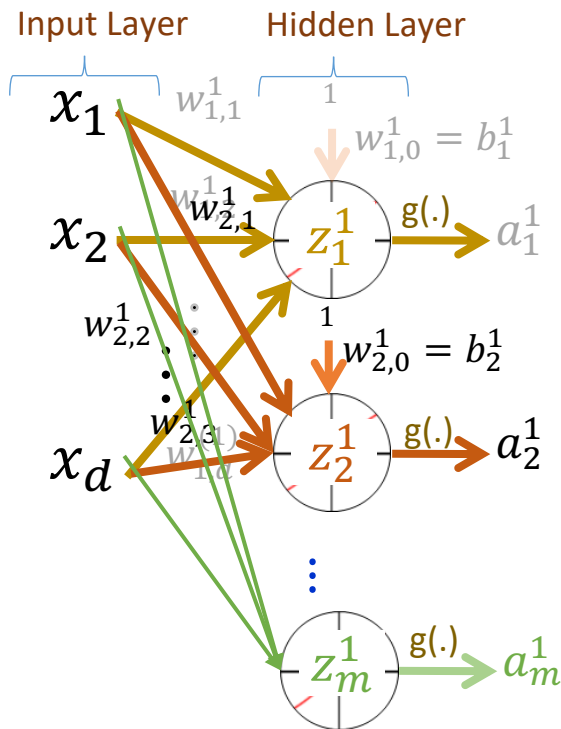
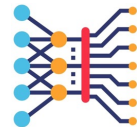
Multilayer Neural Network



$$\begin{aligned} a^{(0)} &= x \\ z^{(1)} &= \mathbf{w}^{(1)} \mathbf{a}^{(0)} \\ a^{(1)} &= g(z^{(1)}) \end{aligned}$$

$W^1 : m \times n$ matrix
 $b^1 : m \times 1$ column vector
 $X : d \times 1$ column vector
 $Z^1 : m \times 1$ column vector
 $A^1 : m \times 1$ column vector

Multilayer Neural Network



$$\begin{aligned} a^{(0)} &= x \\ z^{(1)} &= \mathbf{w}^{(1)} a^{(0)} \\ a^{(1)} &= g(z^{(1)}) \end{aligned}$$

$$\begin{bmatrix} a_1^1 \\ a_2^1 \\ \vdots \\ a_m^1 \end{bmatrix} = \begin{bmatrix} g(z_1^1) \\ g(z_2^1) \\ \vdots \\ g(z_m^1) \end{bmatrix}$$

$$\mathbf{a}^1 = g(\mathbf{z}^{(1)})$$

$$\mathbf{z}^1 = [\mathbf{W}^1 \mathbf{b}^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

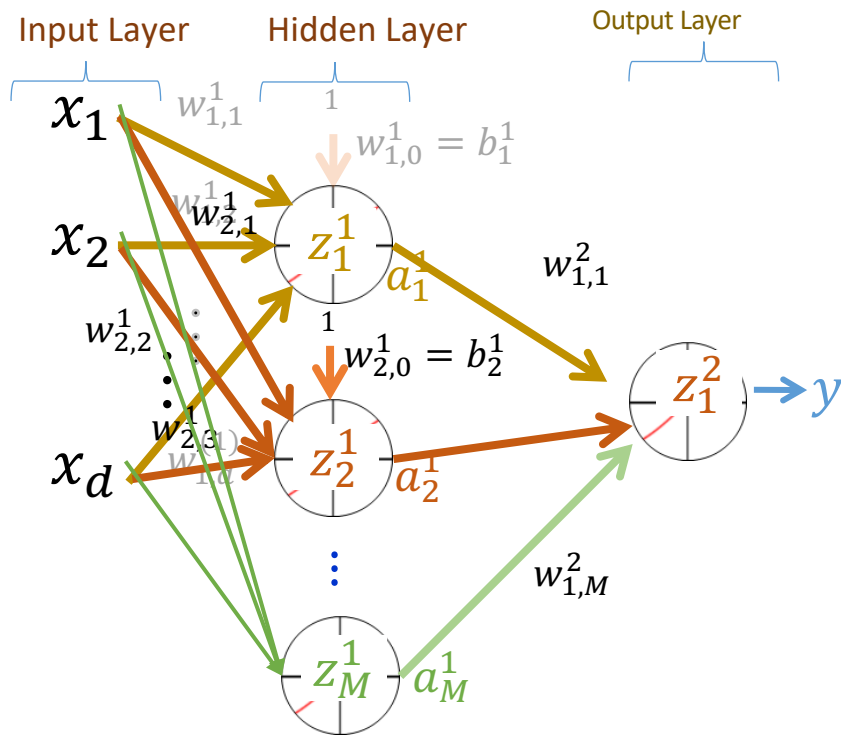
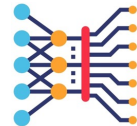
$$z_1^1 = [\mathbf{w}_1^1 b_1^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$z_2^1 = [\mathbf{w}_2^1 b_2^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$z_M^1 = [\mathbf{w}_m^1 b_m^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} z_1^1 \\ z_2^1 \\ \vdots \\ z_m^1 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^1 & b_1^1 \\ \mathbf{w}_2^1 & b_2^1 \\ \vdots & \vdots \\ \mathbf{w}_m^1 & b_m^1 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

