

# CAP 4453

# Robot Vision

Dr. Gonzalo Vaca-Castaño  
[gonzalo.vacacastano@ucf.edu](mailto:gonzalo.vacacastano@ucf.edu)



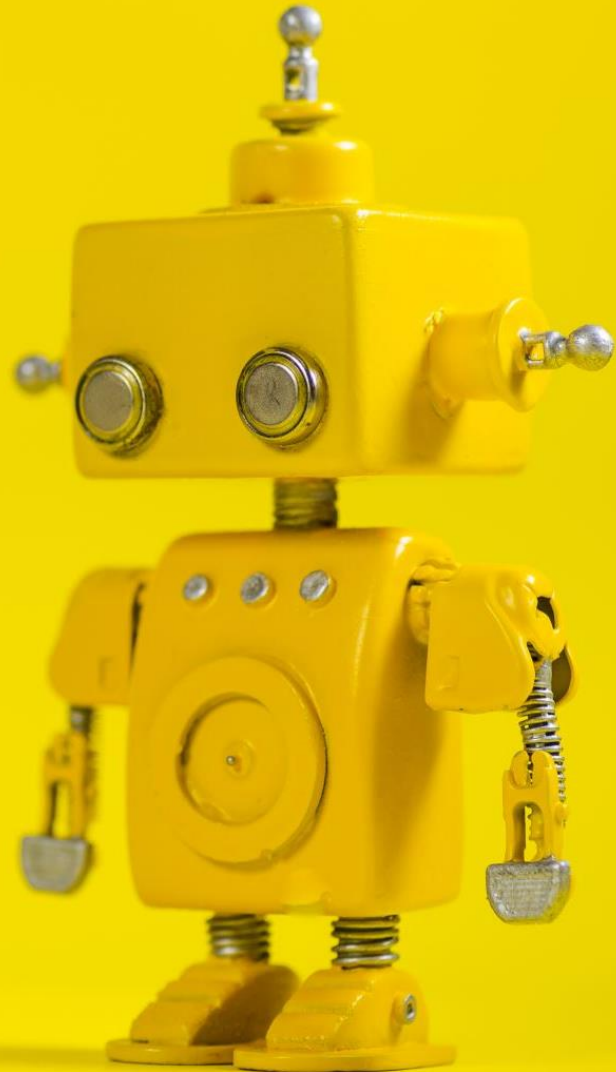
# Administrative details

- Correction of the midterm exam



# Credits

- Some slides comes directly from:
  - Yosesh Rawat
  - Andrew Ng



# Robot Vision

## 17. Introduction to Deep Learning II



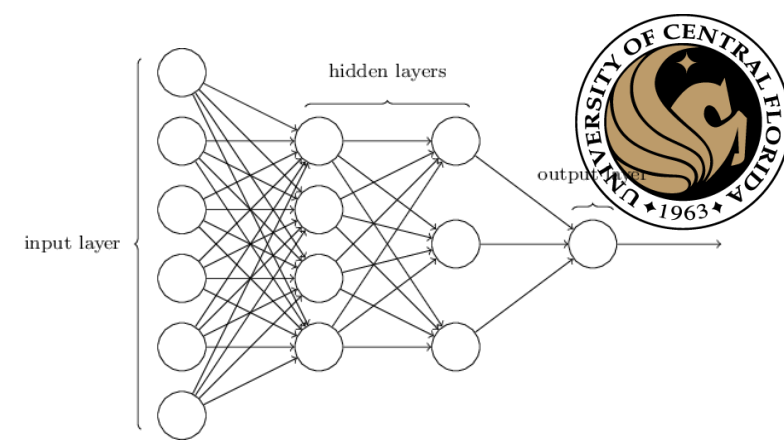
# Outline

- Fully connected Neural network
  - Activation functions:
    - Forward and backward
  - Back propagation
  - Network definitions
  - Initialization
  - Training
    - Hyper parameters
    - Gradient updates: RMS prop,
    - Amount of training data
    - Batch normalization
  - Dataset
    - Train set, test set, validation set
    - Bias and variance
- Implementation network to solve digit identification

# Fully connected networks: The math

## A REVIEW

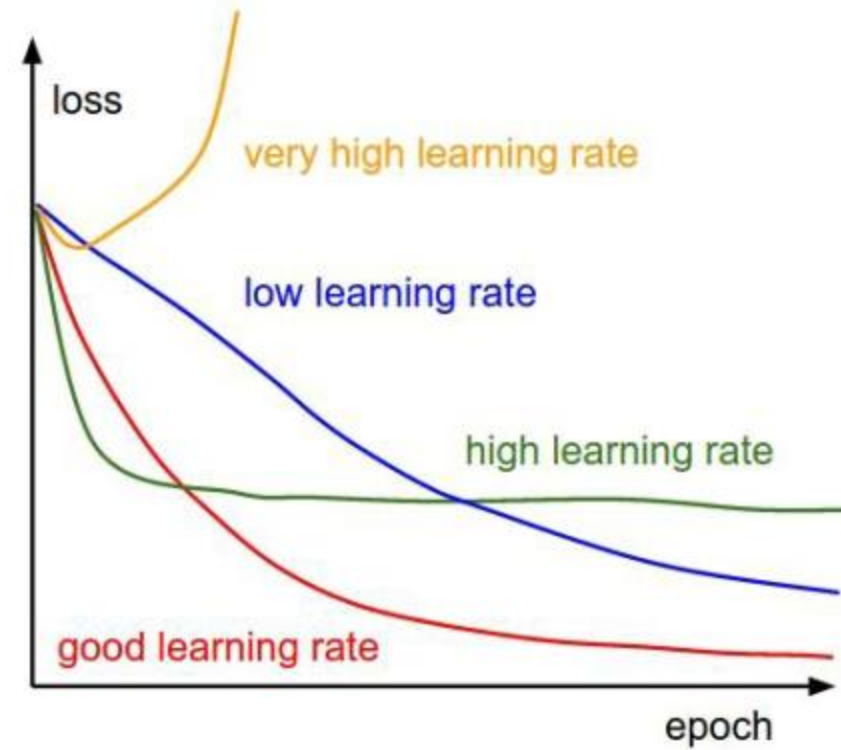
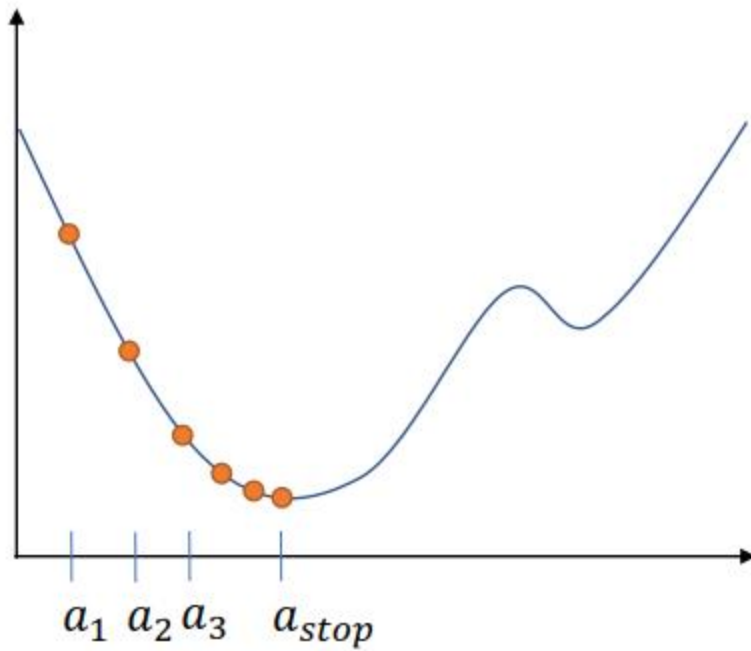
# Fully connected Neural network



- A deep network is a neural network with many layers
- A neuron in a linear function followed for an activation function
- Activation function must be non-linear
- A loss function measures how close is the created function (network) from a desired output
- The “training” is the process of find parameters (‘weights’) that reduces the loss functions
- Updating the weights as  $w_{new} = w_{prev} - \alpha \frac{dJ}{dW}$  reduces the loss
- An algorithm named back-propagation allows to compute  $\frac{dJ}{dW}$  for all the weights of the network in 2 steps: 1 forward, 1 backward

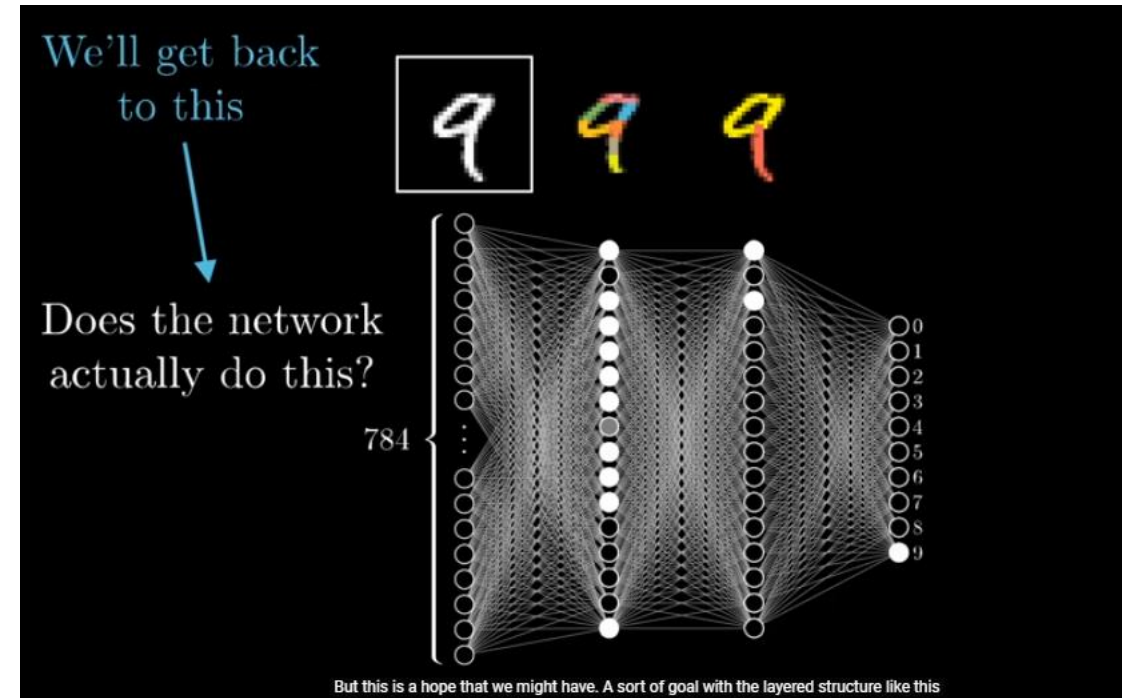
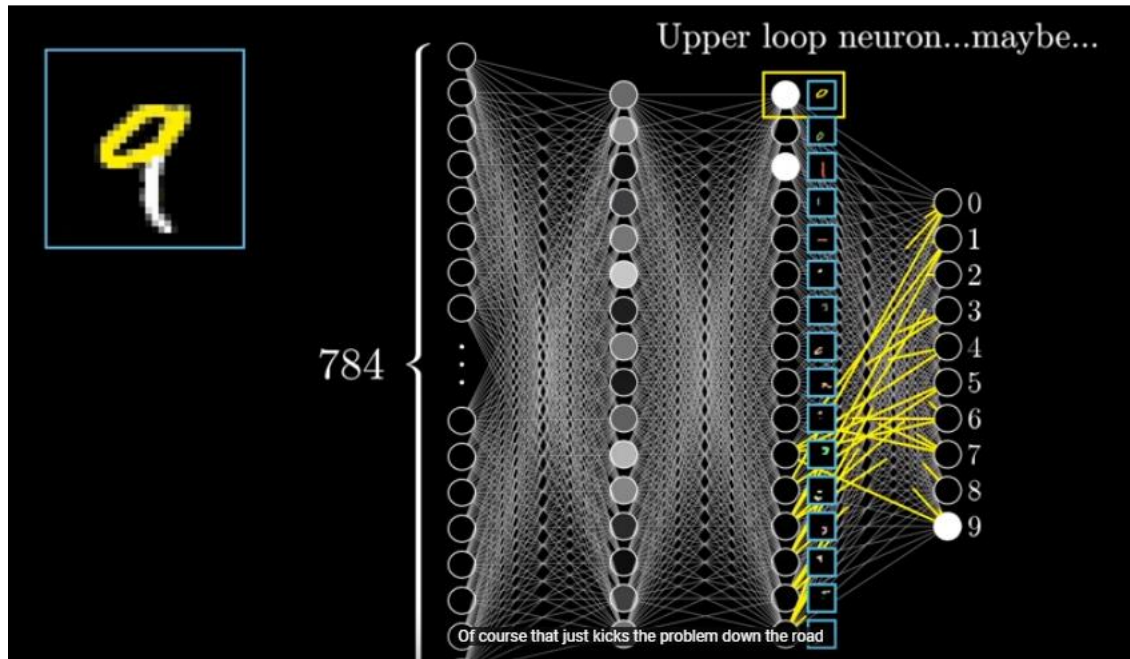
# Learning rate

