



CS60010: Deep Learning

Spring 2023

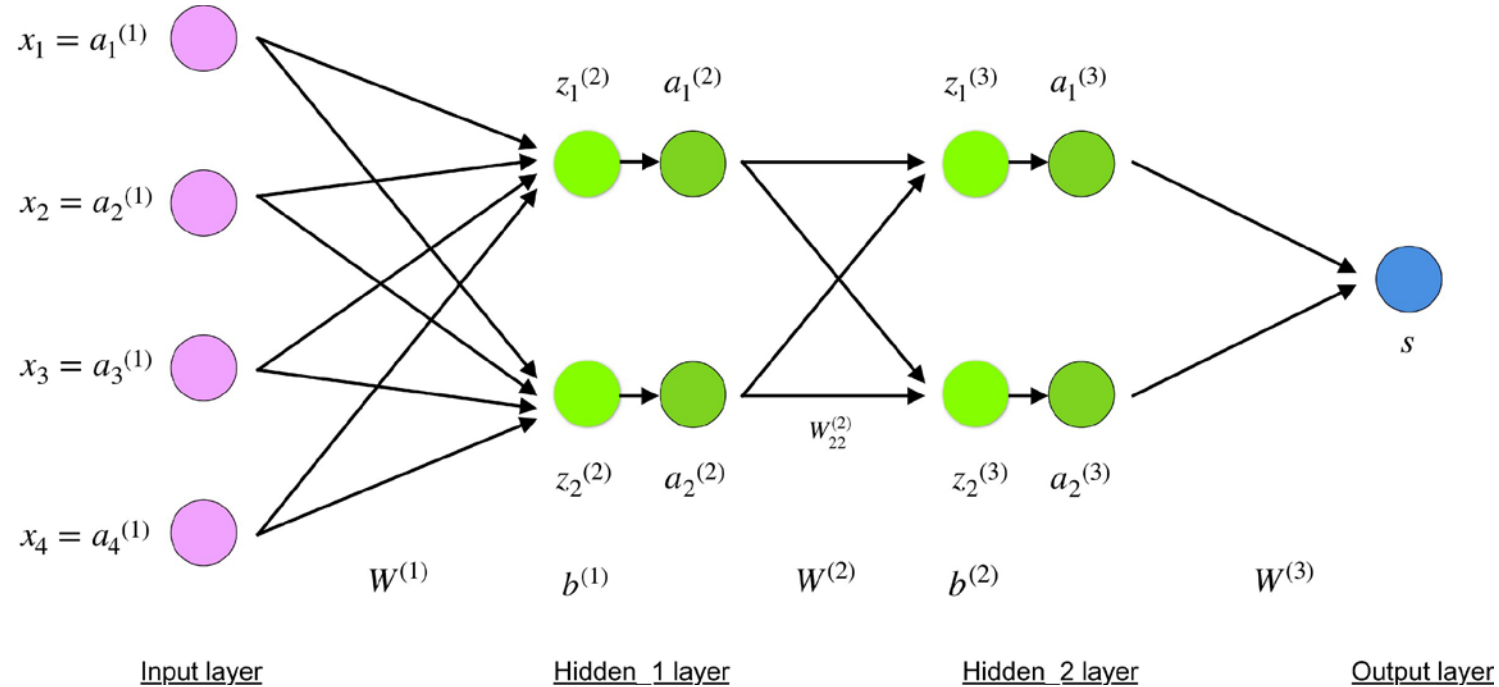
Sudeshna Sarkar

Multilayer Perceptron – Part 2

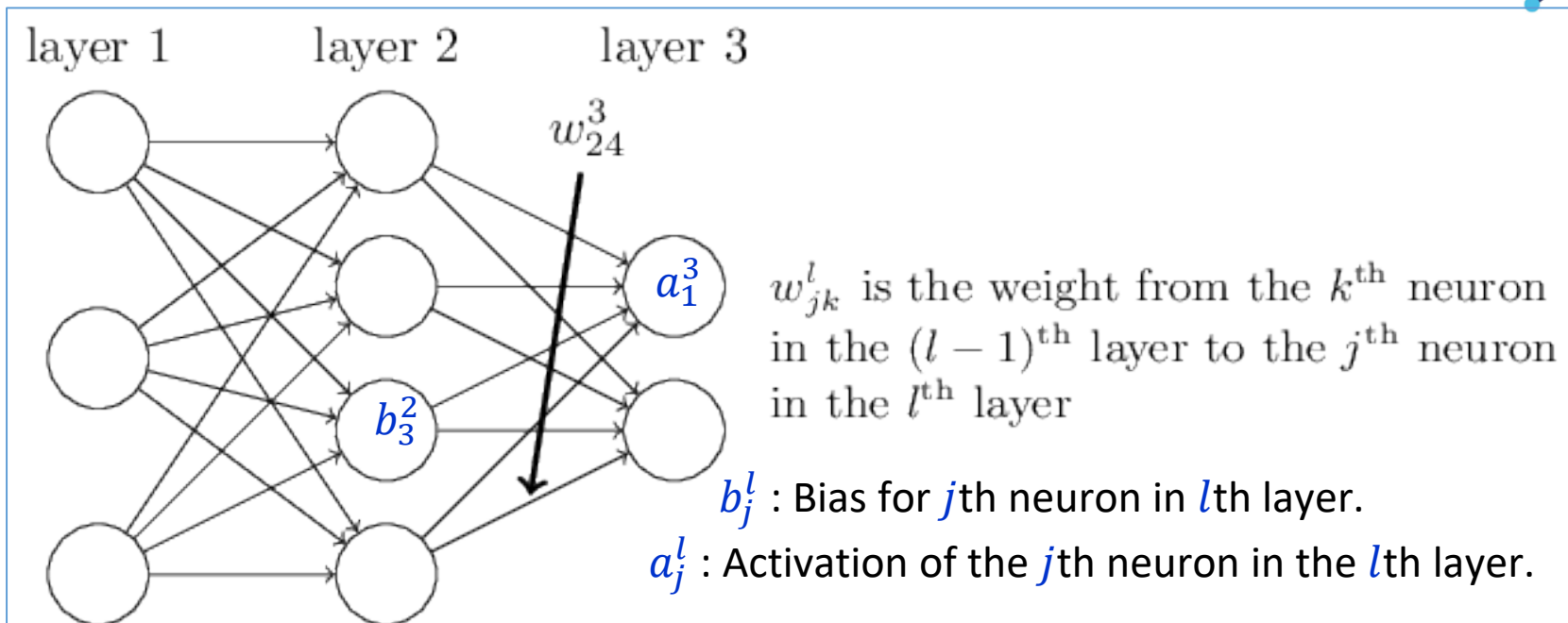
25 Jan 2023



Multilayer Neural Network



Notations



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

Vectorized form: $a^l = g(w^l a^{l-1} + b^l)$

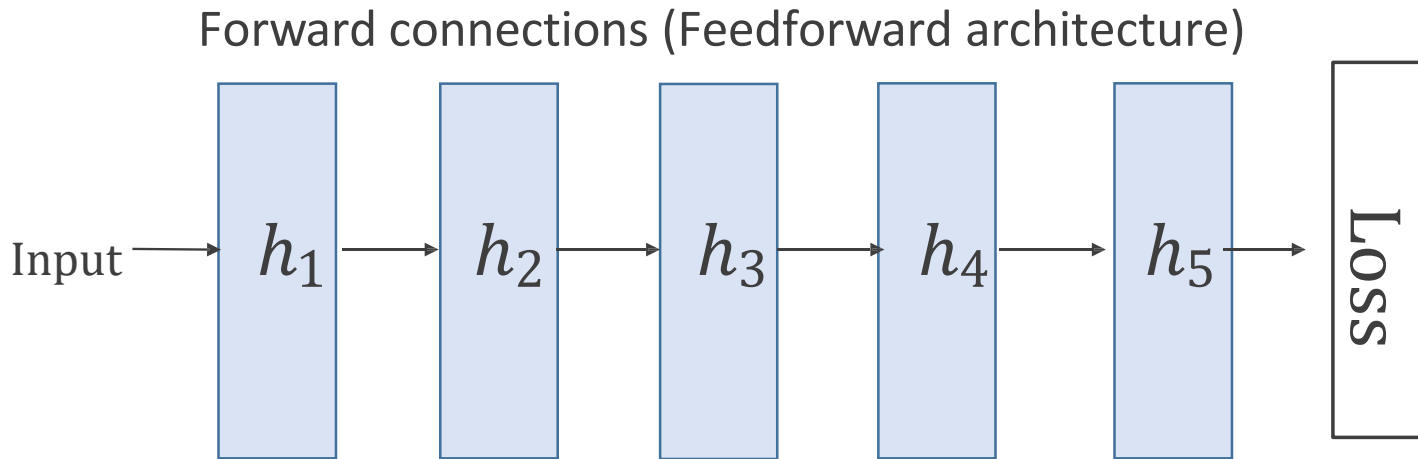
$$z^l = w^l a^{l-1} + b^l$$

$$a^l = g(z^l)$$

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

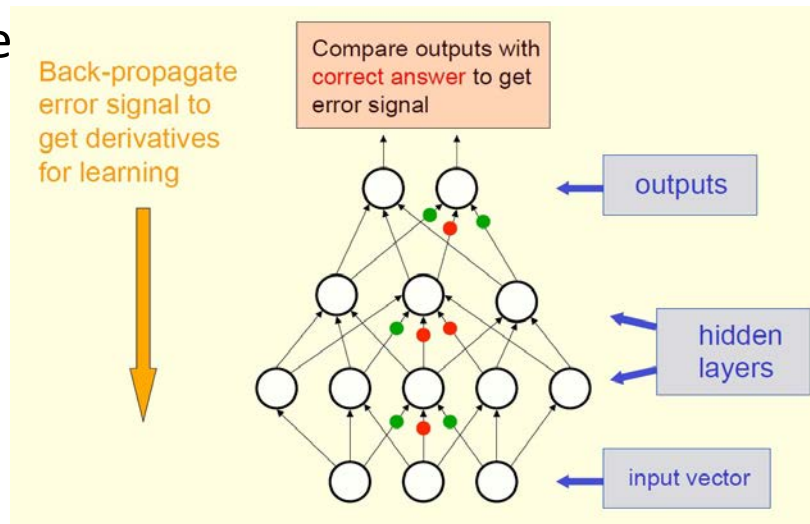


- Feedforward Propagation: Accept input $x^{(i)}$, pass through intermediate stages and obtain output $\hat{y}^{(i)}$

- During Training: Compute scalar cost $J(\theta)$

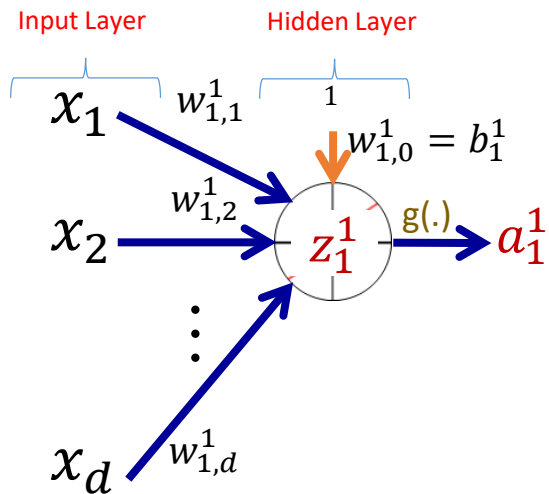
$$J(\theta) = \sum_i L(NN(x^{(i)}; \theta), y^{(i)})$$

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





Multilayer Neural Network

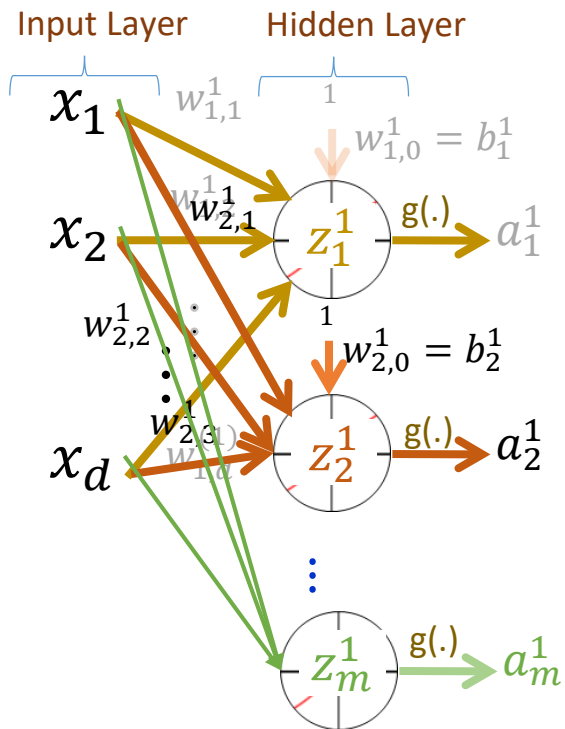


$$z_1^1 = b_1^1 + \sum_{i=1}^d w_{1,i}^1 x_i = [\mathbf{w}_1^1 b_1^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

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

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