Multilayer Neural Network



Output Layer Pre-activation

$$z_1^{(2)} = \begin{bmatrix} \mathbf{w}_1^{(2)} & b_1^{(2)} \end{bmatrix} \begin{bmatrix} \mathbf{a}^{(1)} \\ 1 \end{bmatrix}$$

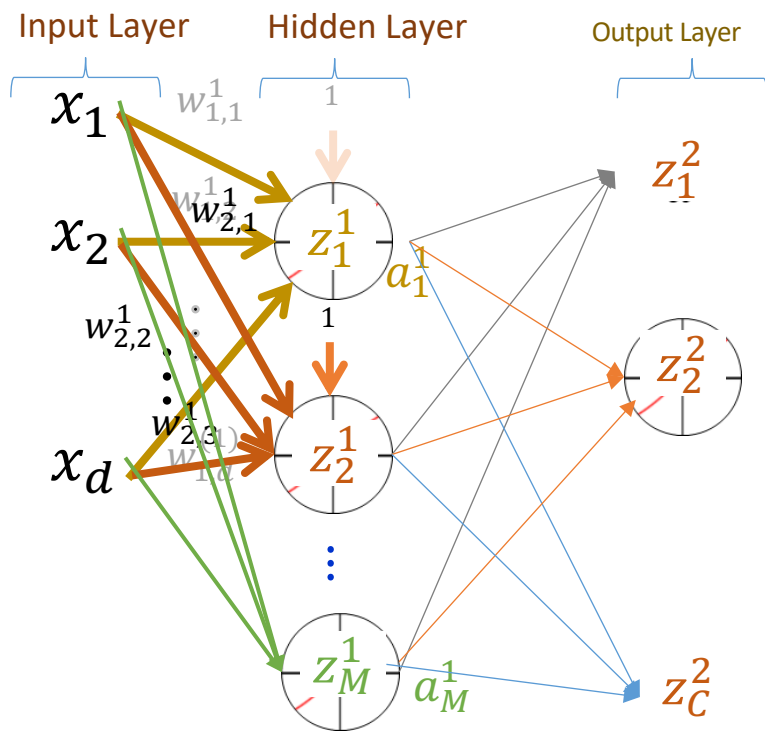
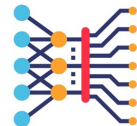
Output Layer Activation

$$y_1 = o(z_1^{(2)})$$

output

- Sigmoid for 2-class classification
- Softmax for multi-class classification
- Linear for regression

Multilayer Neural Network

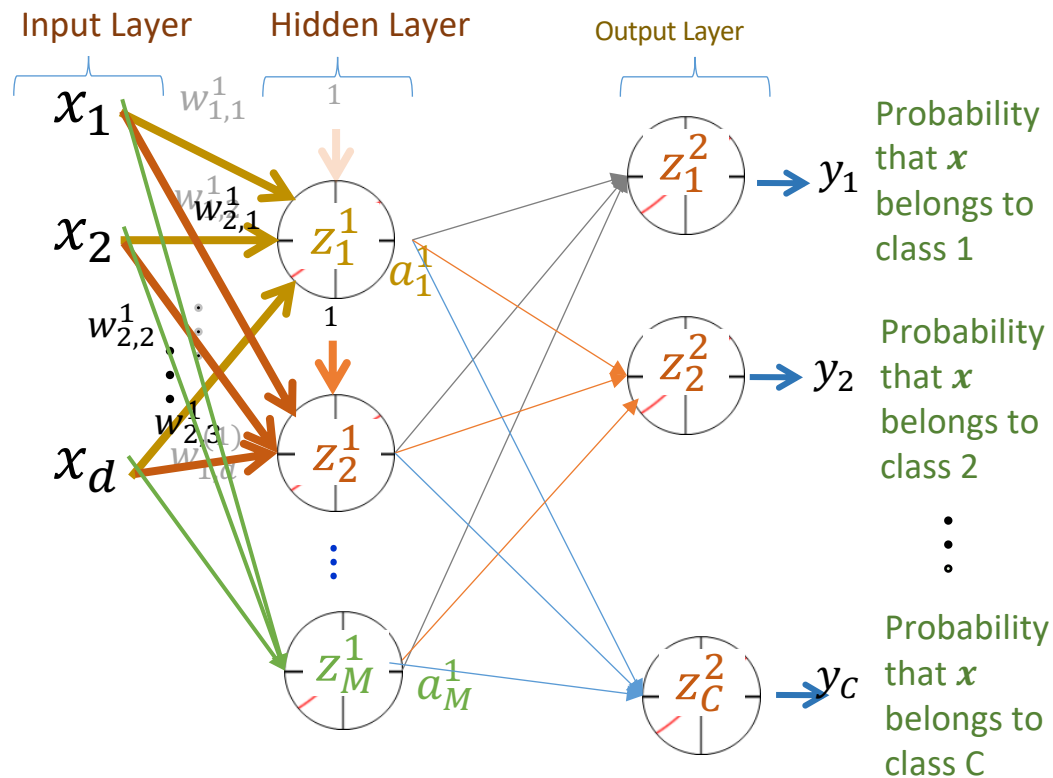
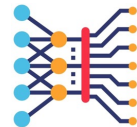


$$\rightarrow y_1 = o_1 \left(z_1^{(2)} \right) = \frac{\exp(z_1^{(2)})}{\sum_c \exp(z_c^{(2)})}$$

$$\rightarrow y_2 = o_2 \left(z_2^{(2)} \right) = \frac{\exp(z_2^{(2)})}{\sum_c \exp(z_c^{(2)})}$$

$$\rightarrow y_c = o_c \left(z_c^{(2)} \right) = \frac{\exp(z_c^{(2)})}{\sum_c \exp(z_c^{(2)})}$$