$$w_{new} = w_{prev} - \alpha \frac{dJ}{dW}$$





# An example





# Softmax

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

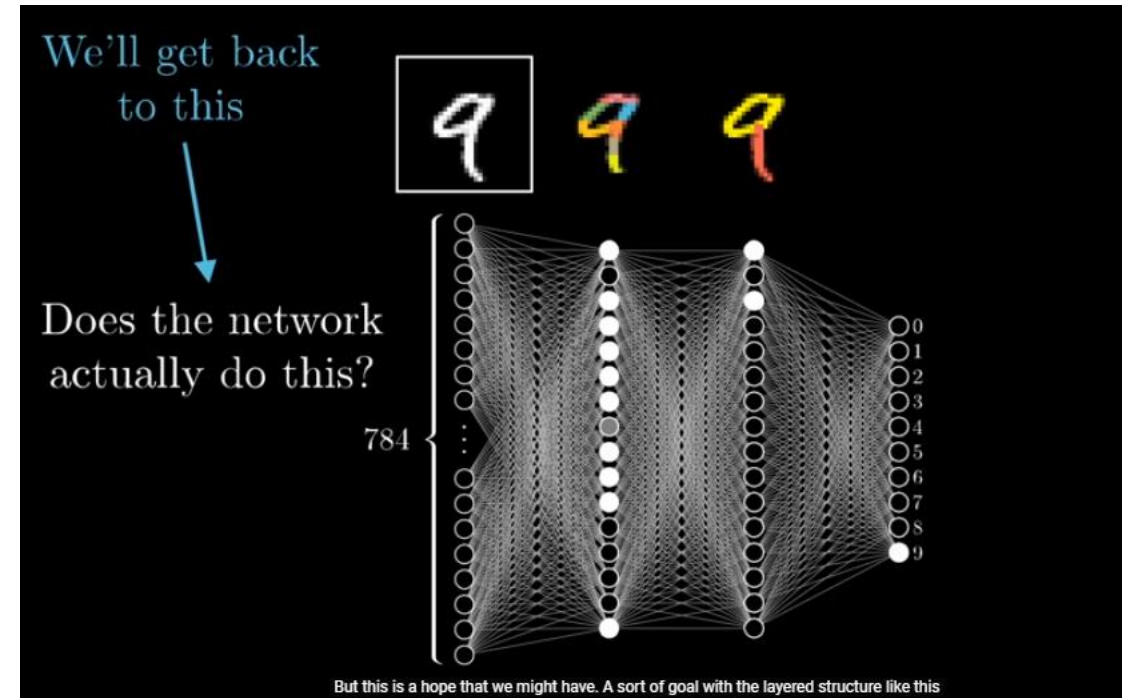
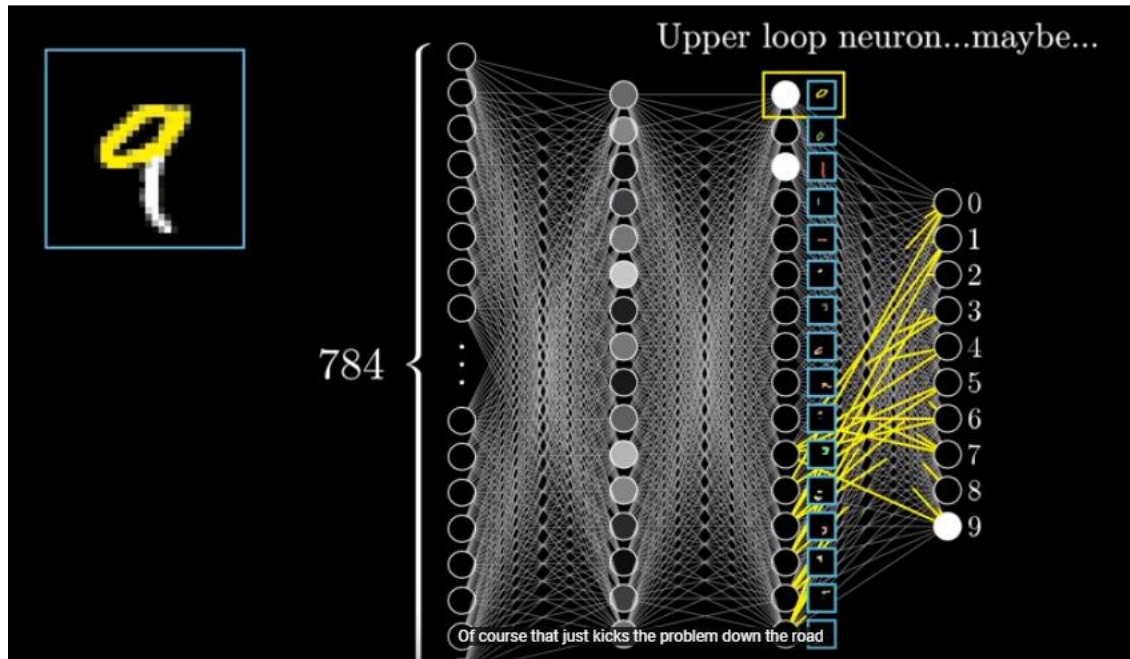
Used to interpret outputs as probabilities

$\vec{z}$	The input vector to the softmax function, made up of (z0, ... zK)
$z_i$	All the $z_i$ values are the elements of the input vector to the softmax function, and they can take any real value, positive, zero or negative. For example a neural network could have output a vector such as (-0.62, 8.12, 2.53), which is not a valid probability distribution, hence why the softmax would be necessary.
$e^{z_i}$	The standard exponential function is applied to each element of the input vector. This gives a positive value above 0, which will be very small if the input was negative, and very large if the input was large. However, it is still not fixed in the range (0, 1) which is what is required of a probability.
$\sum_{j=1}^K e^{z_j}$	The term on the bottom of the formula is the normalization term. It ensures that all the output values of the function will sum to 1 and each be in the range (0, 1), thus constituting a valid probability distribution.
$K$	The number of classes in the multi-class classifier.

$$\begin{aligned} \begin{bmatrix} P(\text{cat}) \\ P(\text{dog}) \end{bmatrix} &= \sigma\left(\begin{bmatrix} 1.2 \\ 0.3 \end{bmatrix}\right) \\ &= \begin{bmatrix} \frac{e^{1.2}}{e^{1.2} + e^{0.3}} \\ \frac{e^{0.3}}{e^{1.2} + e^{0.3}} \end{bmatrix} \\ &= \begin{bmatrix} 0.71 \\ 0.29 \end{bmatrix} \end{aligned}$$

[The Softmax function and its derivative - Eli Bendersky's website \(thegreenplace.net\)](http://thegreenplace.net)

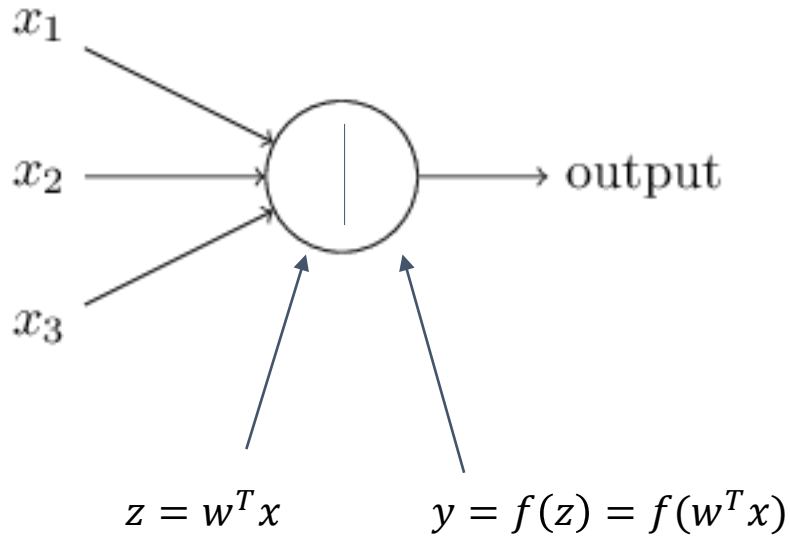
# An example



# A Neuron

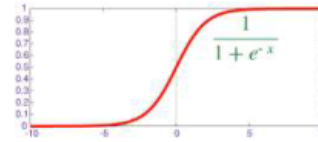
## A REVIEW

# Activations and their derivatives



$$x = [x_1, x_2, x_3, 1]$$

SIGMOID



LOGISTIC FUNCTION

$$f(z) = \frac{1}{1 + \exp(-z)}$$

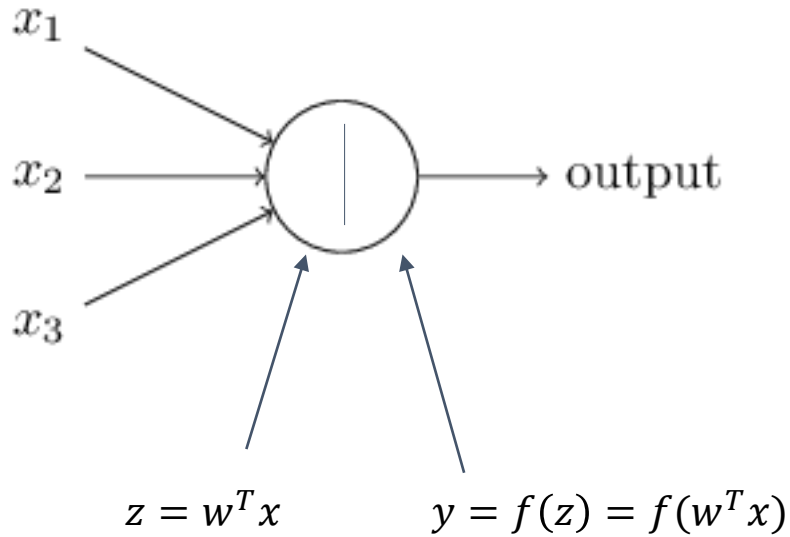
$$f'(z) = f(z)(1 - f(z))$$

$$f(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x},$$

$$\frac{d}{dx} f(x) = \frac{e^x \cdot (1 + e^x) - e^x \cdot e^x}{(1 + e^x)^2} = \frac{e^x}{(1 + e^x)^2} = f(x)(1 - f(x))$$