Multilayer Neural Network

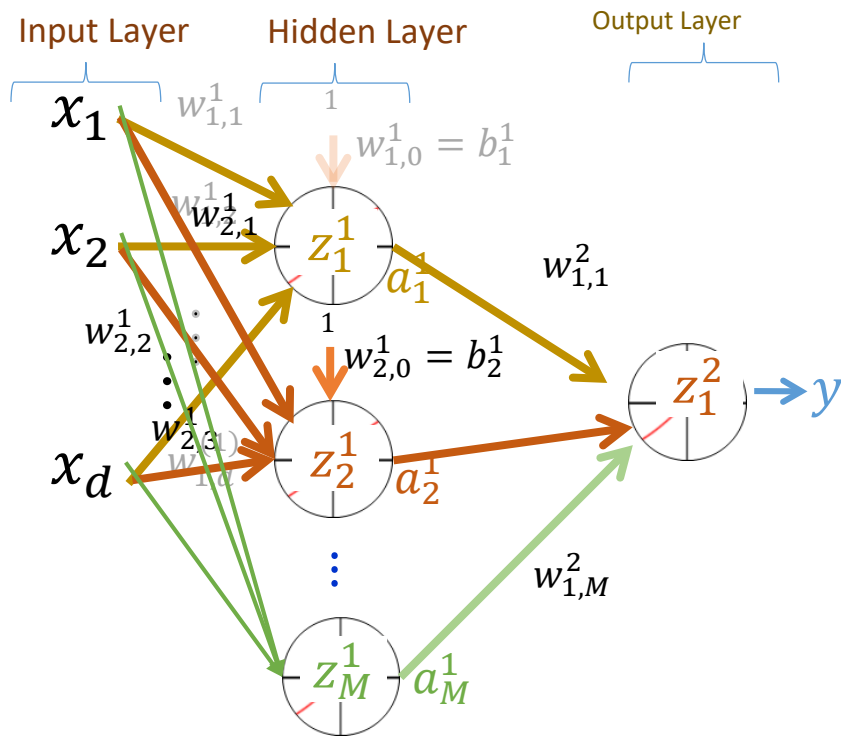


$$\left. \begin{aligned} z_1^1 &= [\mathbf{w}_1^1 b_1^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \\ z_2^1 &= [\mathbf{w}_2^1 b_2^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \\ &\vdots \\ z_m^1 &= [\mathbf{w}_m^1 b_m^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \end{aligned} \right\} \underbrace{\begin{bmatrix} z_1^1 \\ z_2^1 \\ \vdots \\ z_m^1 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^1 & b_1^1 \\ \mathbf{w}_2^1 & b_2^1 \\ \vdots & \vdots \\ \mathbf{w}_m^1 & b_m^1 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}}_{\mathbf{z}^1 = [\mathbf{W}^1 \mathbf{b}^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}}$$

$$\left. \begin{bmatrix} a_1^1 \\ a_2^1 \\ \vdots \\ a_m^1 \end{bmatrix} = \begin{bmatrix} g(z_1^1) \\ g(z_2^1) \\ \vdots \\ g(z_m^1) \end{bmatrix} \right\} \mathbf{a}^1 = \mathbf{g}(\mathbf{z}^{(1)})$$

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

Multilayer Neural Network



Output Layer Pre-activation

$$z_1^2 = [\mathbf{w}_1^2 \ b_1^2] \begin{bmatrix} \mathbf{a}^1 \\ 1 \end{bmatrix}$$