Training a Neural Network – Loss Function



Aim to maximize the probability corresponding to the correct class for any example x

$$\begin{aligned} \max y_c \\ \equiv \max (\log y_c) \\ \equiv \min (-\log y_c) \end{aligned}$$

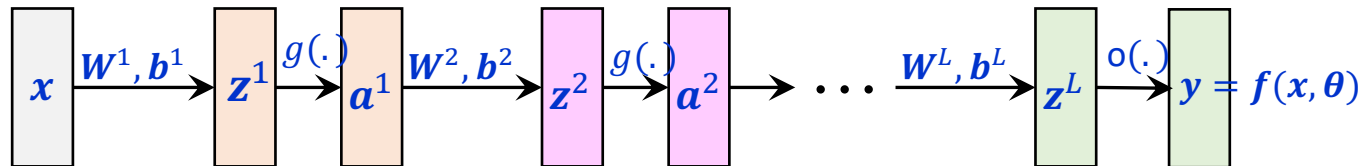
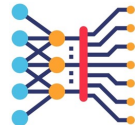
Can be equivalently expressed as

$$-\sum_i \prod_{i=c} \log(y_i)$$

known as cross-entropy loss

Forward Pass

θ is the collection
of all learnable
parameters i.e., all
 W and b



Hidden layer pre-activation:

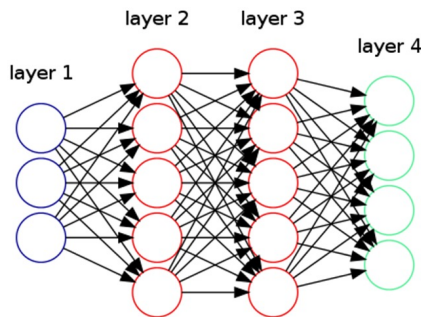
For $l = 1, \dots, L$; $\mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}$

Hidden layer activation:

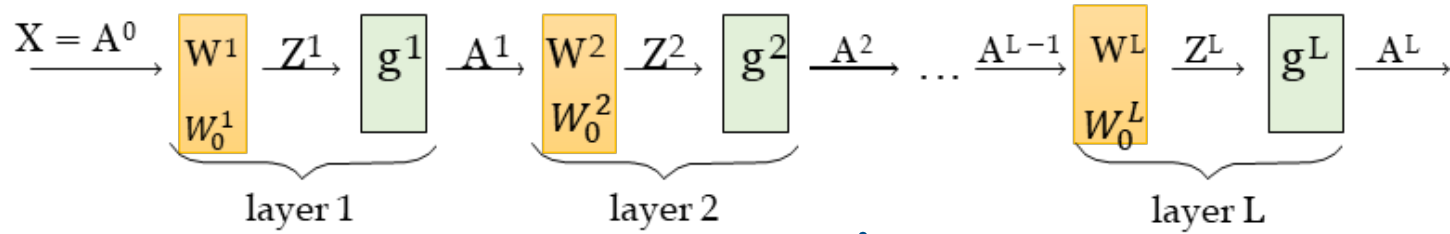
For $l = 1, \dots, L - 1$; $\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)})$

Output layer activation:

For $l = L$; $\mathbf{y} = \mathbf{a}^{(L)} = o(\mathbf{z}^{(L)}) = f(\mathbf{x}, \theta)$



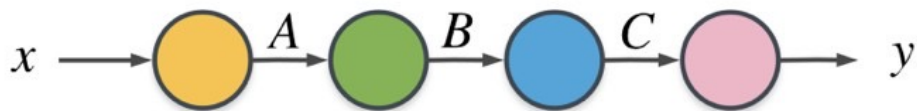
Error back-propagation



1. Compute Loss
2. Compute the derivative of the L w.r.t. the final output of the network A^L
3. Compute the derivative of L w.r.t. the pre-activation Z^L
4. Compute the derivative of L wrt W^L
5. Then compute the derivative of the L w.r.t. the final output of the network A^{L-1}
6. Then compute the derivative of L w.r.t. the pre-activation Z^{L-1}
7. Compute the derivative of L wrt W^{L-1}

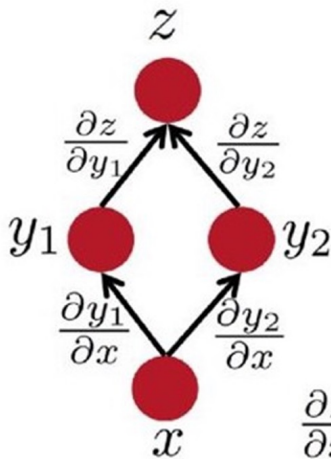
....