# A Neuron

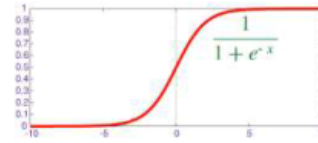
## A REVIEW



$$x = [x_1, x_2, x_3, 1]$$

## Activations and their derivatives

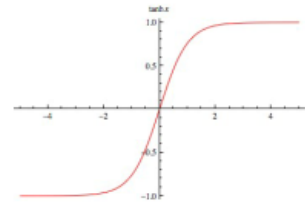
SIGMOID



LOGISTIC FUNCTION

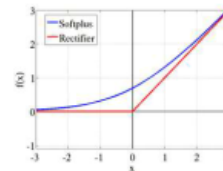
$$f(z) = \frac{1}{1 + \exp(-z)}$$

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



$$f(z) = \tanh(z)$$

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



$$f(z) = \begin{cases} 0, & z < 0 \\ z, & z \geq 0 \end{cases}$$

This space left intentionally (kind of) blank

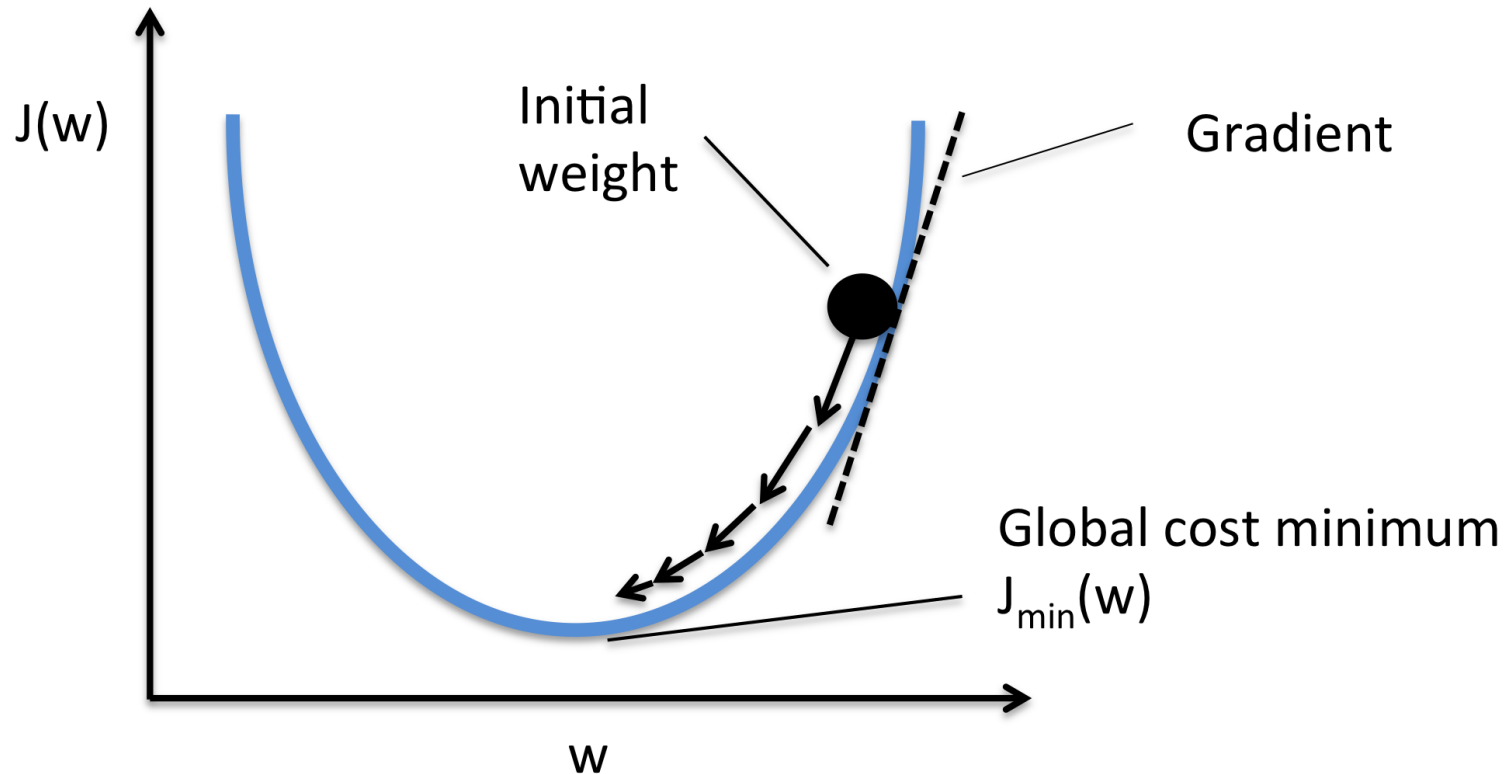
$$f(z) = \log(1 + \exp(z))$$

$$f'(z) = \frac{1}{1 + \exp(-z)}$$

*softplus* or *SmoothReLU* function

# IN OUR CASE THE LOSS FUNCTION

## How to minimize a function ?

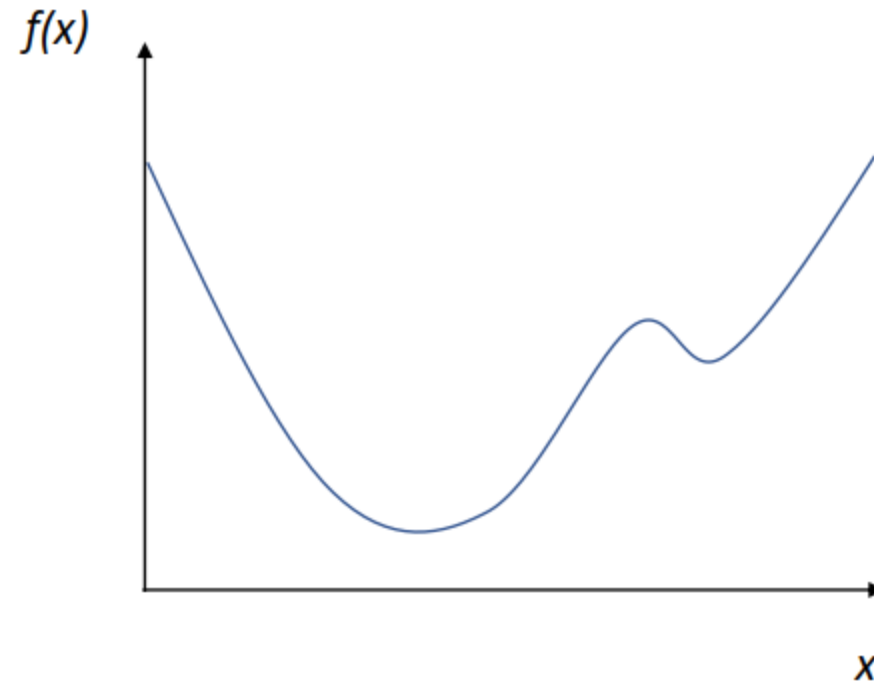


Repeat until there is almost not change

$$w_{new} = w_{prev} - \alpha \frac{dJ}{dW}$$

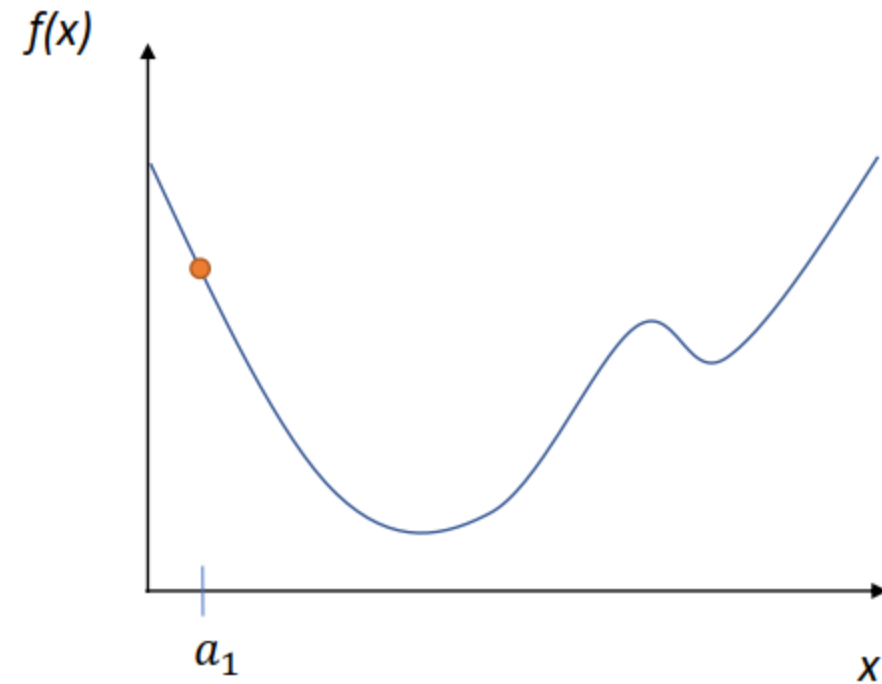
HOW TO COMPUTE THIS GRADIENT?

# Gradient descent



# General approach

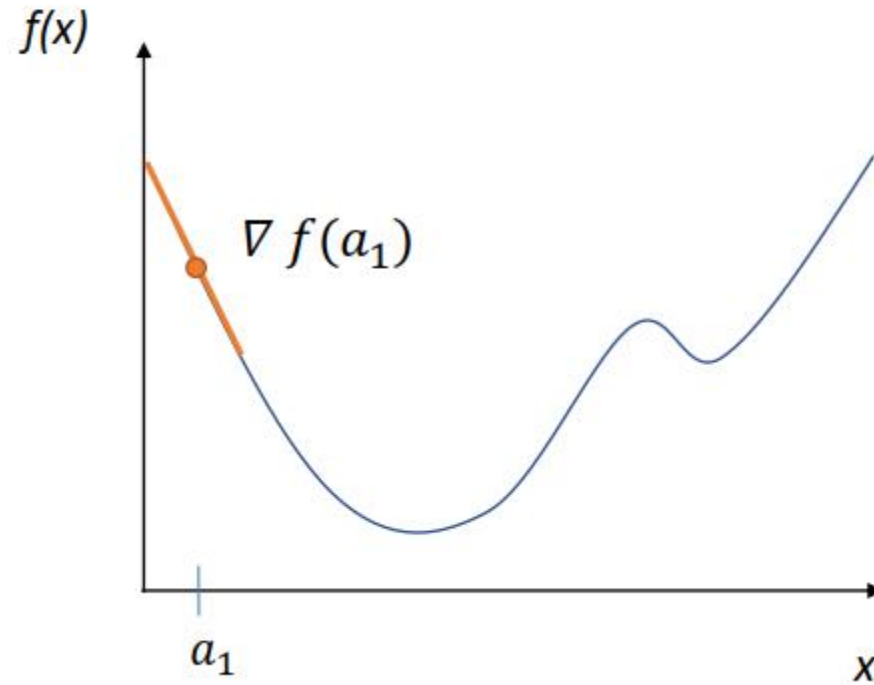
Pick random starting point.





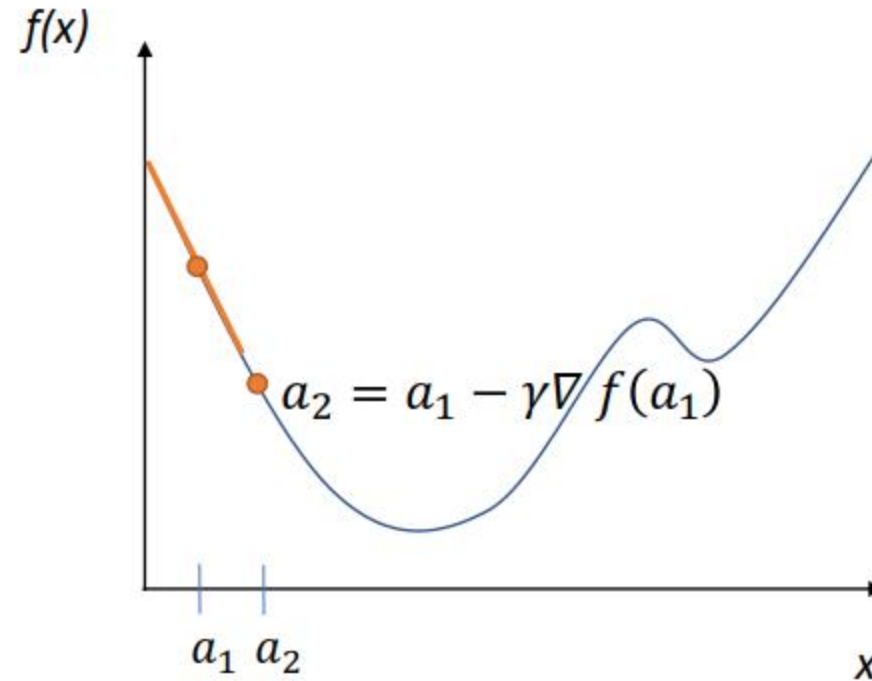
# General approach

Compute gradient at point (analytically or by finite differences)



# General approach

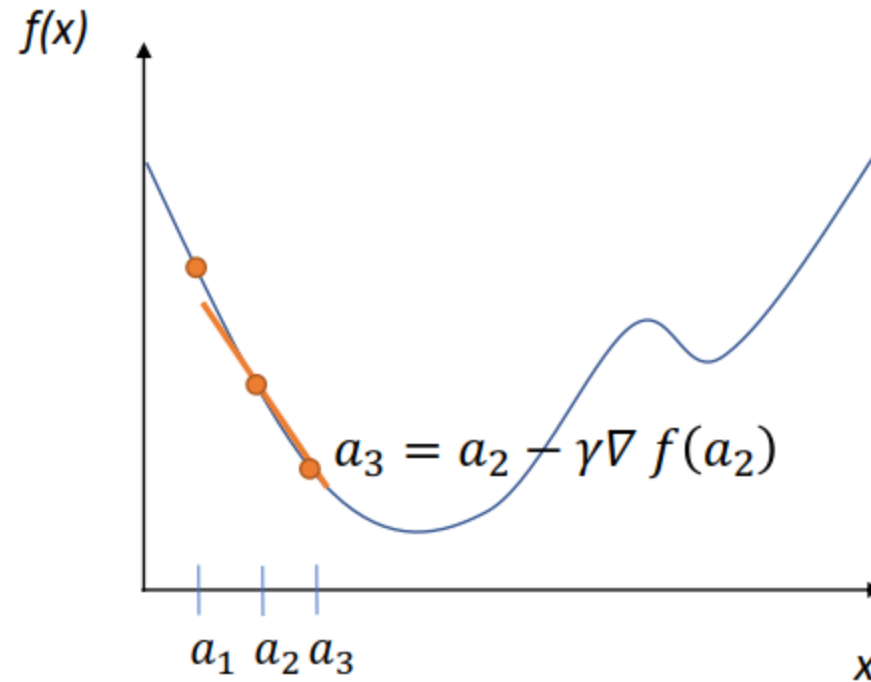
Move along parameter space in direction of negative gradient



$\gamma$  = amount to move  
= *learning rate*

# General approach

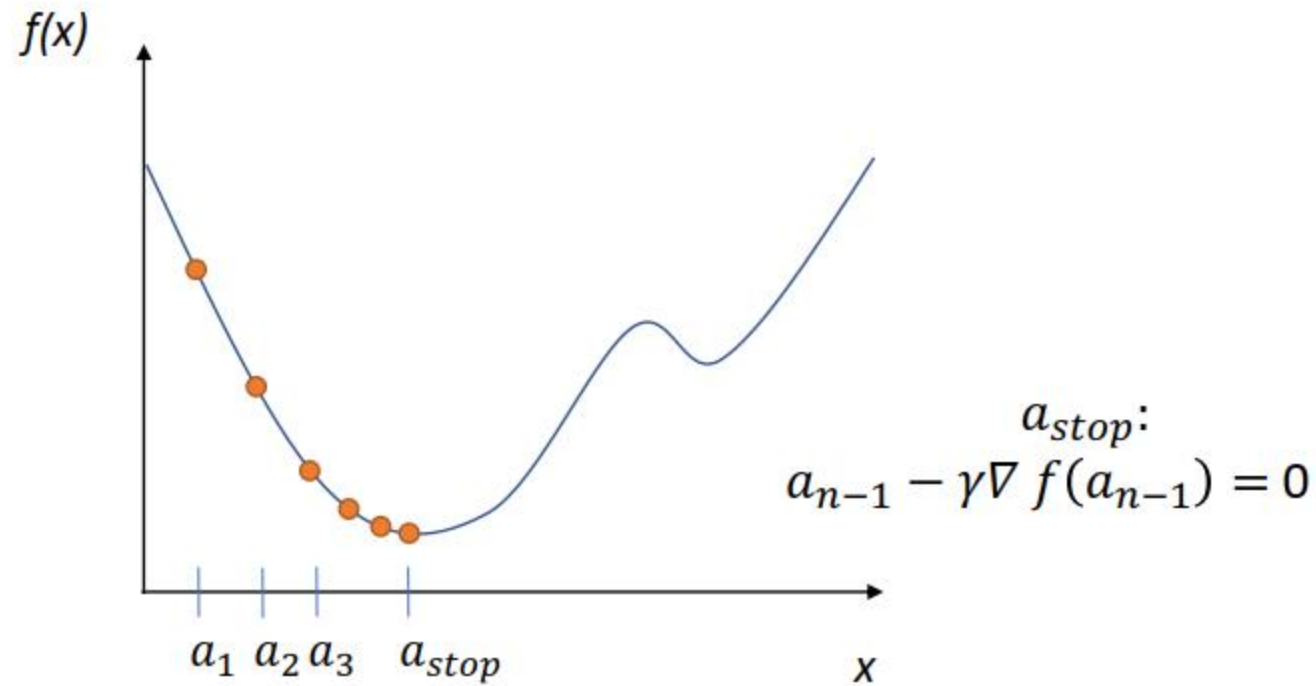
Move along parameter space in direction of negative gradient.