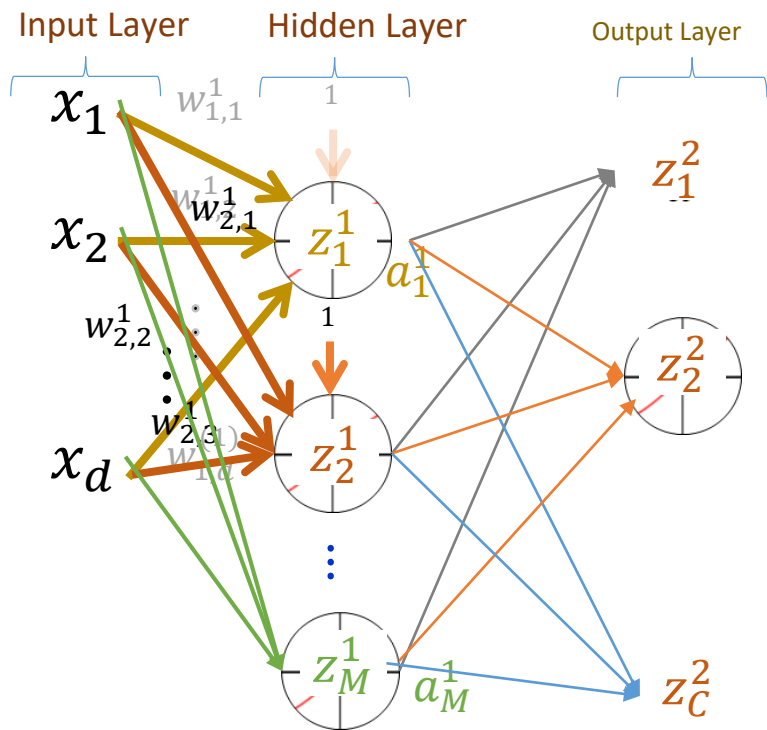
Output Layer Activation

$$y_1 = o(\mathbf{z}_1^2)$$

output

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

Multilayer Neural Network

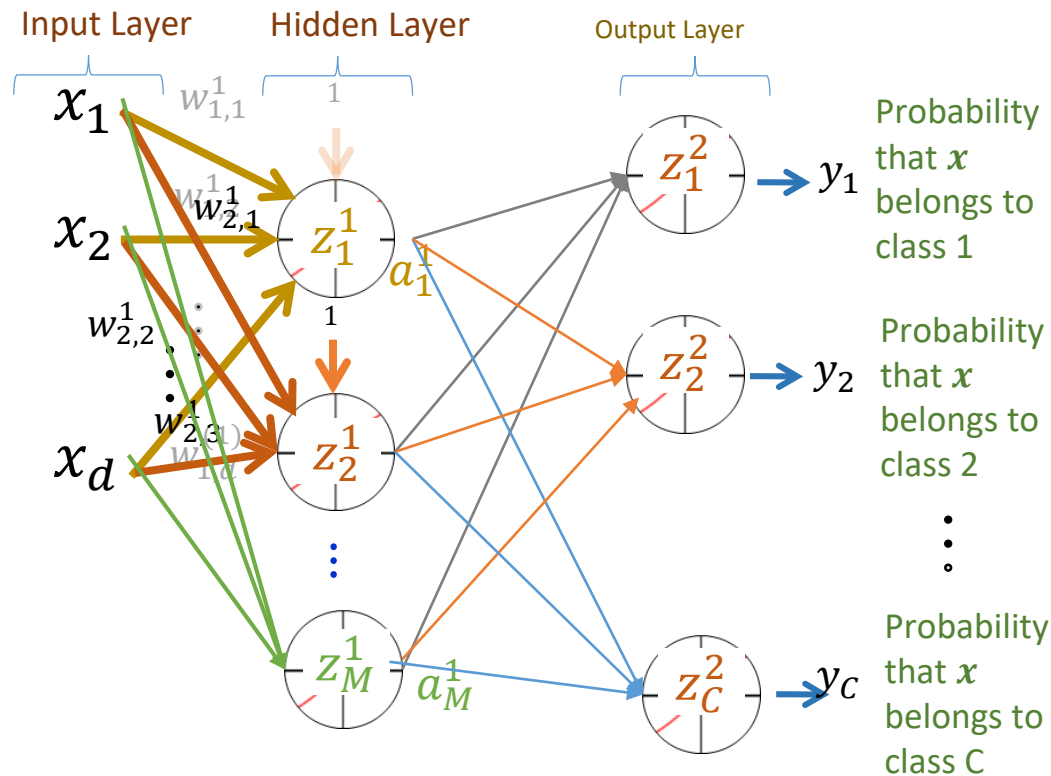


$$\rightarrow y_1 = o_1(z_1^2) = \frac{\exp(z_1^2)}{\sum_c \exp(z_c^2)}$$

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

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

Training a Neural Network – Loss Function



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

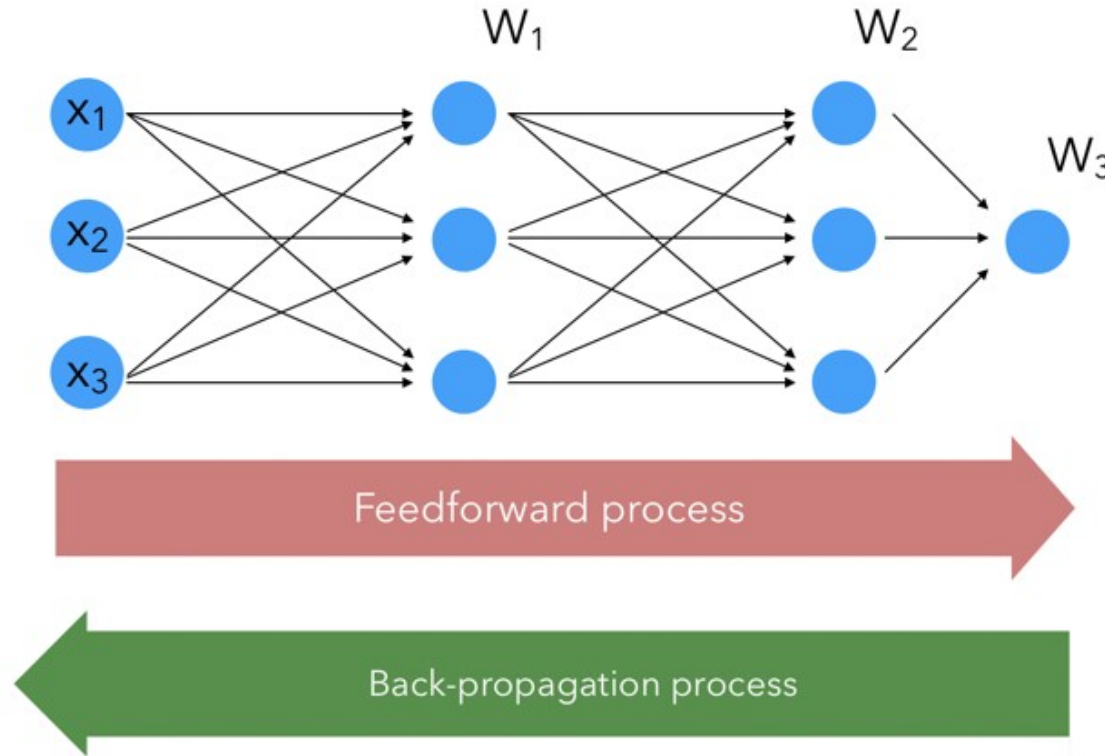
$$\begin{aligned} \max y_c \\ \equiv \max (\log y_c) \\ \equiv \min (-\log y_c) \end{aligned}$$

Can be equivalently expressed as

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

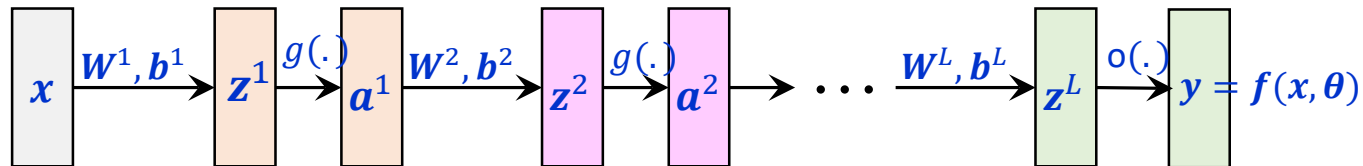
known as cross-entropy loss

Forward-Backward Passes



Forward Pass

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



Hidden layer pre-activation:

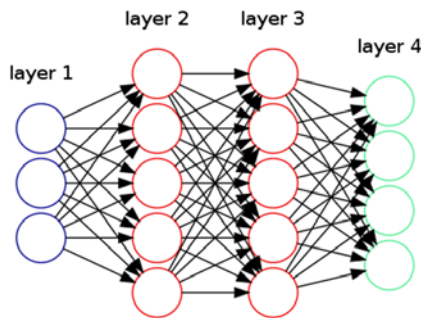
For $l = 1, \dots, L$; $\mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}$