Chain Rule



$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial C} \times \frac{\partial C}{\partial B} \times \frac{\partial B}{\partial A} \times \frac{\partial A}{\partial x}$$

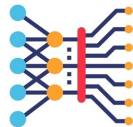
Multiple Path



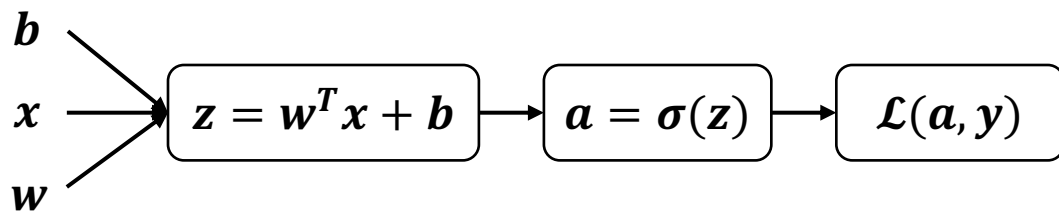
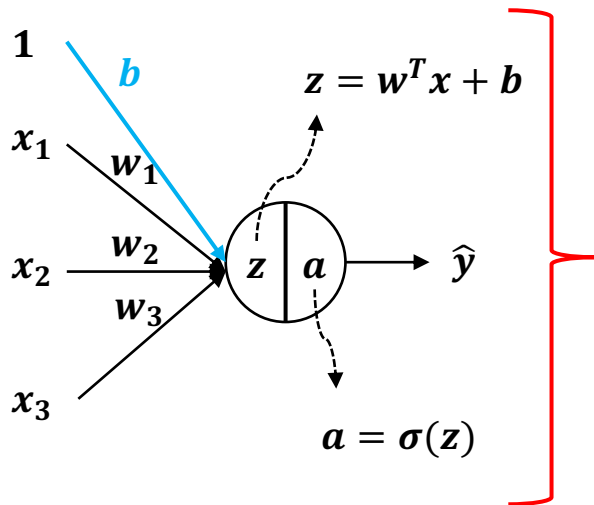
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x}$$

Backpropagation is just repeated application of the chain rule.

The Computation Graph of Logistic Regression

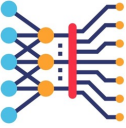


Let us translate logistic regression into a computation graph

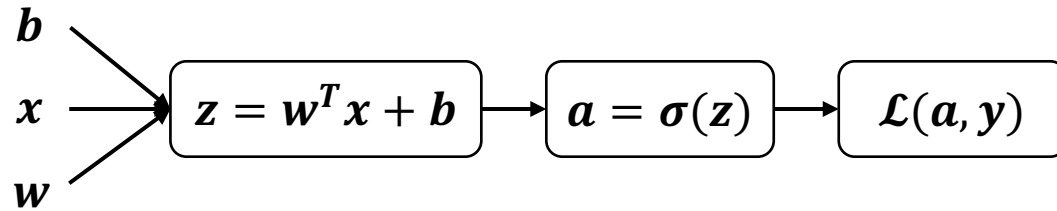


Where $\mathbf{b} = 1$, $\mathbf{w} = [w_1, w_2, w_3]$, $\mathbf{x} = [x_1, x_2, x_3]$,
and $\mathcal{L}(\mathbf{a}, \mathbf{y})$ is the cost (or *loss*) function

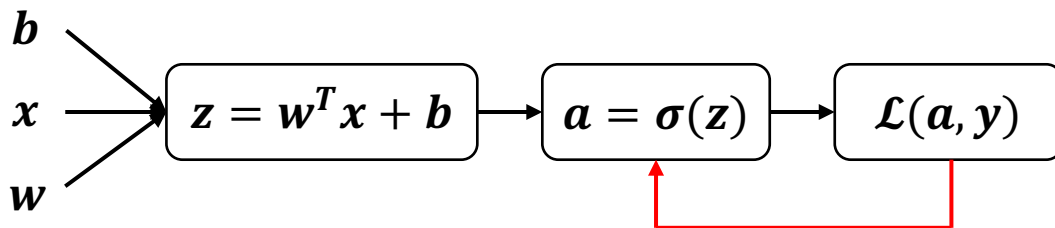
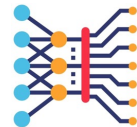
Forward Propagation



- The loss function can be computed by moving from left to right



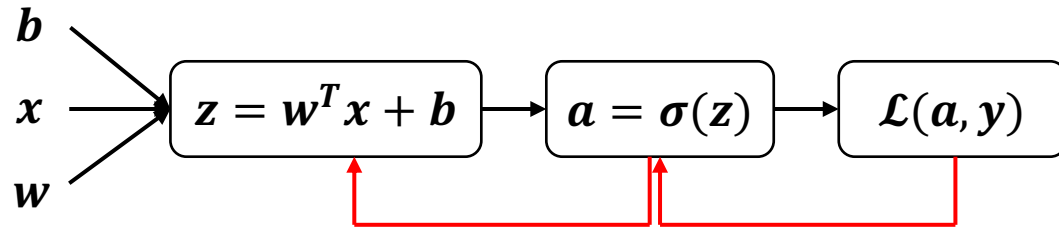
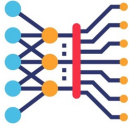
Backward Propagation