$\gamma$  = amount to move  
= learning rate

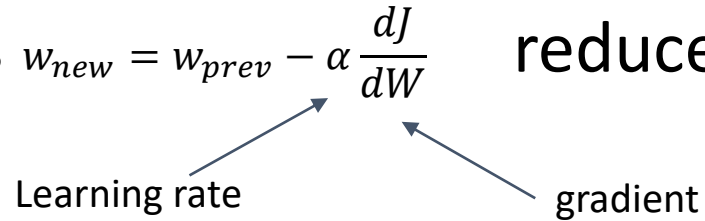
# General approach

Stop when we don't move any more.



# Gradient Descent

- The gradient is the direction of fastest increase in  $J(X)$
- Updating the weights as  $w_{new} = w_{prev} - \alpha \frac{dJ}{dW}$  reduces the loss



## The Approach of Gradient Descent



- Iterative solution:
  - Start at some point
  - Find direction in which to shift this point to decrease error
    - This can be found from the derivative of the function
      - A positive derivative  $\rightarrow$  moving left decreases error
      - A negative derivative  $\rightarrow$  moving right decreases error
  - Shift point in this direction

## Overall Gradient Descent Algorithm

- Initialize:
  - $x^0$
  - $k = 0$
- While  $|f(x^{k+1}) - f(x^k)| > \epsilon$ 
  - $x^{k+1} = x^k - \eta^k \nabla f(x^k)^T$
  - $k = k + 1$



# Train with Gradient Descent

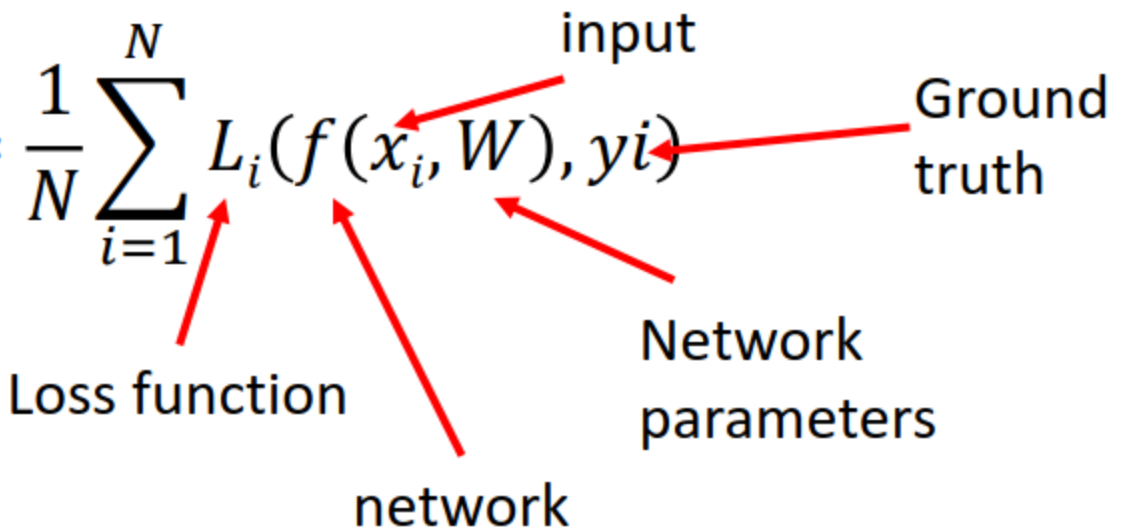
- $x^i, y^i = n$  training examples
- $f(\mathbf{x})$  = feed forward network
- $L(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})$  = some *loss function*

*Loss function* measures how 'good' our network is at classifying the training examples wrt. the parameters of the model (the perceptron weights).

# Loss Function

- Way to define how good the network is performing
  - In terms of prediction
- Network training (Optimization)
  - Find the best network parameters to minimize the loss

Total Error  $(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)$



input

Ground truth

Loss function

network

Network parameters

# Loss Functions and total Error

## Cross-Entropy

a.k.a. log loss

$$\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^n t_j \log(p_j)$$

*n classes*

*t<sub>j</sub> is the truth label*

*p<sub>j</sub> is the Softmax probability for the j<sup>th</sup> class*

N samples

Binary Cross Entropy

$$-\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

Ground-truth

Predicted value

- Mean squared error (MSE)



$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$



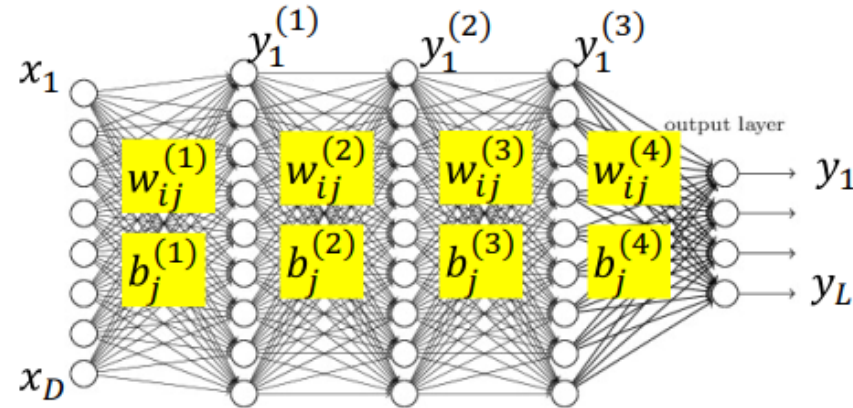
$$L_{CE} = - \sum_{i=1}^n t_i \log(p_i), \text{ for } n \text{ classes,}$$

where  $t_i$  is the truth label and  $p_i$  is the Softmax probability for the  $i^{th}$  class.

**The lower the loss, the more accurate the model** 🤖

<p><b>DOG</b></p> 	<p><math>y = [0.4, 0.4, 0.2]</math></p> <p><math>t = [0, 1, 0]</math></p>	<p><math>L(\mathbf{y}, \mathbf{t}) = -0 \times \ln 0.4 - 1 \times \ln 0.4 - 0 \times \ln 0.2</math></p> <p><math>= 0.92</math></p>
<p><b>HORSE</b></p> 	<p><math>y = [0.1, 0.2, 0.7]</math></p> <p><math>t = [0, 0, 1]</math></p>	<p><math>L(\mathbf{y}, \mathbf{t}) = -0 \times \ln 0.1 - 0 \times \ln 0.2 - 1 \times \ln 0.7</math></p> <p><math>= 0.36</math></p>

# Notation



- The input layer is the 0<sup>th</sup> layer
- We will represent the output of the  $i$ -th perceptron of the  $k$ <sup>th</sup> layer as  $y_i^{(k)}$ 
  - **Input to network:**  $y_i^{(0)} = x_i$
  - **Output of network:**  $y_i = y_i^{(N)}$
- We will represent the weight of the connection between the  $i$ -th unit of the  $k-1$ th layer and the  $j$ th unit of the  $k$ -th layer as  $w_{ij}^{(k)}$ 
  - The bias to the  $j$ th unit of the  $k$ -th layer is  $b_j^{(k)}$

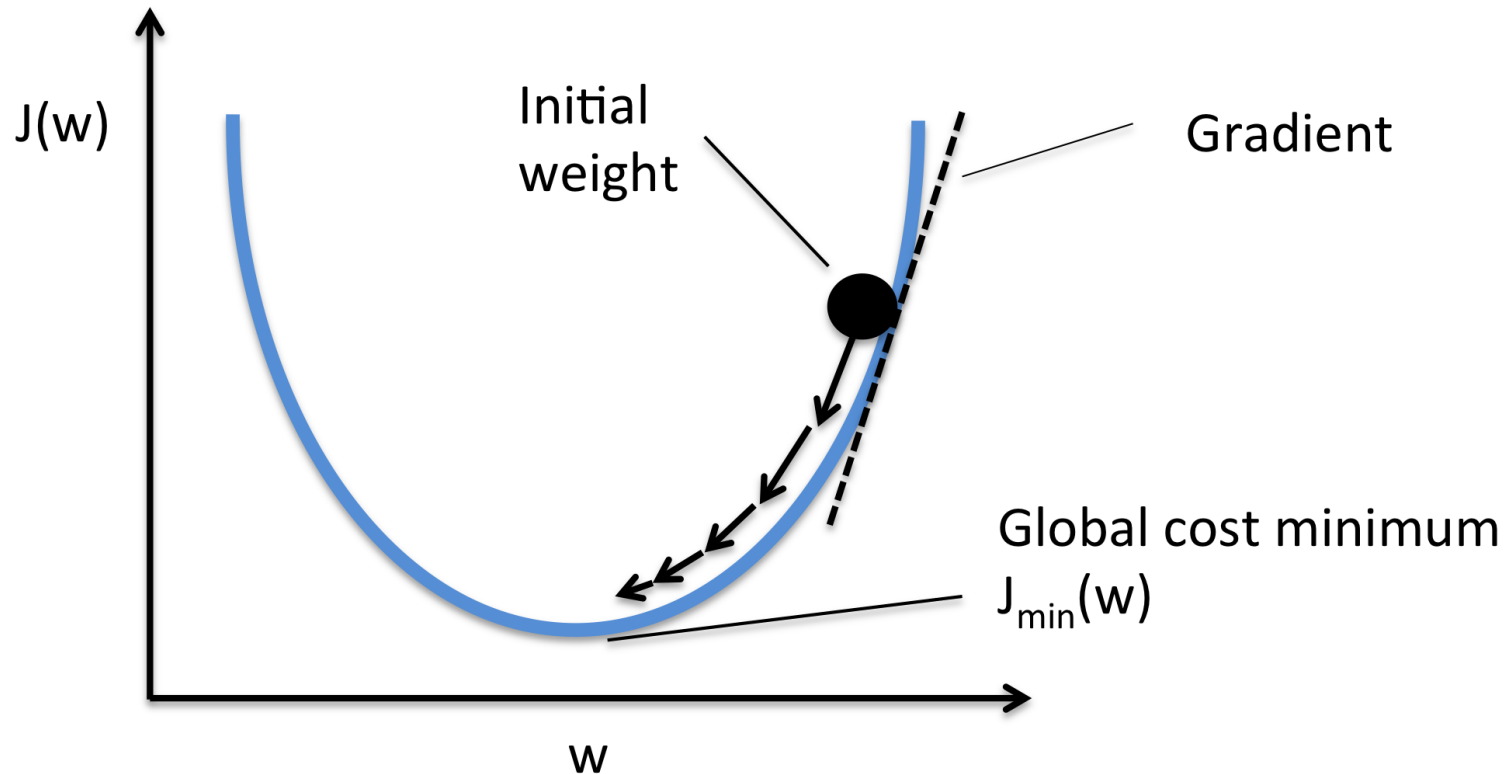


# Training steps

- Define network
- Loss function
- Initialize network parameters
- Get training data
  - Prepare batches
- Feedforward one batch
  - Compute loss
  - Update network parameters
  - Repeat

IN OUR CASE THE LOSS FUNCTION

# How to minimize a function ?



Repeat until there is almost not change

$$w_{new} = w_{prev} - \alpha \frac{dJ}{dW}$$

HOW TO COMPUTE THIS GRADIENT?