Hidden layer activation:

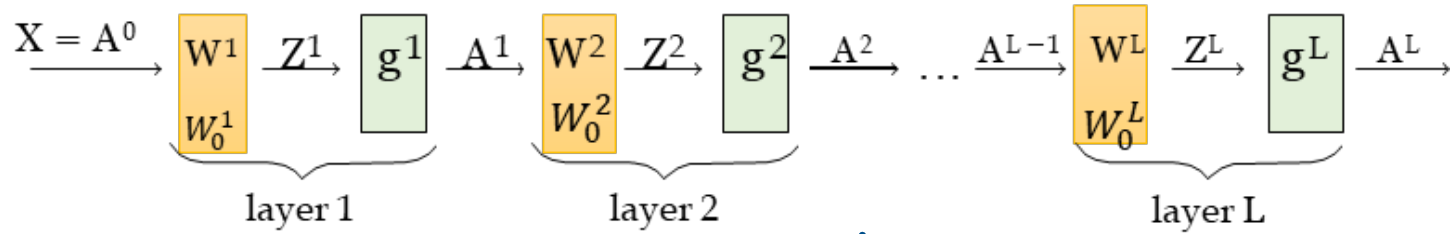
For $l = 1, \dots, L - 1$; $\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)})$

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}

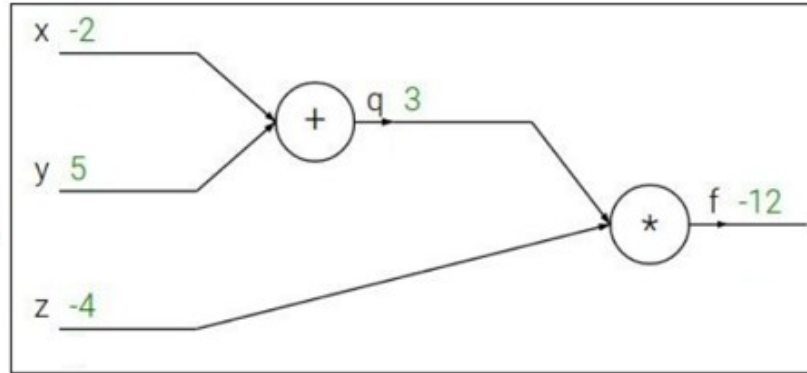
....



Consider following computational graph:

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

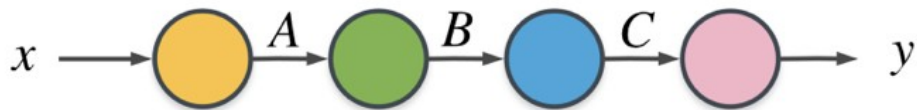


$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

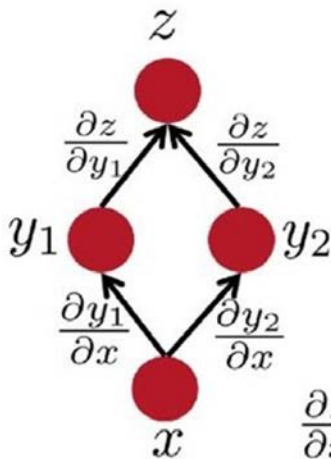
Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

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.

Computation Graph and Backpropagation



$$f(x, y, z) = (x + y)z$$

$$x = -2, y = 5, z = -4$$

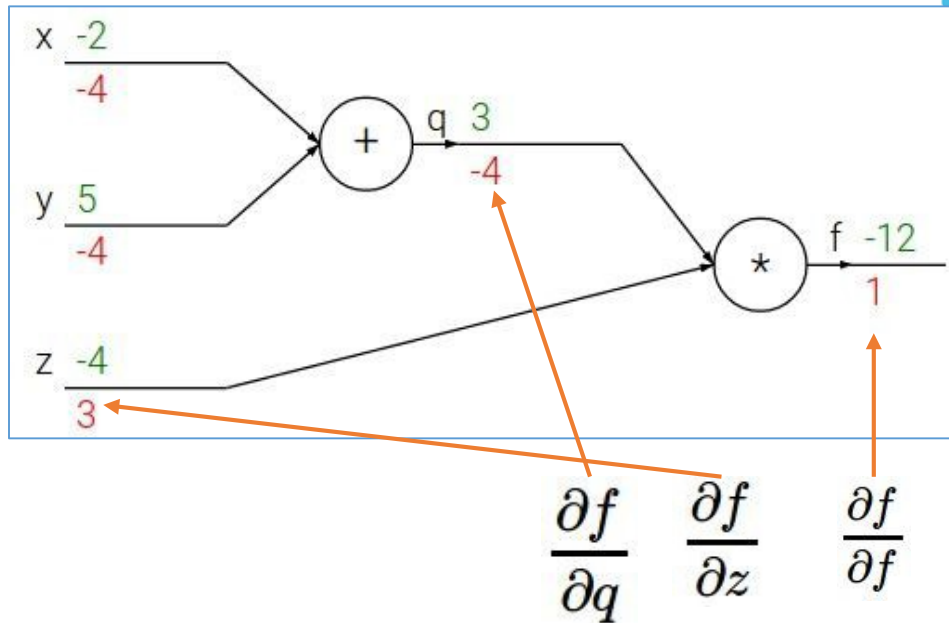
$$q = x + y$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Find $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Computation Graph and Backpropagation



$$f(x, y, z) = (x + y)z$$

$$x = -2, y = 5, z = -4$$

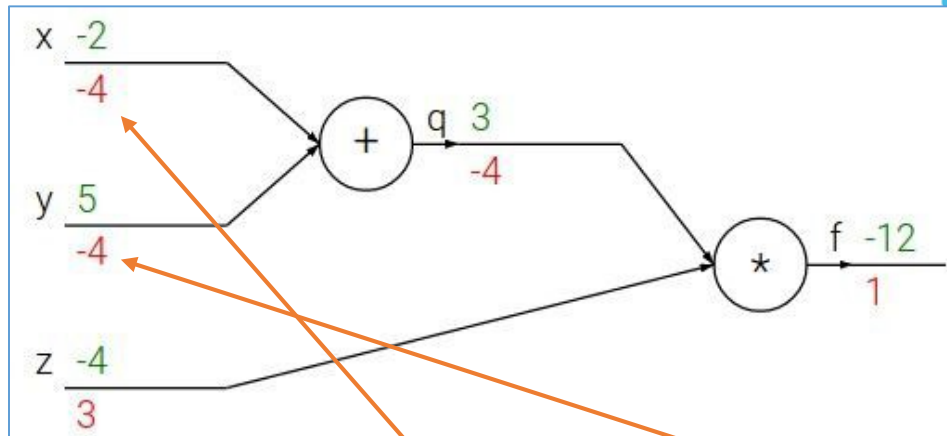
$$q = x + y$$

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Find $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