Partial derivative of \mathcal{L} with respect to a

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial a} &= \frac{\partial}{\partial a} (-y \log(a) - (1 - y) \log(1 - a)) \\ &= \frac{-y}{a} + \frac{(1 - y)}{(1 - a)}\end{aligned}$$

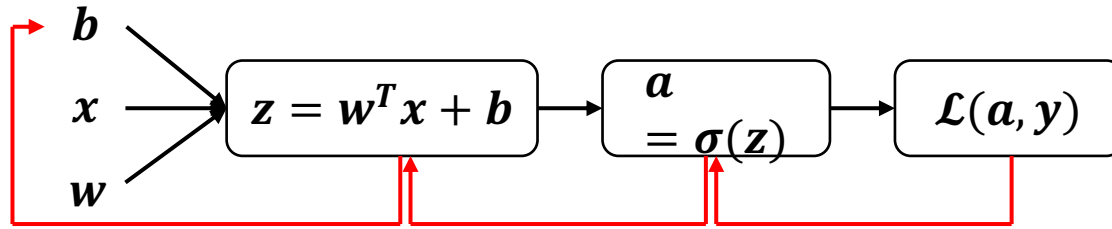
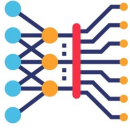
Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} = \left(\frac{-y}{a} + \frac{(1-y)}{(1-a)} \right) \times \frac{\partial a}{\partial z} = \left(\frac{-y}{a} + \frac{(1-y)}{(1-a)} \right) \times a(1-a)$$

$= a - y$

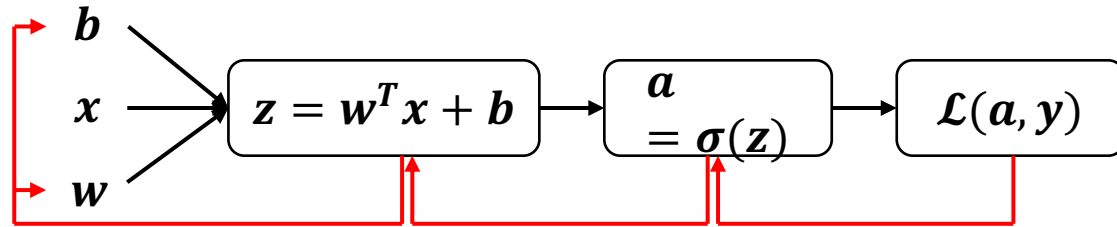
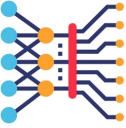
Backward Propagation



$\frac{\partial \mathcal{L}}{\partial b}$ = Partial derivative of \mathcal{L} with respect to b

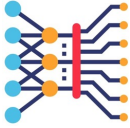
$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b} = (a - y) \times \frac{\partial z}{\partial b} = (a - y) \times 1 = (a - y)$$

Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}} \times \frac{\partial \mathbf{a}}{\partial \mathbf{z}} \times \frac{\partial \mathbf{z}}{\partial \mathbf{w}} = (\mathbf{a} - \mathbf{y}) \times \frac{\partial \mathbf{z}}{\partial \mathbf{w}} = (\mathbf{a} - \mathbf{y})\mathbf{x}$$

Backward Propagation: Summary

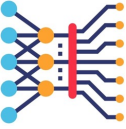


- Here is the summary of the gradients in logistic regression:

$$dz = \frac{\partial \mathcal{L}}{\partial z} = a - y$$

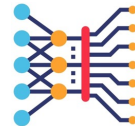
$$db = \frac{\partial \mathcal{L}}{\partial b} = a - y$$

$$dw = \frac{\partial \mathcal{L}}{\partial w} = (a - y)x$$



Backward Pass MLP

Optimizing the Neural Network



$$J(\Theta) = -\frac{1}{n} \left[\sum_{i=1}^n \sum_{k=1}^K y_{ik} \log(h_{\Theta}(\mathbf{x}_i))_k + (1 - y_{ik}) \log(1 - (h_{\Theta}(\mathbf{x}_i))_k) \right] \\ + \frac{\lambda}{2n} \sum_{l=1}^{L-1} \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} \left(\Theta_{ji}^{(l)} \right)^2$$

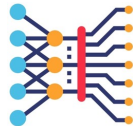
Solve via: $\min_{\Theta} J(\Theta)$

$J(\Theta)$ is not convex, so GD on a neural net yields a local optimum

- But, tends to work well in practice

Need code to compute:

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

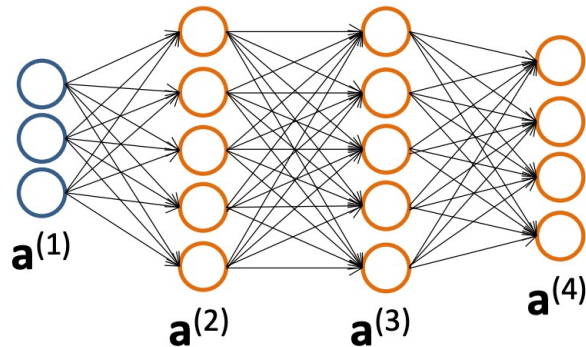


Forward Propagation

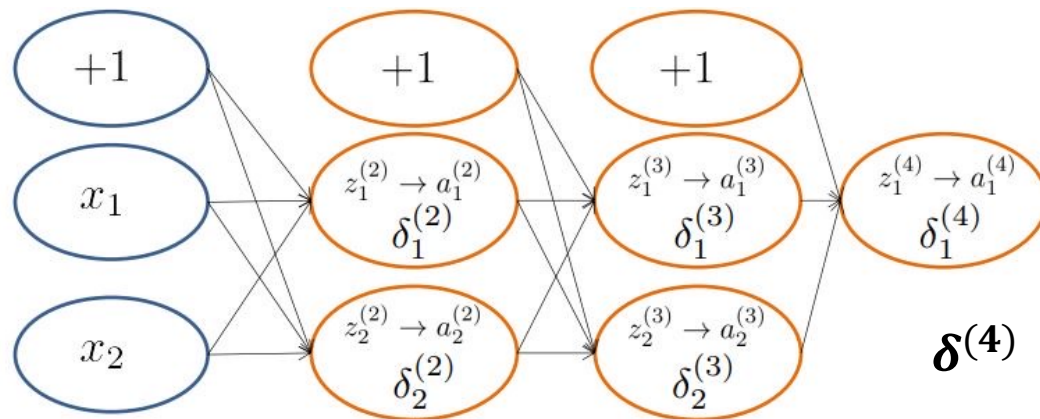
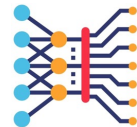
- Given one labeled training instance (\mathbf{x}, y) :

Forward Propagation

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$ [add $a_0^{(2)}$]
- $\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$ [add $a_0^{(3)}$]
- $\mathbf{z}^{(4)} = \Theta^{(3)}\mathbf{a}^{(3)}$
- $\mathbf{a}^{(4)} = h_{\Theta}(\mathbf{x}) = g(\mathbf{z}^{(4)})$



Backpropagation

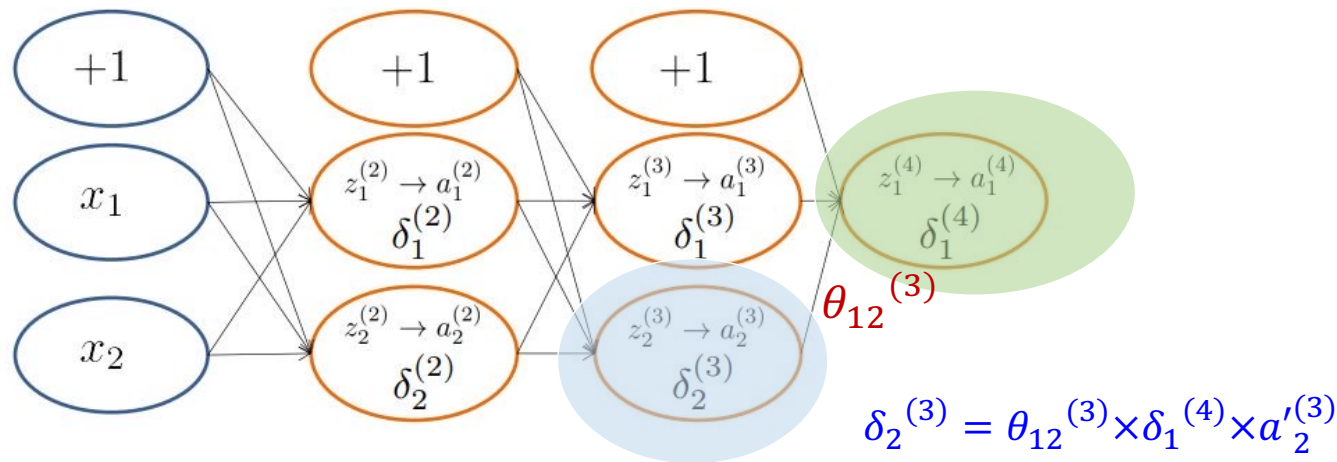
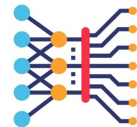


$\delta_j^{(l)}$ = “error” of node j in layer l

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = y_i \log h_{\Theta}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\Theta}(\mathbf{x}_i))$