# Training Neural Nets through Gradient Descent

Total training error:

$$Err = \frac{1}{T} \sum_t Div(\mathbf{Y}_t, \mathbf{d}_t)$$

- Gradient descent algorithm:
- Initialize all weights and biases  $\{w_{ij}^{(k)}\}$ 
  - Using the extended notation: the bias is also a weight
- Do:
  - For every layer  $k$  for all  $i, j$ , update:
    - $w_{i,j}^{(k)} = w_{i,j}^{(k)} - \eta \frac{dErr}{dw_{i,j}^{(k)}}$
- Until  $Err$  has converged

Assuming the bias is also represented as a weight

Example: L2

$$Div = \frac{1}{2} (y_t - d_t)^2$$

$$\frac{dDiv}{dy_i} = (y_t - d_t)$$

# The derivative

Total training error:

$$Err = \frac{1}{T} \sum_t Div(\mathbf{Y}_t, \mathbf{d}_t)$$

- Computing the derivative

Total derivative:

$$\frac{dErr}{dw_{i,j}^{(k)}} = \frac{1}{T} \sum_t \frac{dDiv(\mathbf{Y}_t, \mathbf{d}_t)}{dw_{i,j}^{(k)}}$$

# The derivative

Total training error:

$$Err = \frac{1}{T} \sum_t Div(\mathbf{Y}_t, \mathbf{d}_t)$$

Total derivative:

$$\frac{dErr}{dw_{i,j}^{(k)}} = \frac{1}{T} \sum_t \frac{dDiv(\mathbf{Y}_t, \mathbf{d}_t)}{dw_{i,j}^{(k)}}$$

- So we must first figure out how to compute the derivative of divergences of individual training inputs


# Calculus Refresher: Basic rules of calculus

For any differentiable function

$$y = f(x)$$

with derivative

$$\frac{dy}{dx}$$

the following must hold for sufficiently small  $\Delta x$    $\Delta y \approx \frac{dy}{dx} \Delta x$

For any differentiable function

$$y = f(x_1, x_2, \dots, x_M)$$

with partial derivatives

$$\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_M}$$

the following must hold for sufficiently small  $\Delta x_1, \Delta x_2, \dots, \Delta x_M$

$$\Delta y \approx \frac{\partial y}{\partial x_1} \Delta x_1 + \frac{\partial y}{\partial x_2} \Delta x_2 + \dots + \frac{\partial y}{\partial x_M} \Delta x_M$$



# Calculus Refresher: Chain rule

For any nested function  $y = f(g(x))$

$$\frac{dy}{dx} = \frac{\partial y}{\partial g(x)} \frac{dg(x)}{dx}$$

Check - we can confirm that :  $\Delta y = \frac{dy}{dx} \Delta x$

$$z = g(x) \Rightarrow \Delta z = \frac{dg(x)}{dx} \Delta x$$

$$y = f(z) \Rightarrow \Delta y = \frac{dy}{dz} \Delta z = \frac{dy}{dz} \frac{dg(x)}{dx} \Delta x$$





# Calculus Refresher: Distributed Chain rule

$$y = f(g_1(x), g_1(x), \dots, g_M(x))$$

$$\frac{dy}{dx} = \frac{\partial y}{\partial g_1(x)} \frac{dg_1(x)}{dx} + \frac{\partial y}{\partial g_2(x)} \frac{dg_2(x)}{dx} + \dots + \frac{\partial y}{\partial g_M(x)} \frac{dg_M(x)}{dx}$$

Check:  $\Delta y = \frac{dy}{dx} \Delta x$

$$\Delta y = \frac{\partial y}{\partial g_1(x)} \Delta g_1(x) + \frac{\partial y}{\partial g_2(x)} \Delta g_2(x) + \dots + \frac{\partial y}{\partial g_M(x)} \Delta g_M(x)$$

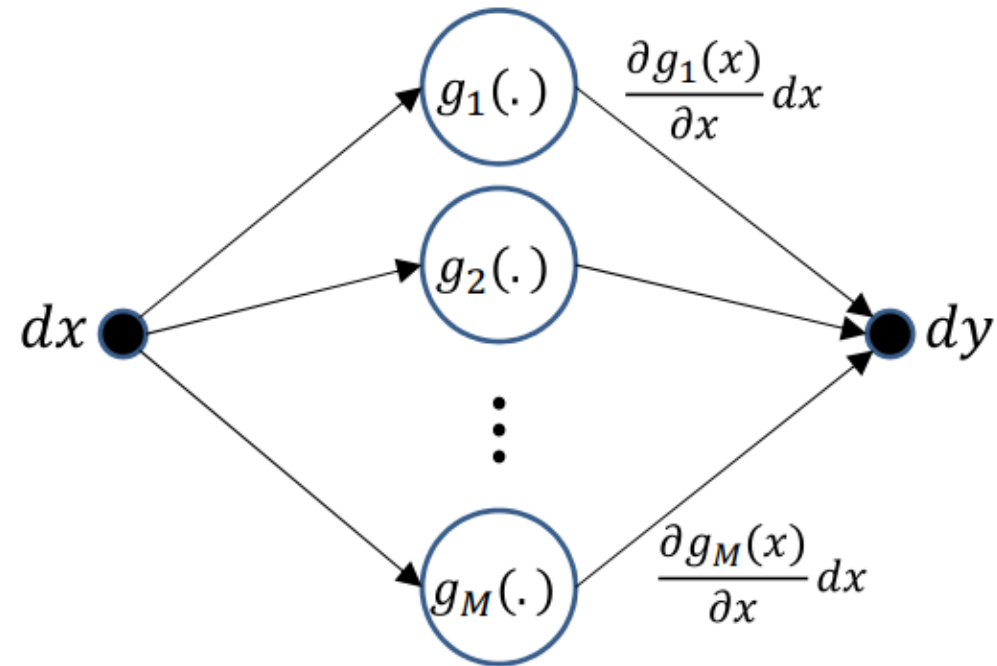
$$\Delta y = \frac{\partial y}{\partial g_1(x)} \frac{dg_1(x)}{dx} \Delta x + \frac{\partial y}{\partial g_2(x)} \frac{dg_2(x)}{dx} \Delta x + \dots + \frac{\partial y}{\partial g_M(x)} \frac{dg_M(x)}{dx} \Delta x$$

$$\Delta y = \left( \frac{\partial y}{\partial g_1(x)} \frac{dg_1(x)}{dx} + \frac{\partial y}{\partial g_2(x)} \frac{dg_2(x)}{dx} + \dots + \frac{\partial y}{\partial g_M(x)} \frac{dg_M(x)}{dx} \right) \Delta x$$



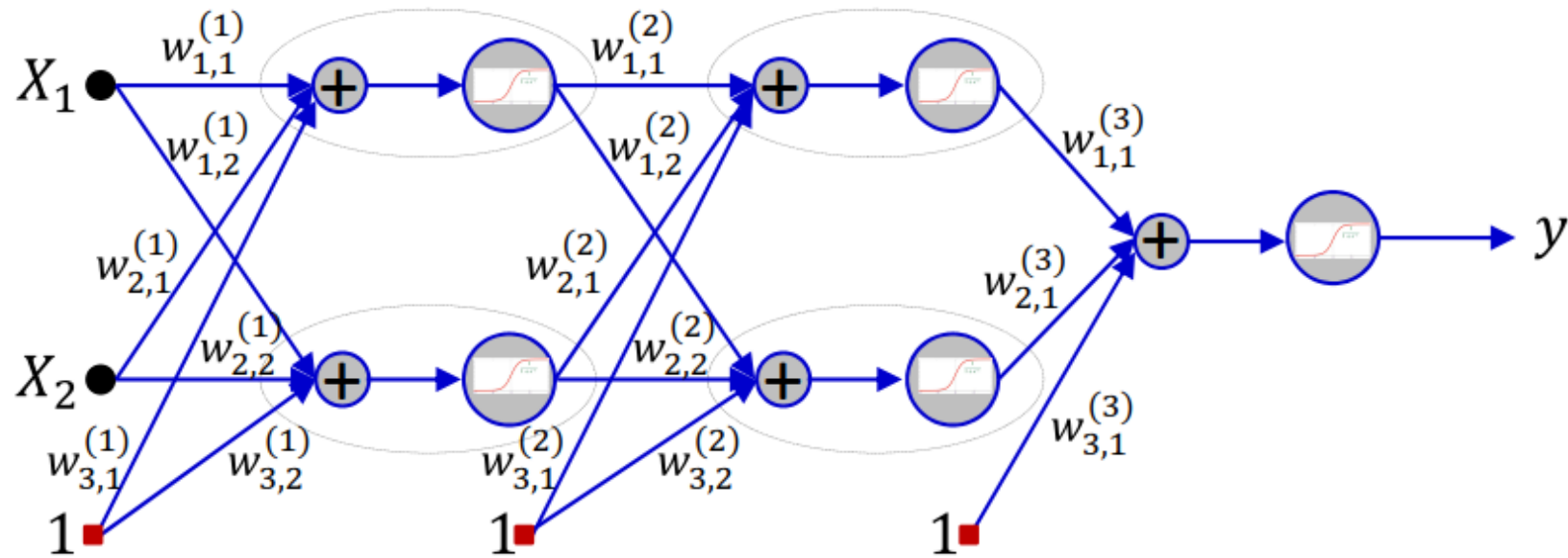
96

# Distributed Chain Rule: Influence Diagram



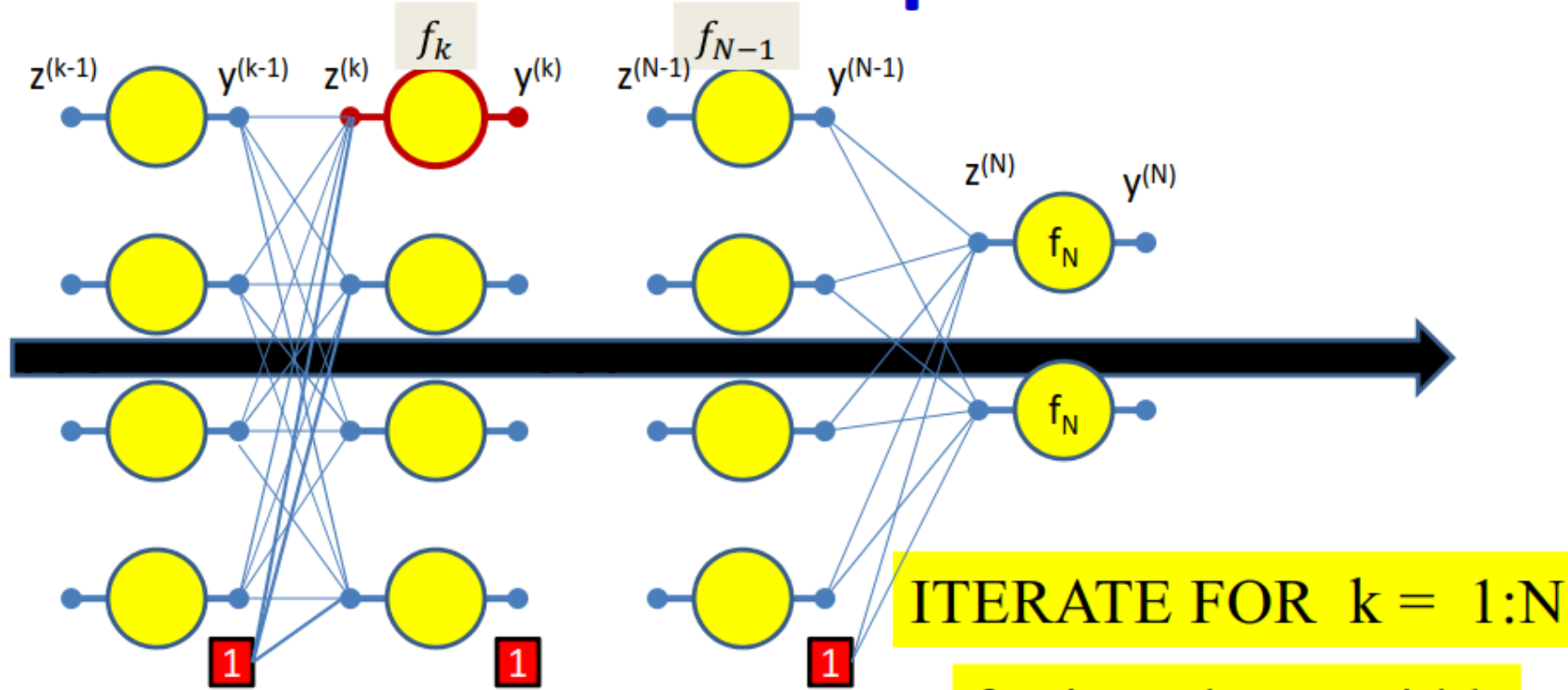
- Small perturbations in  $x$  cause small perturbations in each of  $g_1 \dots g_M$ , each of which individually additively perturbs  $y$

# A first closer look at the network



- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs
- Expanded **with all weights and activations shown**
- The overall function is differentiable w.r.t every weight, bias and input