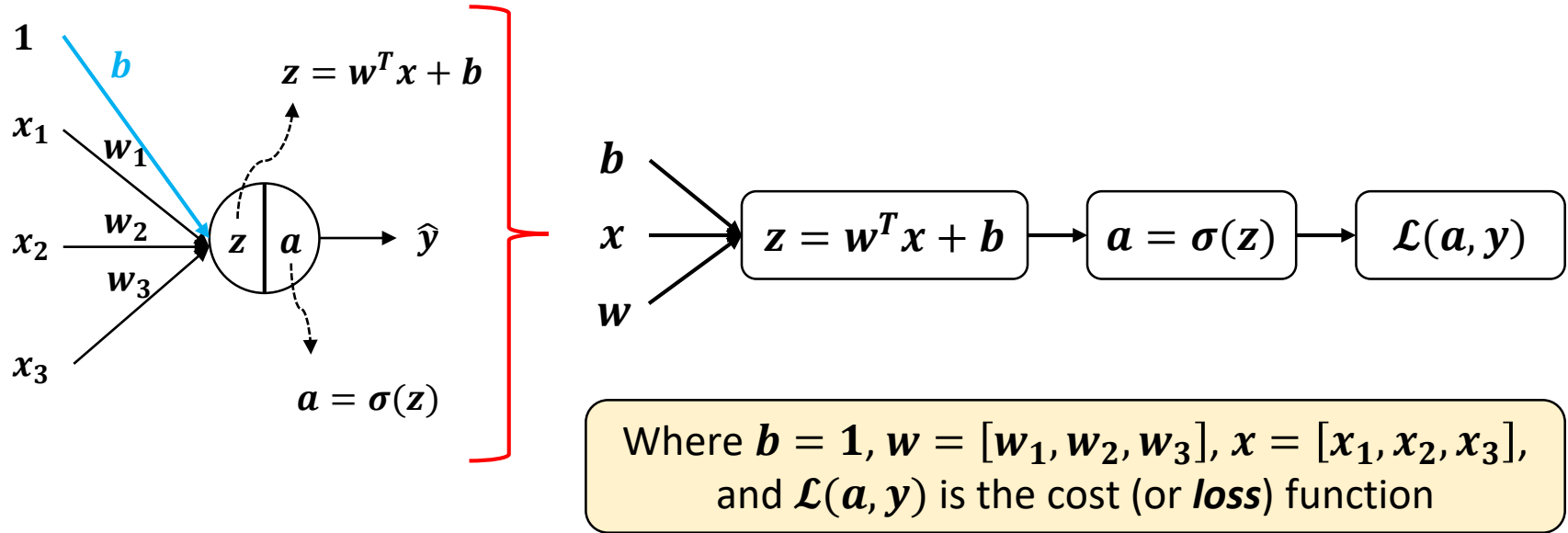
Upstream
gradient

Local
gradient

The Computation Graph of Logistic Regression



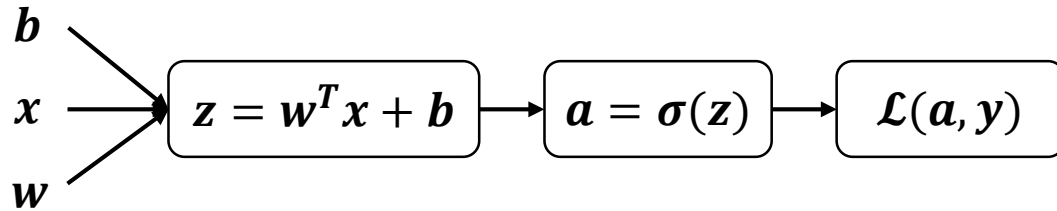
Let us translate logistic regression into a computation graph



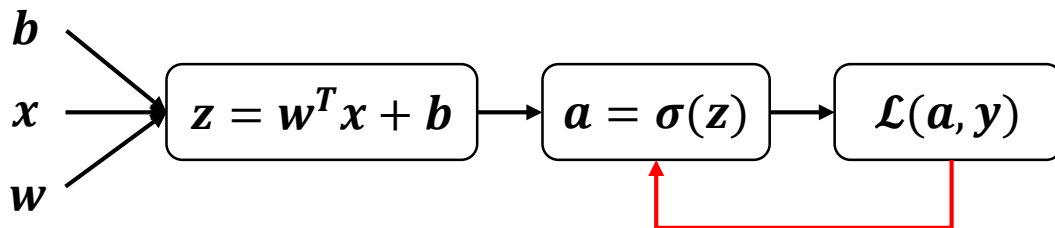
Forward Propagation



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



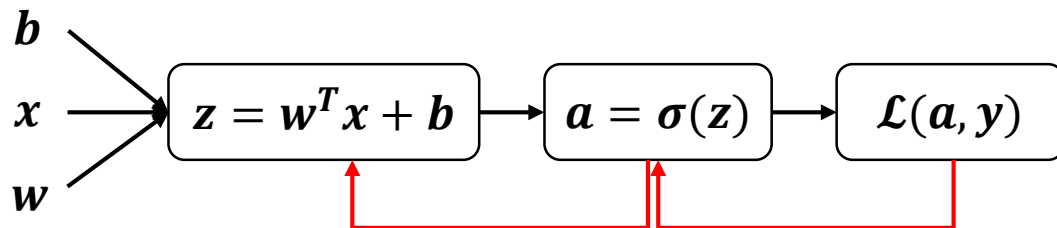
Backward Propagation



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

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

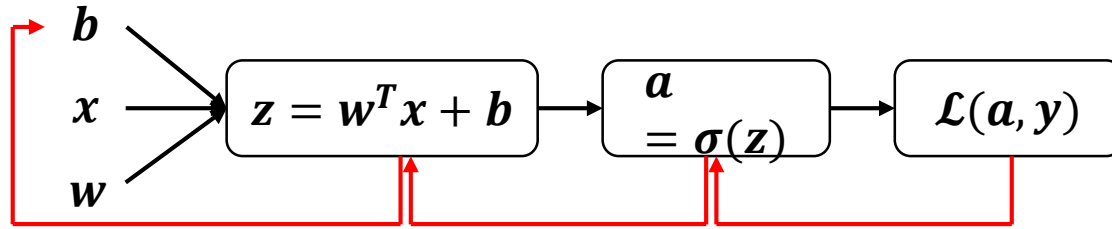
Backward Propagation



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

$= a - y$

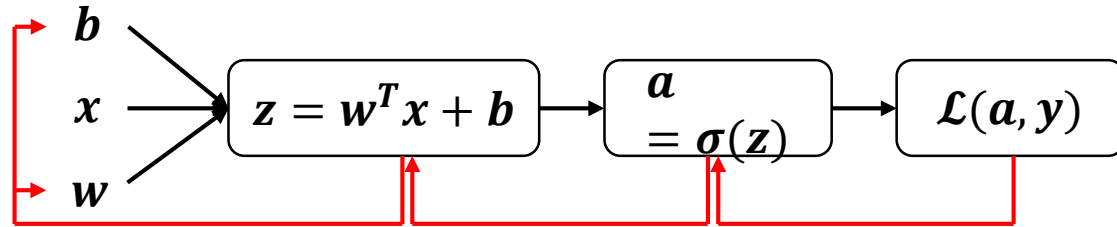
Backward Propagation



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

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

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:

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

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

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

Gradient Descent For Logistic Regression



Outline:

- Have a loss function $\mathcal{L}(\mathbf{w}, b)$
- Start off with some guesses for $\mathbf{w}_1, \dots, \mathbf{w}_m$
- Repeat until convergence{

$$\left. \begin{aligned} w_j &= w_j - \eta \frac{\partial \mathcal{L}(\mathbf{w}, b)}{\partial w_j} \\ b &= b - \eta \frac{\partial \mathcal{L}(\mathbf{w}, b)}{\partial b} \end{aligned} \right\}$$

Gradient Descent For Logistic Regression



- **Outline:**

- Repeat until convergence{

Assuming n examples

for $i = 1$ to n :

Forward
propagation

$$\begin{cases} z^{(i)} = w^T x^{(i)} + b \\ a^{(i)} = \sigma(z^{(i)}) \end{cases}$$

Backward
propagation

$$\begin{cases} dz^{(i)} = a^{(i)} - y^{(i)} \\ dw = dw + dz^{(i)} x^{(i)} \\ db = db + dz^{(i)} \end{cases}$$

Outside
the loop

$$\begin{cases} dw = dw/n \\ db = db/n \\ w = w - \eta dw \\ b = b - \eta db \end{cases}$$

}

$$Z = w^T X + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{n} X dZ^T$$

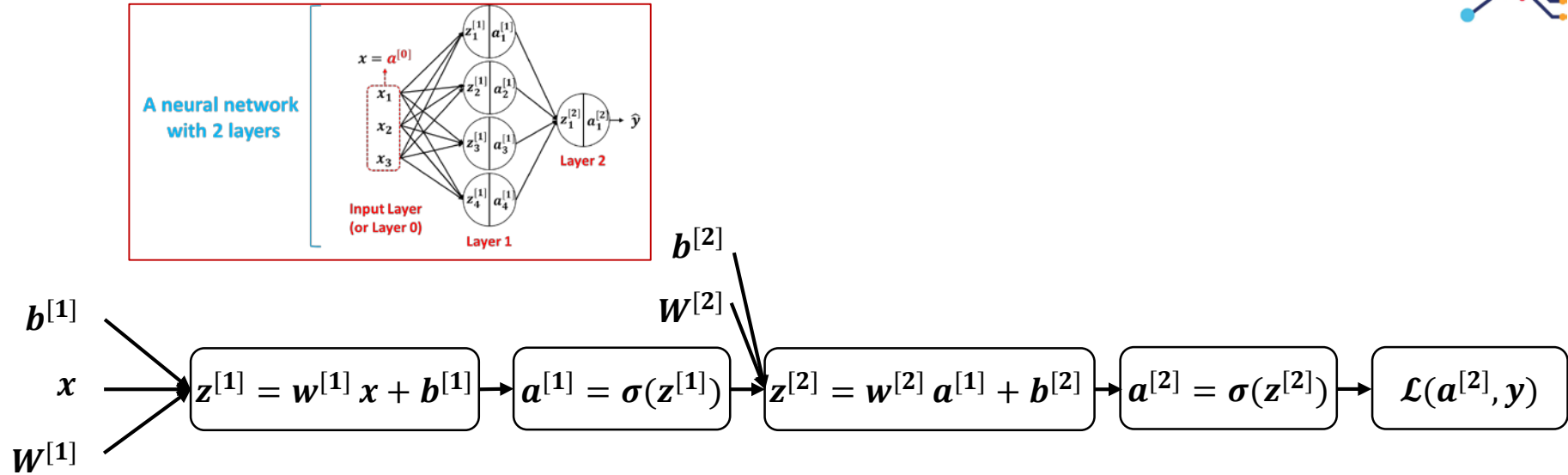
$$db = \frac{1}{n} \sum_{i=1}^n dz^{(i)}$$