Backpropagation

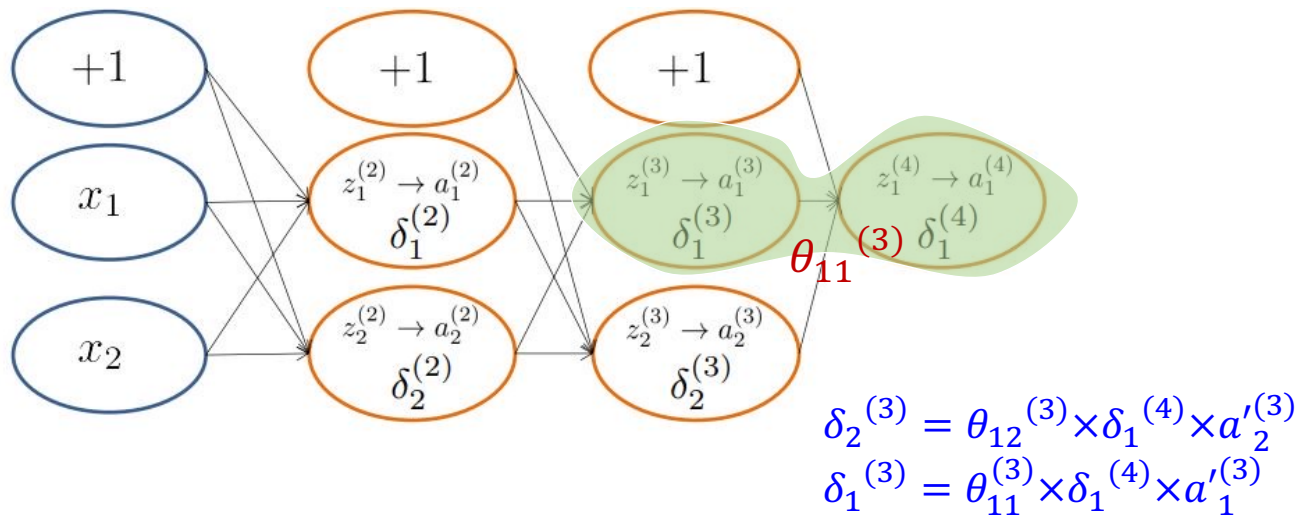
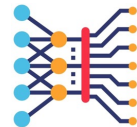


$\delta_j^{(l)}$ = “error” of node j in layer l

Formally,
$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$$

where $\text{cost}(\mathbf{x}_i) = y_i \log h_{\Theta}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\Theta}(\mathbf{x}_i))$

Backpropagation

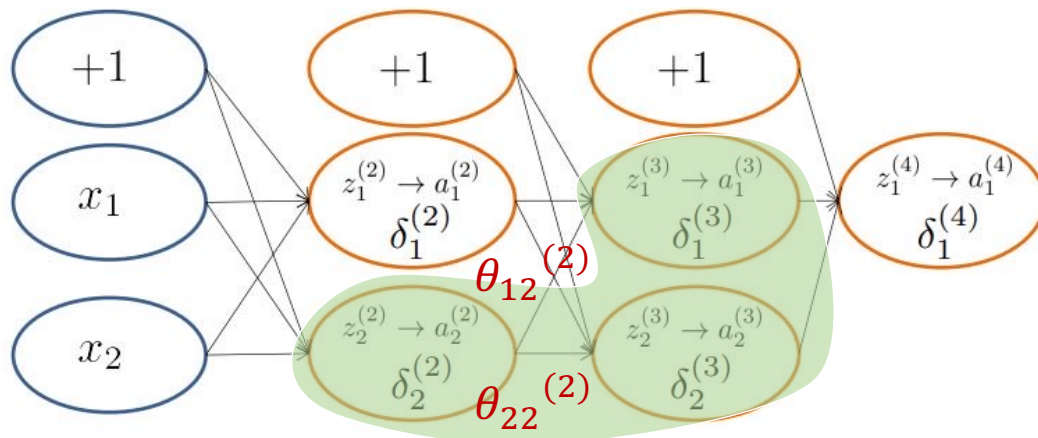
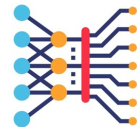


$\delta_j^{(l)}$ = “error” of node j in layer l

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = y_i \log h_{\Theta}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\Theta}(\mathbf{x}_i))$

Backpropagation



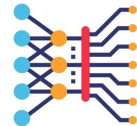
$$\delta_2^{(2)} = \theta_{12}^{(2)} \times \delta_1^{(3)} \times a_2'^{(2)} + \theta_{22}^{(2)} \times \delta_2^{(3)} \times a_2'^{(2)}$$

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

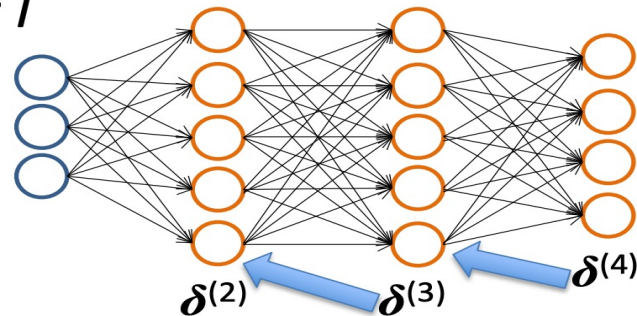
Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = y_i \log h_{\Theta}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\Theta}(\mathbf{x}_i))$

Gradient Computation



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



Backpropagation

- $\delta^{(4)} = \mathbf{a}^{(4)} - \mathbf{y}$
- $\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \text{ .* } g'(\mathbf{z}^{(3)})$
- $\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \text{ .* } g'(\mathbf{z}^{(2)})$
- (No $\delta^{(1)}$)