# Forward Computation

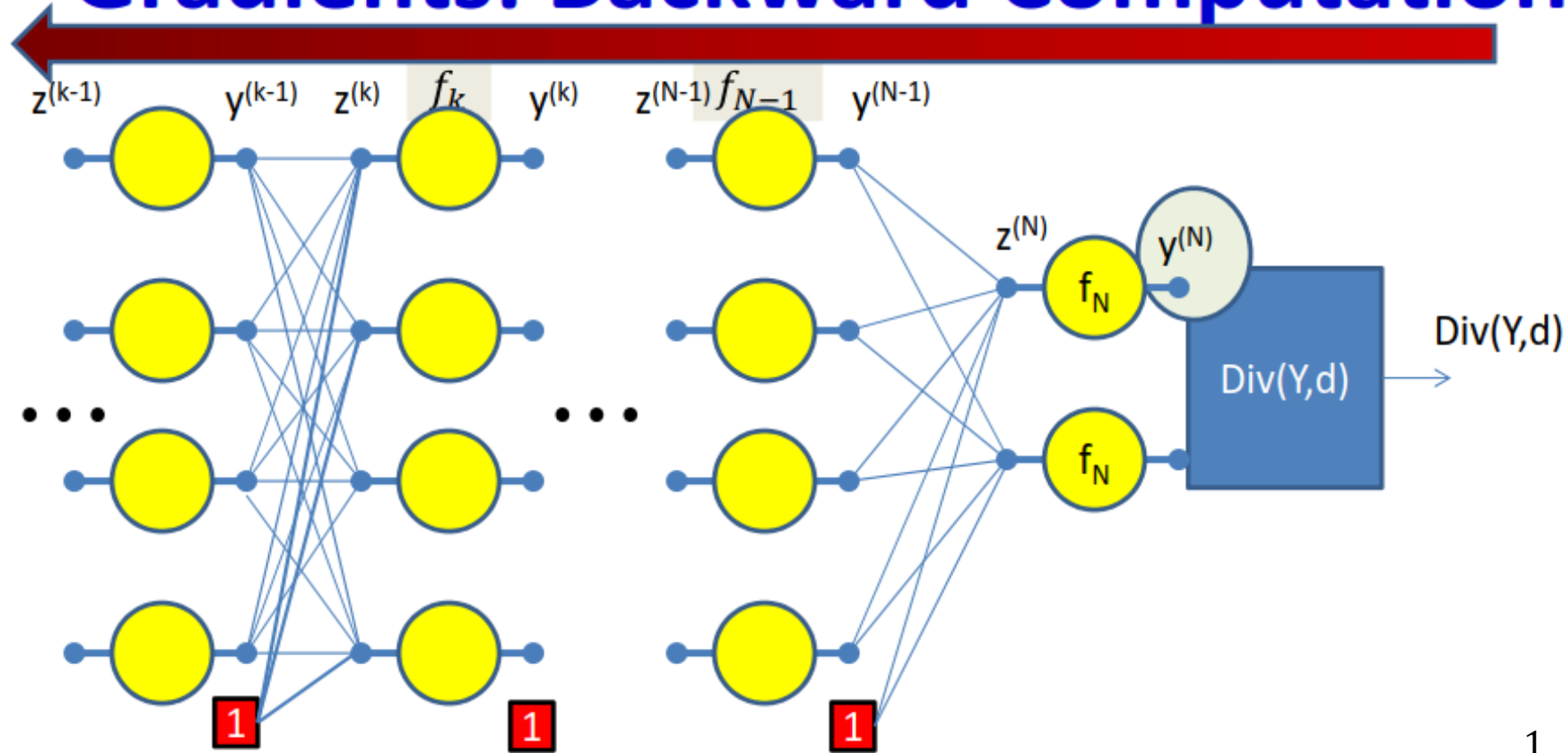


$$y_i^{(0)} = x_i$$

$$z_j^{(k)} = \sum_i w_{ij}^{(k)} y_i^{(k-1)}$$

$$y_j^{(k)} = f_k(z_j^{(k)})$$

# Gradients: Backward Computation

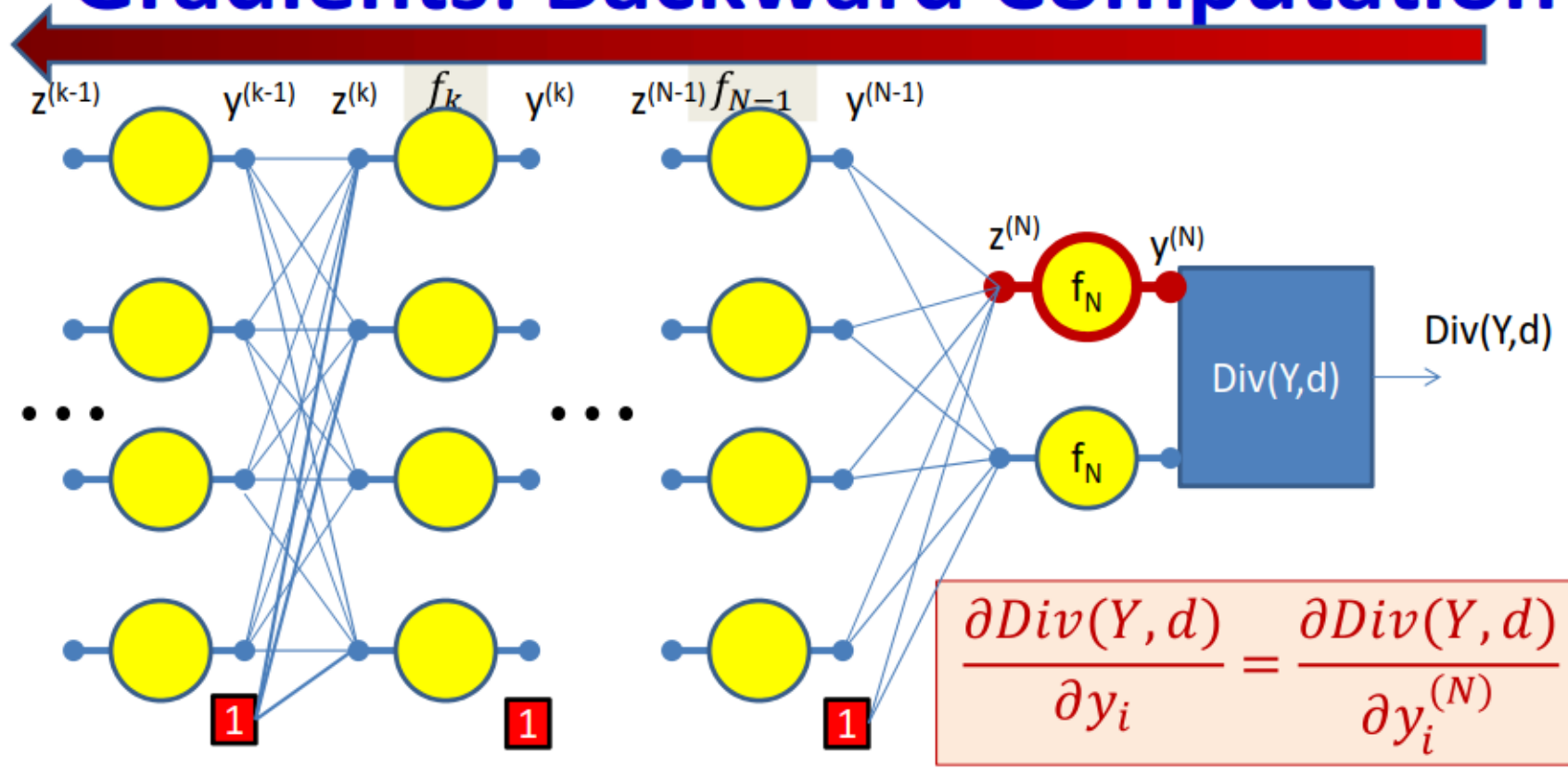


$$Div = \frac{1}{2}(y_t - d_t)^2$$

$$\frac{dDiv}{dy_i} = (y_t - d_t)$$

$$\frac{\partial Div(Y, d)}{\partial y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

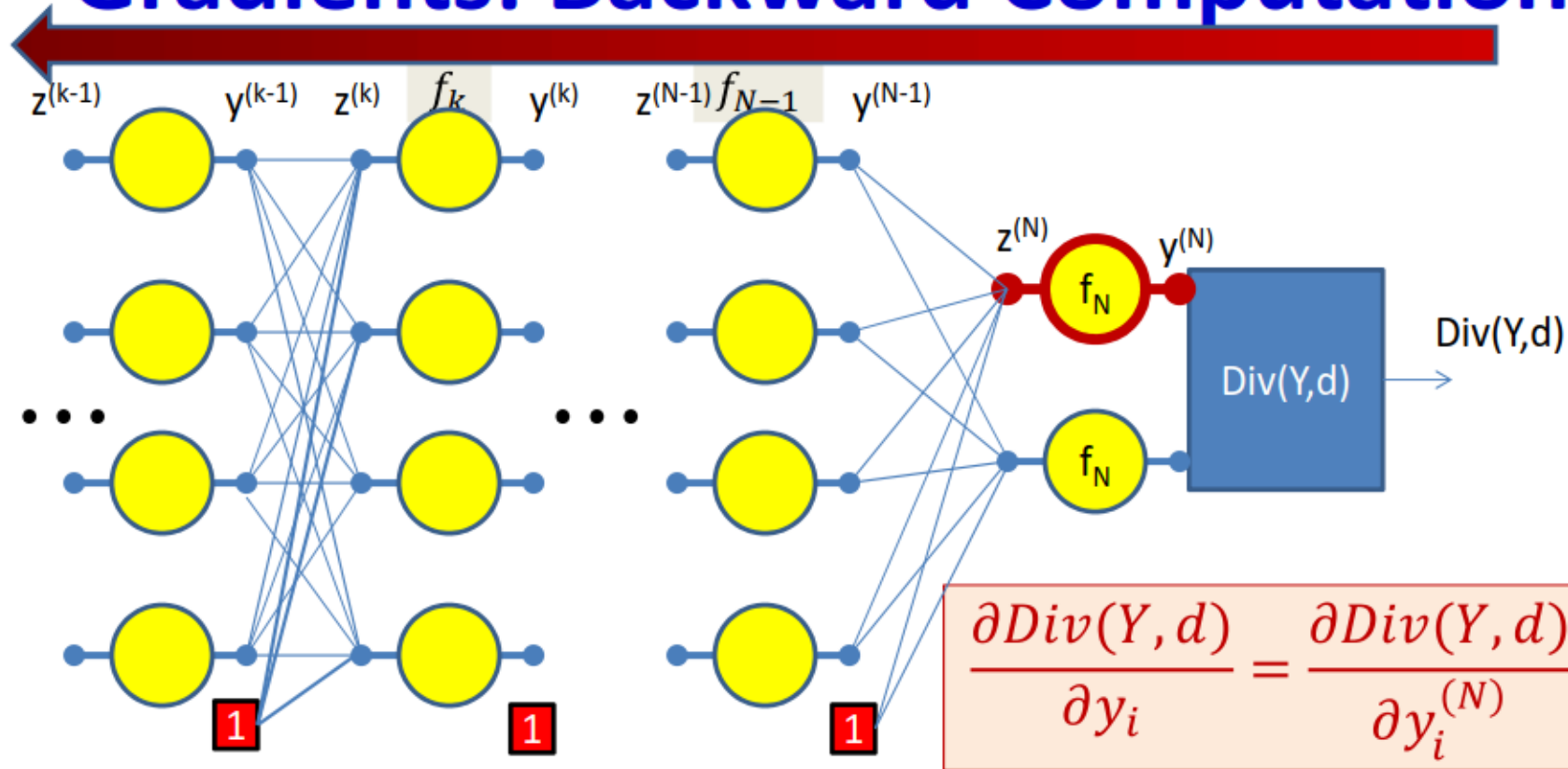
# Gradients: Backward Computation



$$\frac{\partial Div}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial Div}{\partial y_i} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$



# Gradients: Backward Computation



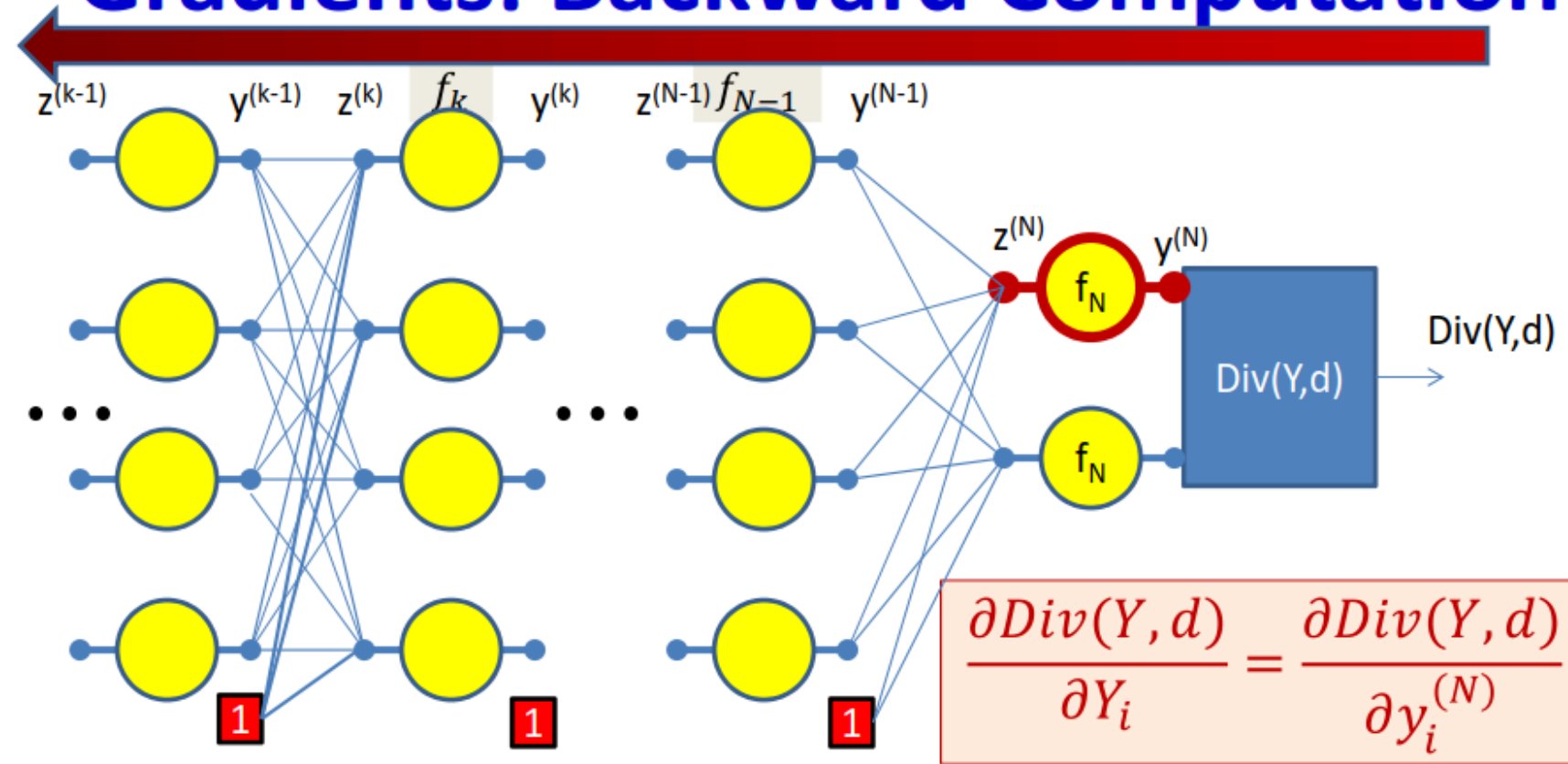
$$y_i^{[N]} = f(z_i^{[N]})$$

$$\frac{\partial y_i^{[N]}}{\partial z_i^{[N]}} = f^{[N]'}(z_i^{[N]})$$

$$\frac{\partial Div}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial Div}{\partial y_i} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$