$$w = w - \eta dw$$

$$b = b - \eta db$$

Vectorized version

Computation Graph of a 2-layer Neural Network

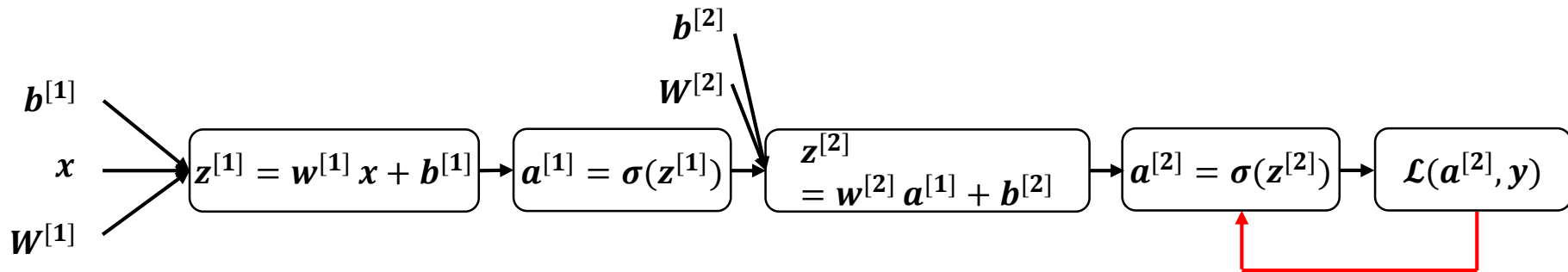


The loss function can be computed by moving from left to right

Backward Propagation



The derivatives can be computed by moving from right to left

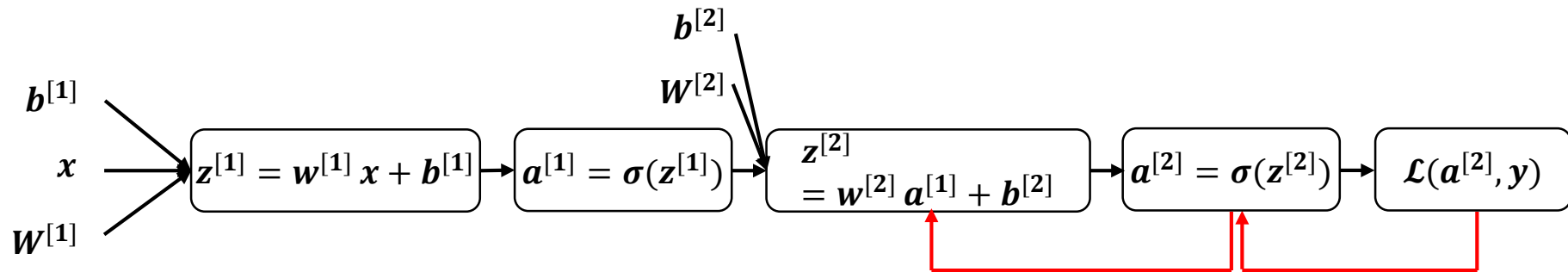


$$\frac{\partial \mathcal{L}}{\partial a^{[2]}} = \frac{-y}{a^{[2]}} + \frac{(1-y)}{(1-a^{[2]})}$$

Backward Propagation

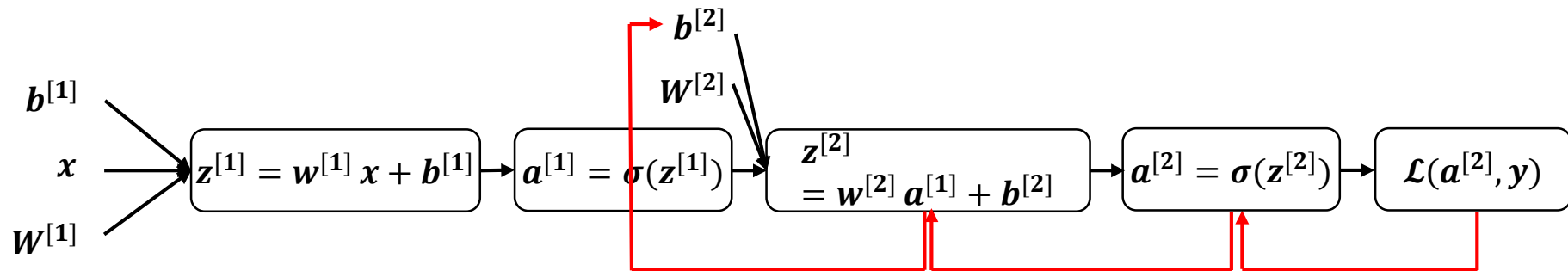


- The derivatives can be computed by moving from right to left



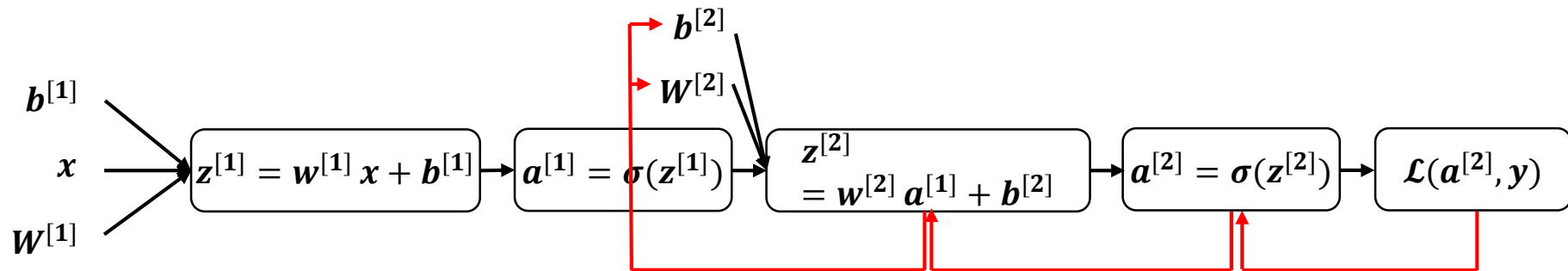
$$\frac{\partial \mathcal{L}}{\partial z^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} = a^{[2]} - y$$

Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial b^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial b^{[2]}} = a^{[2]} - y$$

Backward Propagation

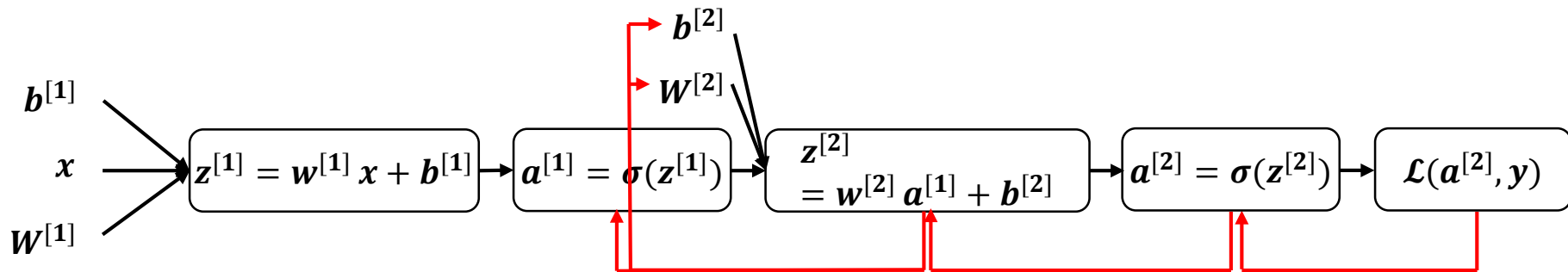


$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial W^{[2]}} = (a^{[2]} - y)a^{[1]T}$$