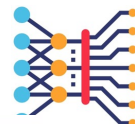
Element-wise
product .*

$$g'(\mathbf{z}^{(3)}) = \mathbf{a}^{(3)} \text{ .* } (1 - \mathbf{a}^{(3)})$$

$$g'(\mathbf{z}^{(2)}) = \mathbf{a}^{(2)} \text{ .* } (1 - \mathbf{a}^{(2)})$$

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

Backpropagation



Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$

(Used to accumulate gradient)

For each training instance (\mathbf{x}_i, y_i) :

Set $\mathbf{a}^{(1)} = \mathbf{x}_i$

Compute $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\delta^{(L)} = \mathbf{a}^{(L)} - y_i$

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

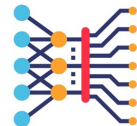
Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

$\mathbf{D}^{(l)}$ is the matrix of partial derivatives of $J(\Theta)$

Note: Can vectorize $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ as $\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} \mathbf{a}^{(l)\top}$

Training a NN via Gradient Descent with Backprop



Given: training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)

Loop // each iteration is called an epoch

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$ (Used to accumulate gradient)

For each training instance (\mathbf{x}_i, y_i) :

Set $\mathbf{a}^{(1)} = \mathbf{x}_i$

Compute $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\delta^{(L)} = \mathbf{a}^{(L)} - y_i$

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

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

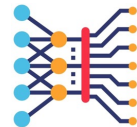
Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

Until weights converge or max #epochs is reached

Backpropagation

Training a Neural Network



- Pick a network architecture (connectivity pattern between nodes)
- # input units = # of features in dataset
- # output units = # classes
- Choose number of hidden layers and number of units in each layer.

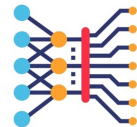
Gradient descent

Partial derivatives give us the slope (i.e. direction to move) in that dimension

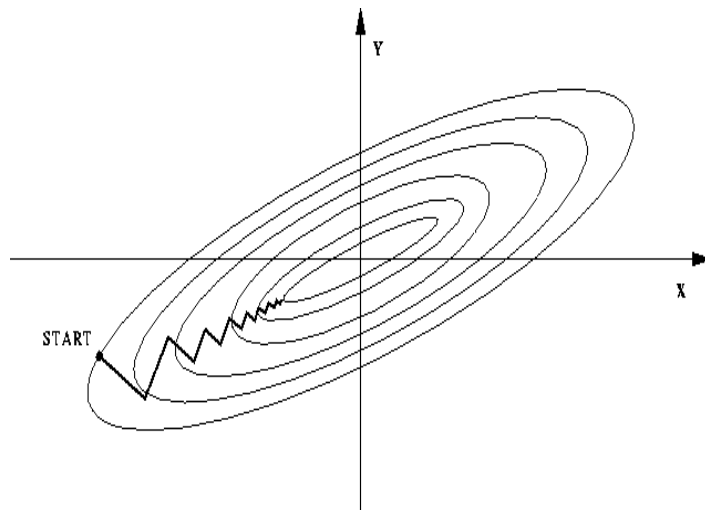
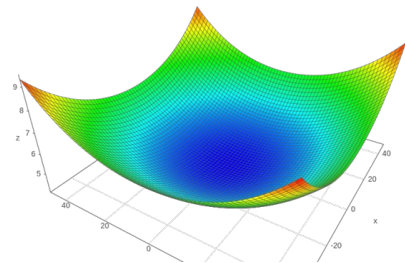
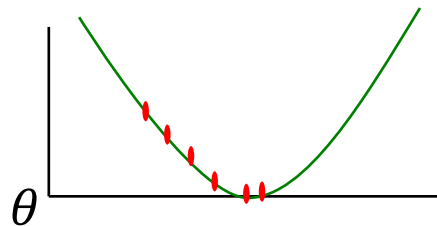
Approach:

- pick a starting point (θ)
- repeat:
 - pick a dimension
 - move a small amount in that dimension towards decreasing loss (using the derivative)

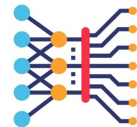
$$\theta_j = \theta_j - \eta \frac{d}{d\theta_j} \text{Loss}(\theta)$$



Loss

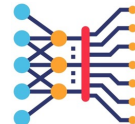


Batch, Stochastic and Minibatch



- Optimization algorithms that use the entire training set to compute the gradient are called batch or deterministic gradient methods. Ones that use a single training example for that task are called stochastic or online gradient methods
- Most of the algorithms we use for deep learning fall somewhere in between!
- These are called minibatch or minibatch stochastic methods

Batch, Stochastic and Mini-batch Stochastic Gradient Descent



Algorithm 1 Batch Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

- 1: **while** stopping criteria not met **do**
- 2: Compute gradient estimate over N examples:
- 3: $\hat{\mathbf{g}} \leftarrow +\frac{1}{N} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
- 4: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
- 5: **end while**

Algorithm 2 Stochastic Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

- 1: **while** stopping criteria not met **do**
- 2: Sample example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ from training set
- 3: Compute gradient estimate:
- 4: $\hat{\mathbf{g}} \leftarrow +\nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
- 5: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
- 6: **end while**

Mini-batch

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

Batch and Stochastic Gradient Descent