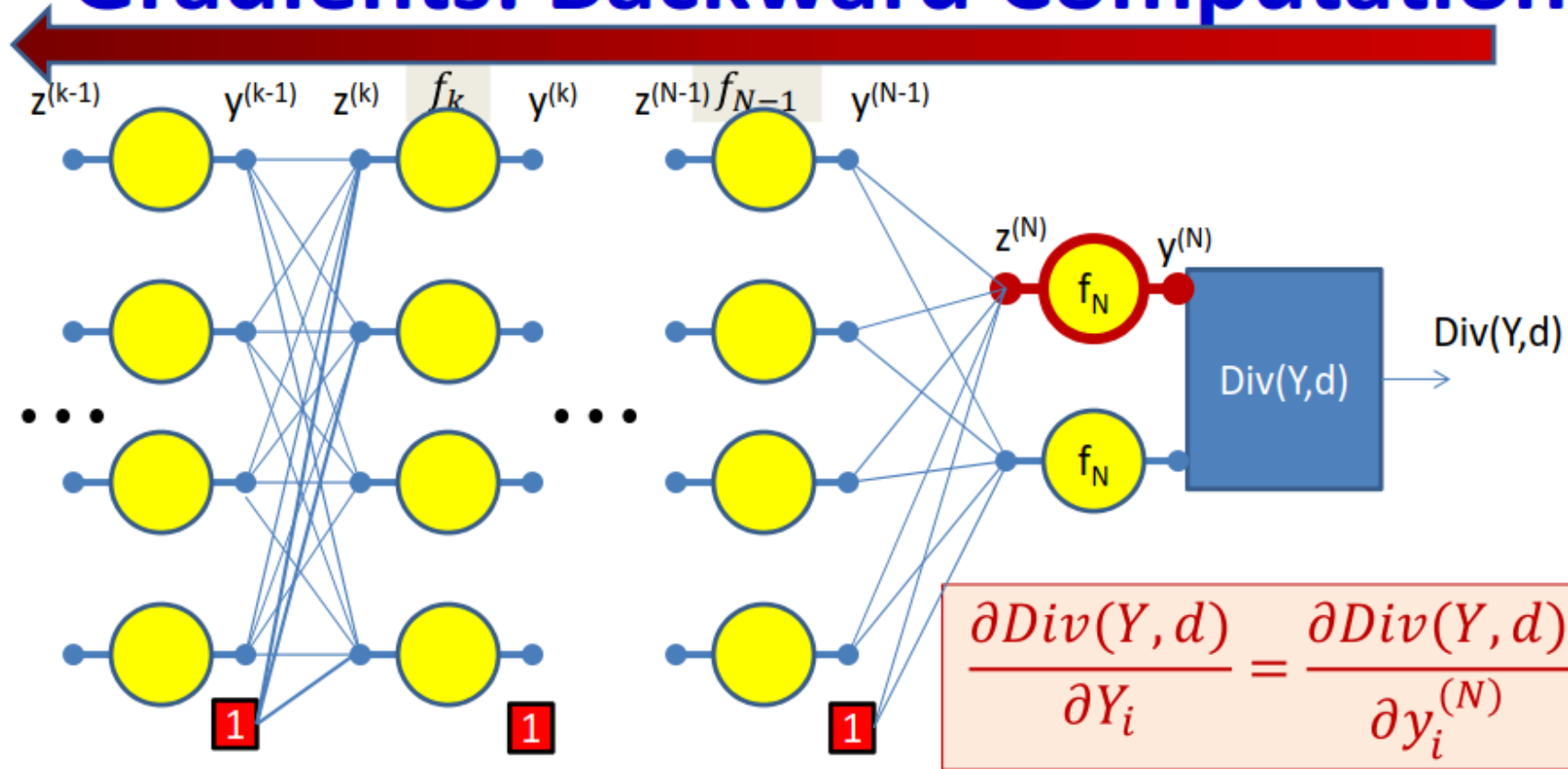
# Gradients: Backward Computation



$z_i^{(N)}$  computed during the forward pass

$$\frac{\partial Div}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial Div}{\partial Y_i} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

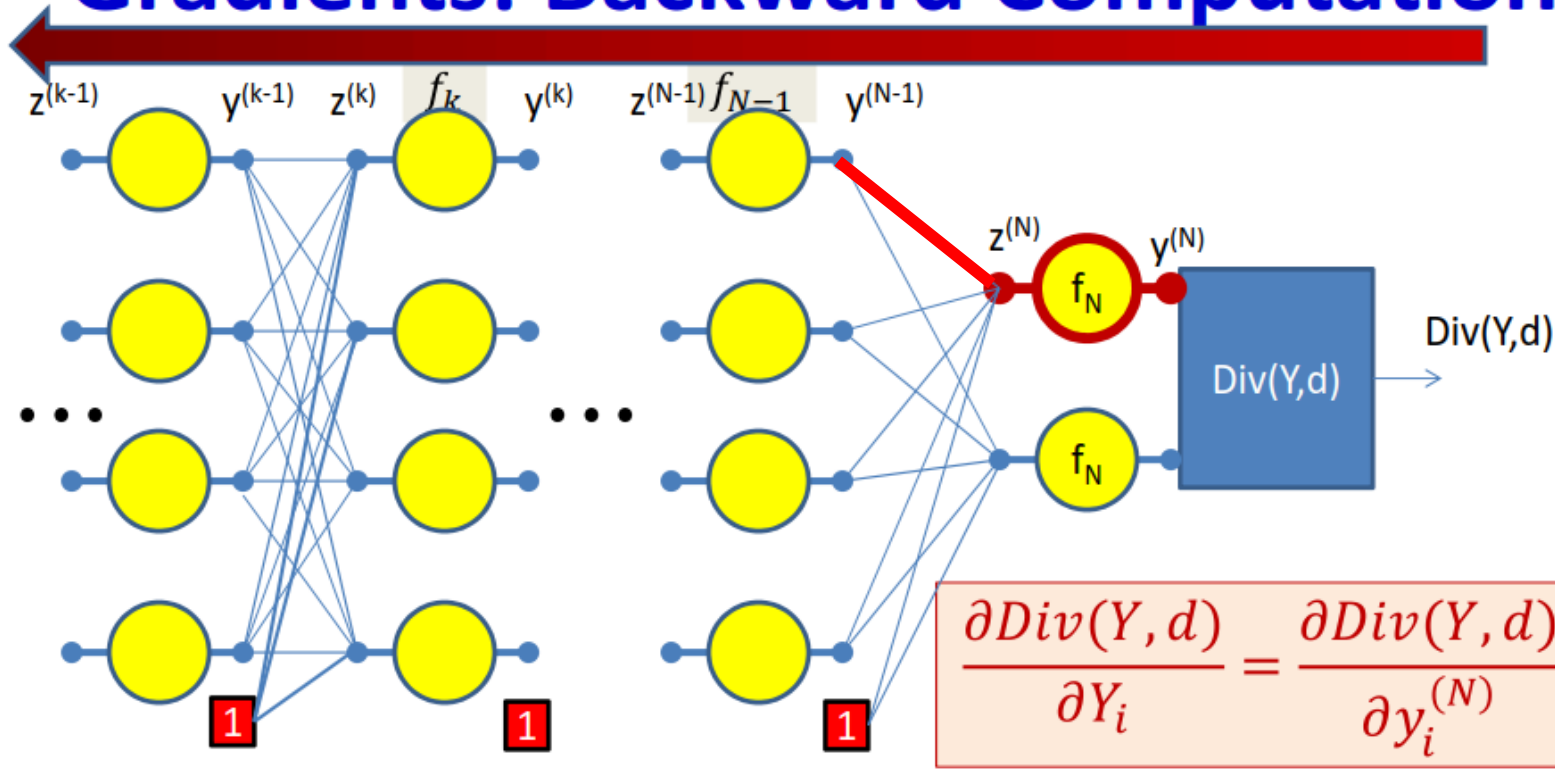
# Gradients: Backward Computation



Derivative of the activation function of Nth layer

$$\frac{\partial Div}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial Div}{\partial Y_i} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

# Gradients: Backward Computation



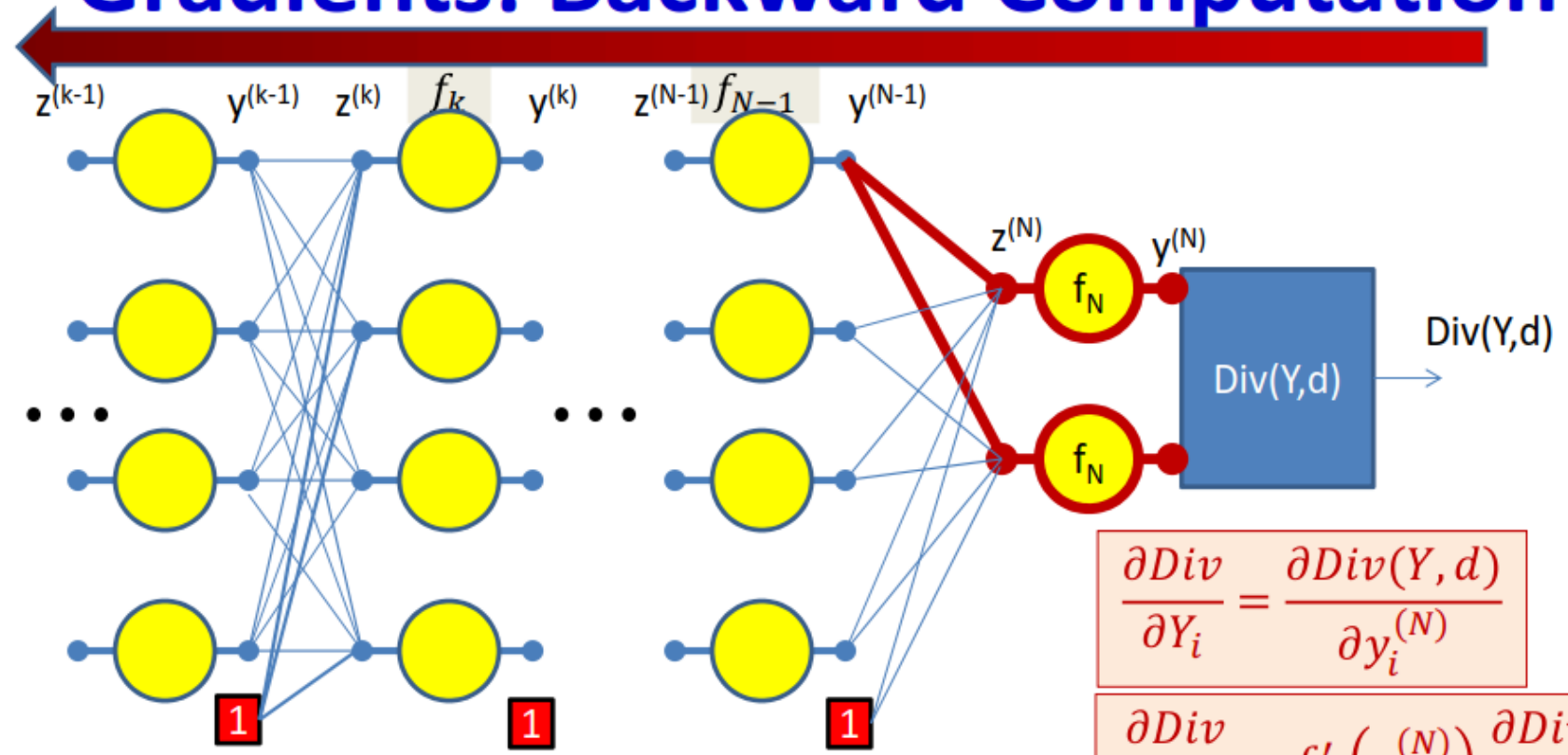
$$z_j^{[N]} = w^{Tj} y_i^{[N-1]}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = \frac{\partial z_j^{(k)}}{\partial w_{ij}^{(k)}} \frac{\partial Div}{\partial z_j^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div(Y, d)}{\partial Y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

$$\frac{\partial Div}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

# Gradients: Backward Computation



$$\frac{\partial Div}{\partial Y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