Backward Propagation



- The derivatives can be computed by moving from right to left

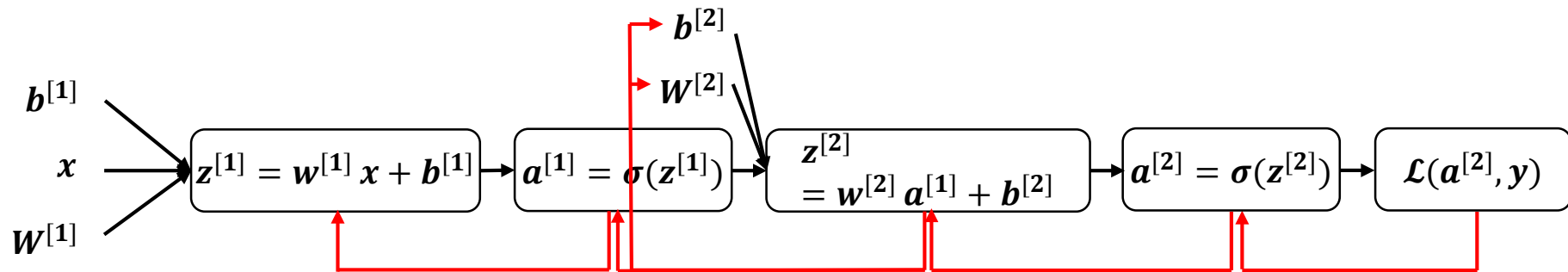


$$\frac{\partial \mathcal{L}}{\partial a^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} = (a^{[2]} - y)w^{[2]T}$$

Backward Propagation



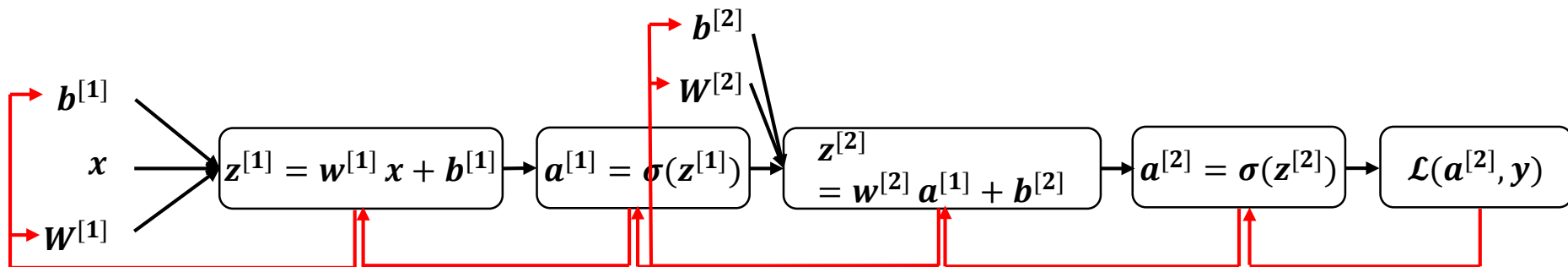
- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial z^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} = (a^{[2]} - y) w^{[2]T} * a^{[1]} (1 - a^{[1]})$$

 Element-wise product

Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} \times \frac{\partial z^{[1]}}{\partial W^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \left((a^{[2]} - y) w^{[2]T} * a^{[1]} (1 - a^{[1]}) \right) x^T$$

$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} \times \frac{\partial z^{[1]}}{\partial b^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = (a^{[2]} - y) w^{[2]T} * a^{[1]} (1 - a^{[1]})$$

Backward Propagation: Summary



$$dz^{[2]} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} = a^{[2]} - y$$

$$db^{[2]} = \frac{\partial \mathcal{L}}{\partial b^{[2]}} = a^{[2]} - y$$

$$dW^{[2]} = \frac{\partial \mathcal{L}}{\partial W^{[2]}} = (a^{[2]} - y) a^{[1]T}$$

$$dz^{[1]} = \frac{\partial \mathcal{L}}{\partial z^{[1]}} = dz^{[2]} w^{[2]T} * a^{[1]} (1 - a^{[1]})$$

$$db^{[1]} = \frac{\partial \mathcal{L}}{\partial b^{[1]}} = dz^{[2]} w^{[2]T} * a^{[1]} (1 - a^{[1]})$$

$$dW^{[1]} = \frac{\partial \mathcal{L}}{\partial W^{[1]}} = \left(dz^{[2]} w^{[2]T} * a^{[1]} (1 - a^{[1]}) \right) x^T$$

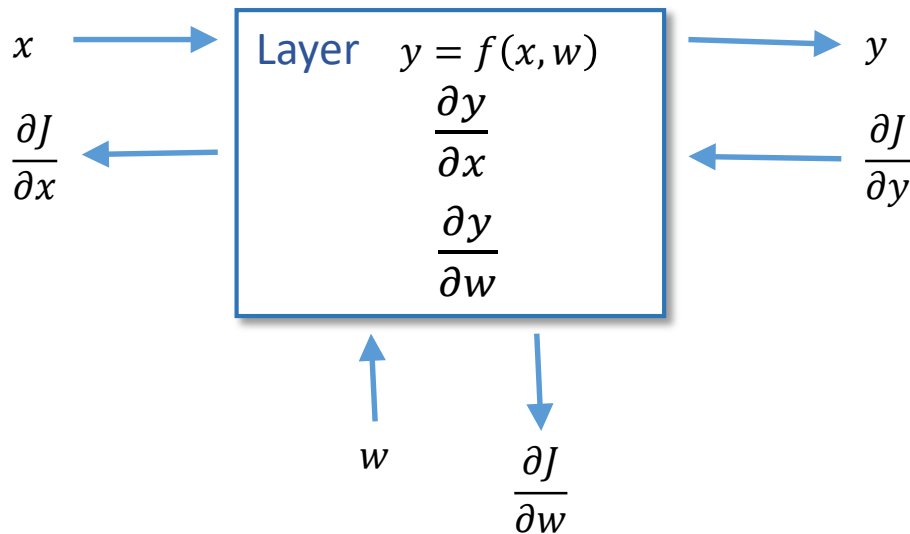
Backpropagation



- Compute derivatives per layer, utilizing previous derivatives
- Objective: $J(\mathbf{w})$
- Arbitrary layer: $y = f(x, w)$
- Need:

- $\frac{\partial J}{\partial x} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x}$

- $\frac{\partial J}{\partial w} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w}$



Backpropagation (layerwise)



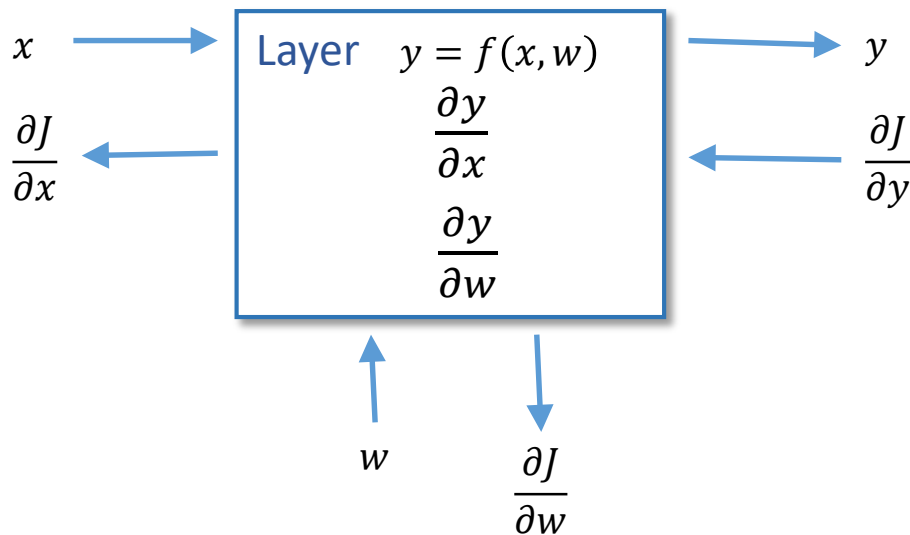
- Compute derivatives per layer, utilizing previous derivatives
- Objective: $J(\mathbf{w})$
- Arbitrary layer: $y = f(x, w)$

- Init:

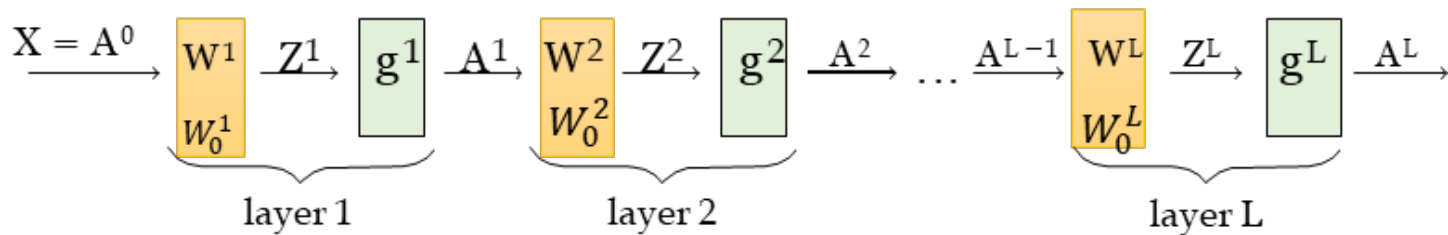
- $\frac{\partial J}{\partial x} = 0$
 - $\frac{\partial J}{\partial w} = 0$

- Compute:

- $\frac{\partial J}{\partial x} \text{ ~~+~~ } = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x}$
 - $\frac{\partial J}{\partial w} \text{ ~~+~~ } = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w}$



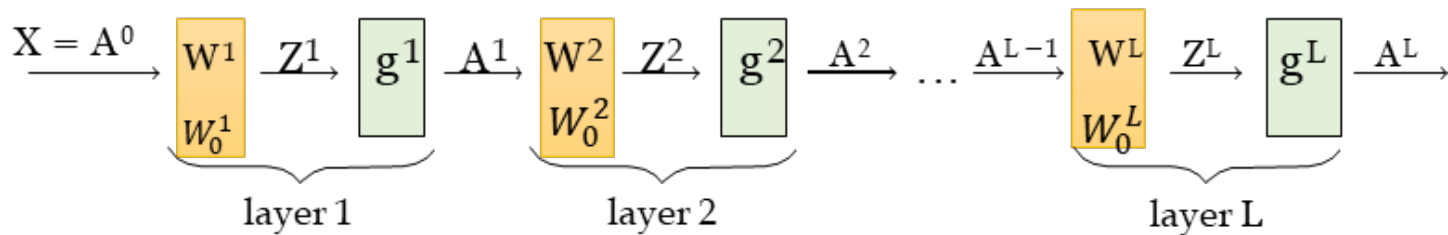
Error back-propagation



To do SGD for a training example (x, y) , we need to compute

$$\nabla_W \text{Loss}(NN(x; W), y)$$
$$\frac{\partial L}{\partial W^L} = \frac{\partial L}{\partial A^L} \cdot \frac{\partial A^L}{\partial Z^L} \cdot \frac{\partial Z^L}{\partial W^L}$$

Error back-propagation



$$\frac{\partial \text{Loss}}{\partial W^L} = \frac{\partial \text{Loss}}{\partial A^L} \cdot \frac{\partial A^L}{\partial Z^L} \cdot \frac{\partial Z^L}{\partial W^L}$$

$$\frac{\partial L}{\partial W^l} = A^{l-1} \left(\frac{\partial L}{\partial Z^l} \right)^T$$

$m^l \times n^l \quad m^l \times 1 \quad 1 \times n^l$

So, in order to find the gradient of the loss with respect to the weights in the other layers of the network,

we just need to be able to find $\frac{\partial \text{Loss}}{\partial Z^l}$

Gradient descent

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

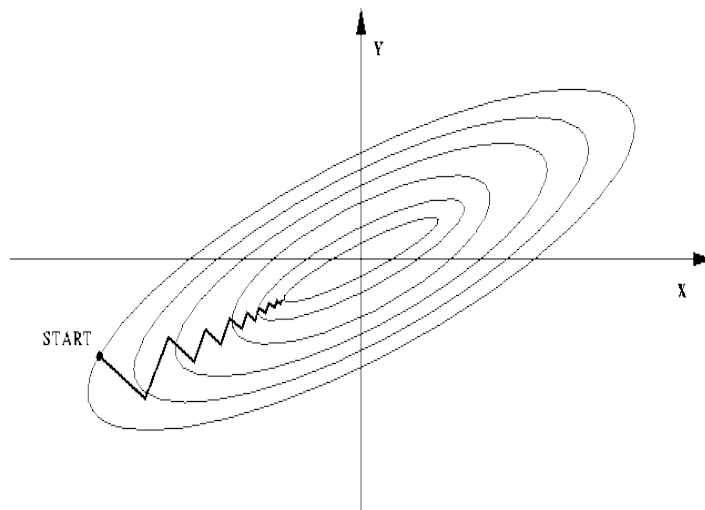
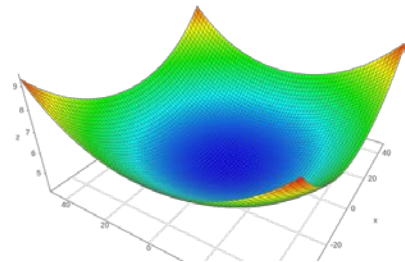
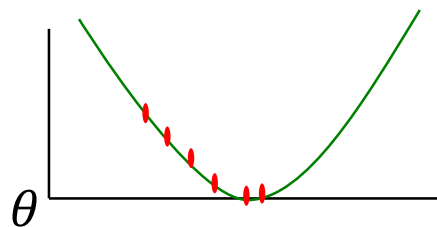
Approach:

- 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

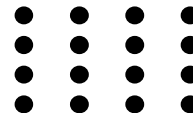
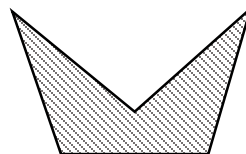
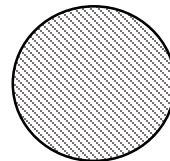
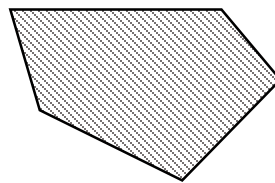
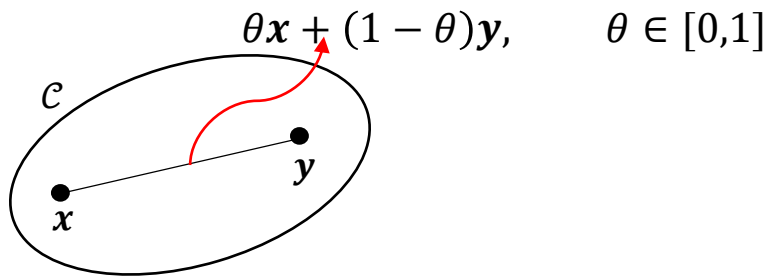


Aside - Convex Sets and Functions



- **Convex Set:** A set $\mathcal{C} \subseteq \mathbb{R}^n$ is a convex set if for all $x, y \in \mathcal{C}$, the line segment connecting x and y is in \mathcal{C} , i.e.,

$$\forall x, y \in \mathcal{C}, \theta \in [0, 1], \theta x + (1 - \theta)y \in \mathcal{C}$$

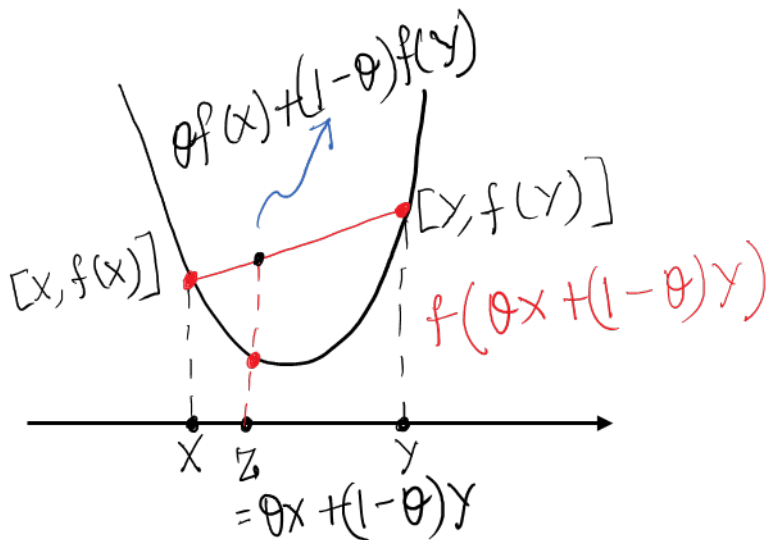


Convex Sets and Functions



- **Convex Function:** A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function if its domain $dom(f)$ is a convex set and for all $x, y \in dom(f)$, and $\theta \in [0,1]$, we have

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$



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