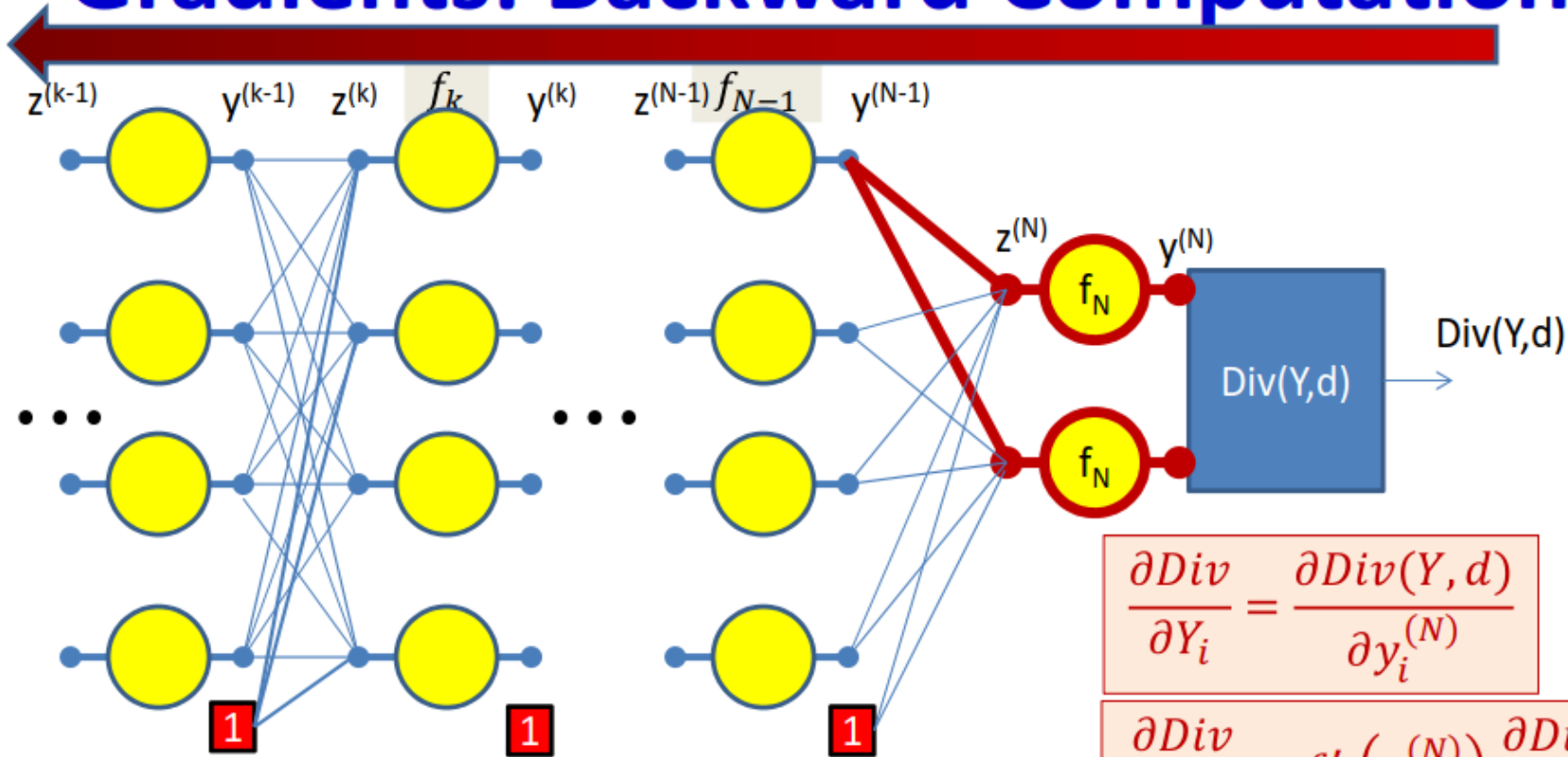
$$\frac{\partial Div}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

$$\frac{\partial Div}{\partial y_i^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}} = \sum_j w_{ij}^{(N)} \frac{\partial Div}{\partial z_j^{(N)}}$$

Because :

$$\frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} = w_{ij}^{(N)}$$

# Gradients: Backward Computation



$$\frac{\partial Div}{\partial Y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

$$\frac{\partial Div}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

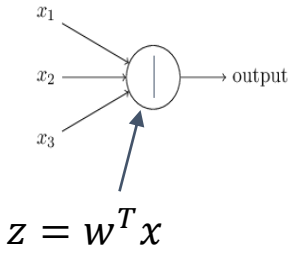
$$\frac{\partial Div}{\partial y_i^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}} = \sum_j w_{ij}^{(N)} \frac{\partial Div}{\partial z_j^{(N)}}$$

Because :

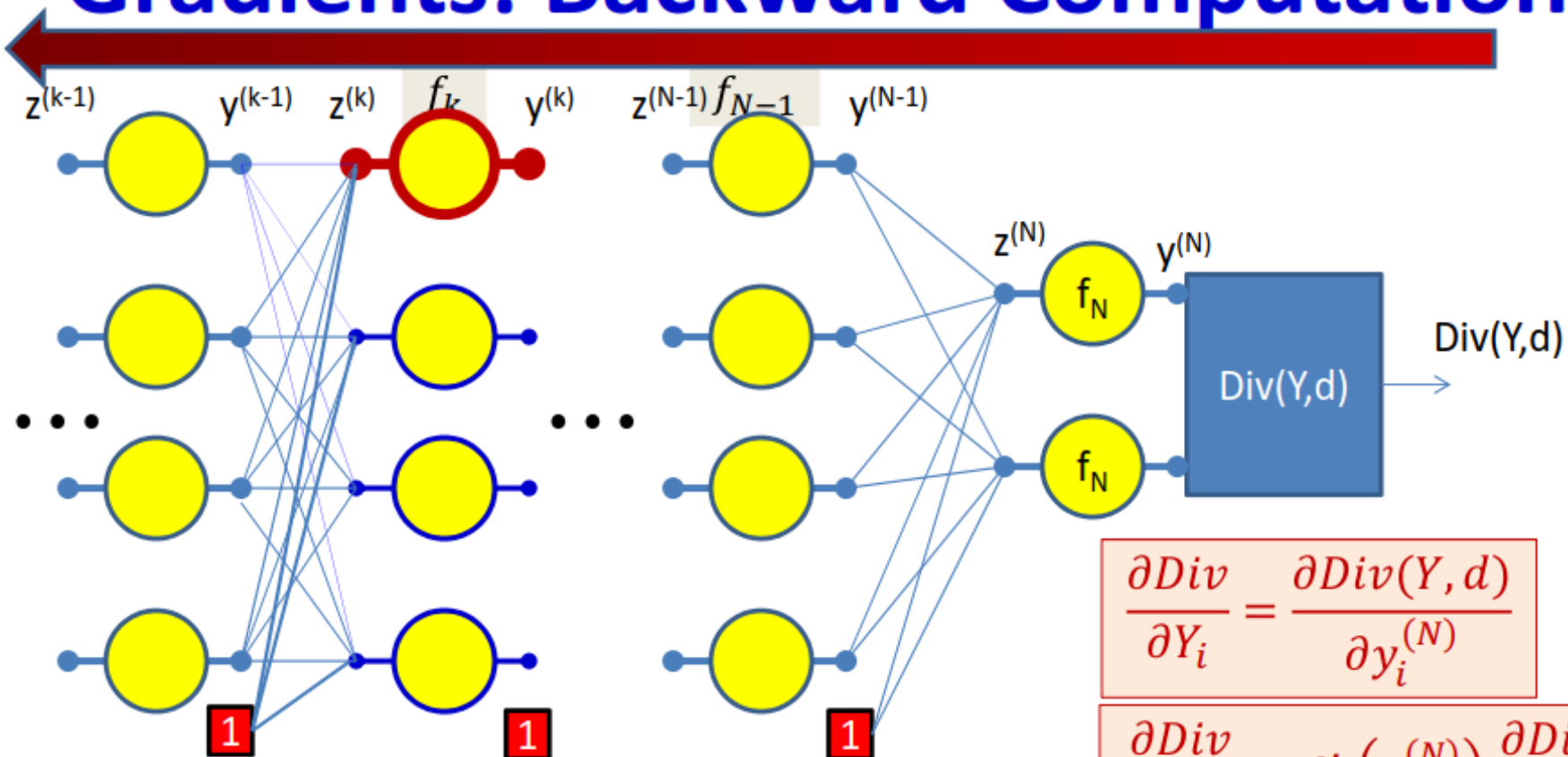
$$\frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} = w_{ij}^{(N)}$$

$$z_j^{[N]} = w^T y_i^{[N-1]}$$

But In this case the input is the output from previous layer



# Gradients: Backward Computation



computed during  
the forward pass

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial \text{Div}}{\partial y_i^{(k)}}$$

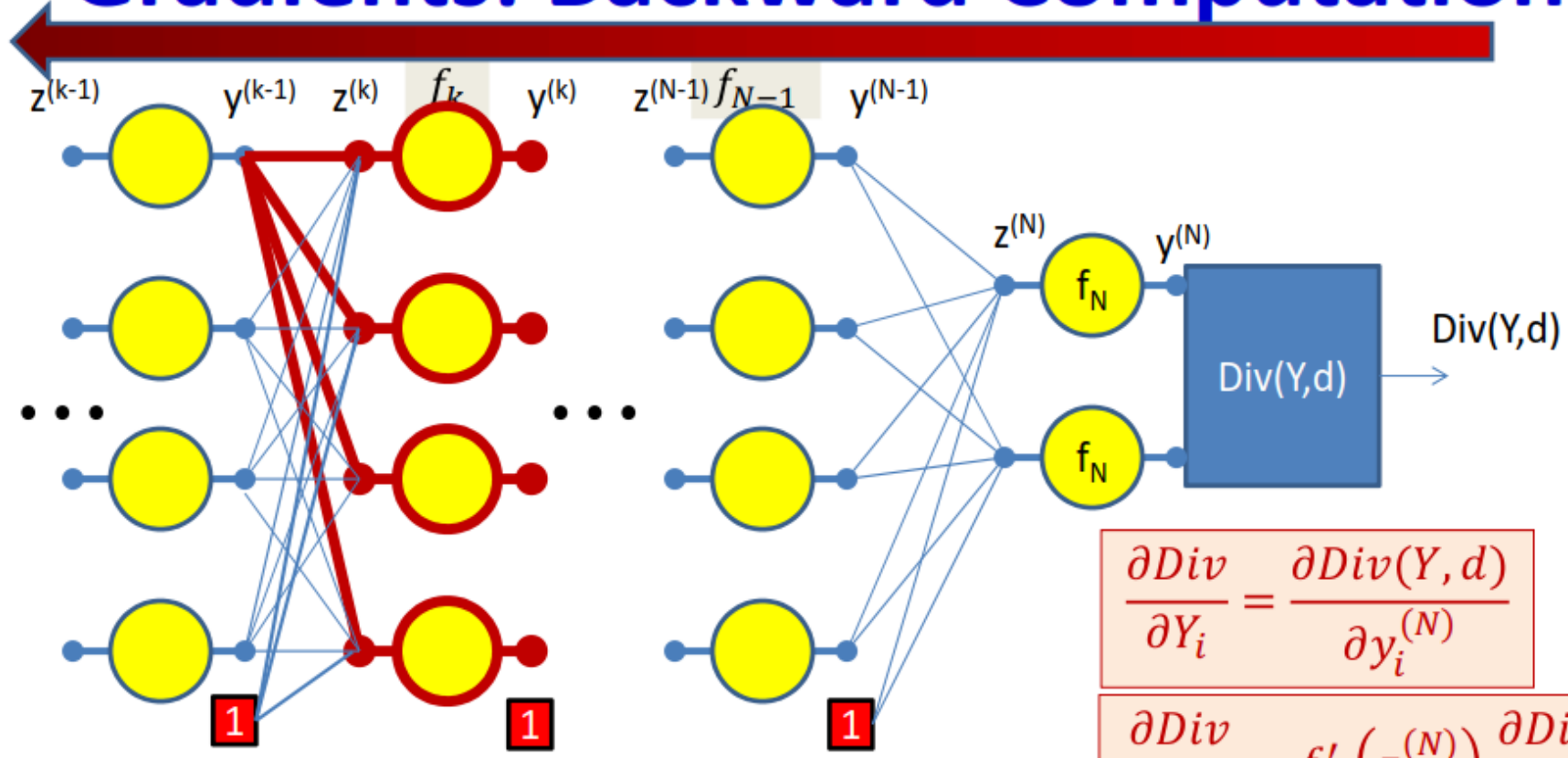
$$\frac{\partial \text{Div}}{\partial Y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(N-1)}} = \sum_j w_{ij}^{(N)} \frac{\partial \text{Div}}{\partial z_j^{(N)}}$$



# Gradients: Backward Computation

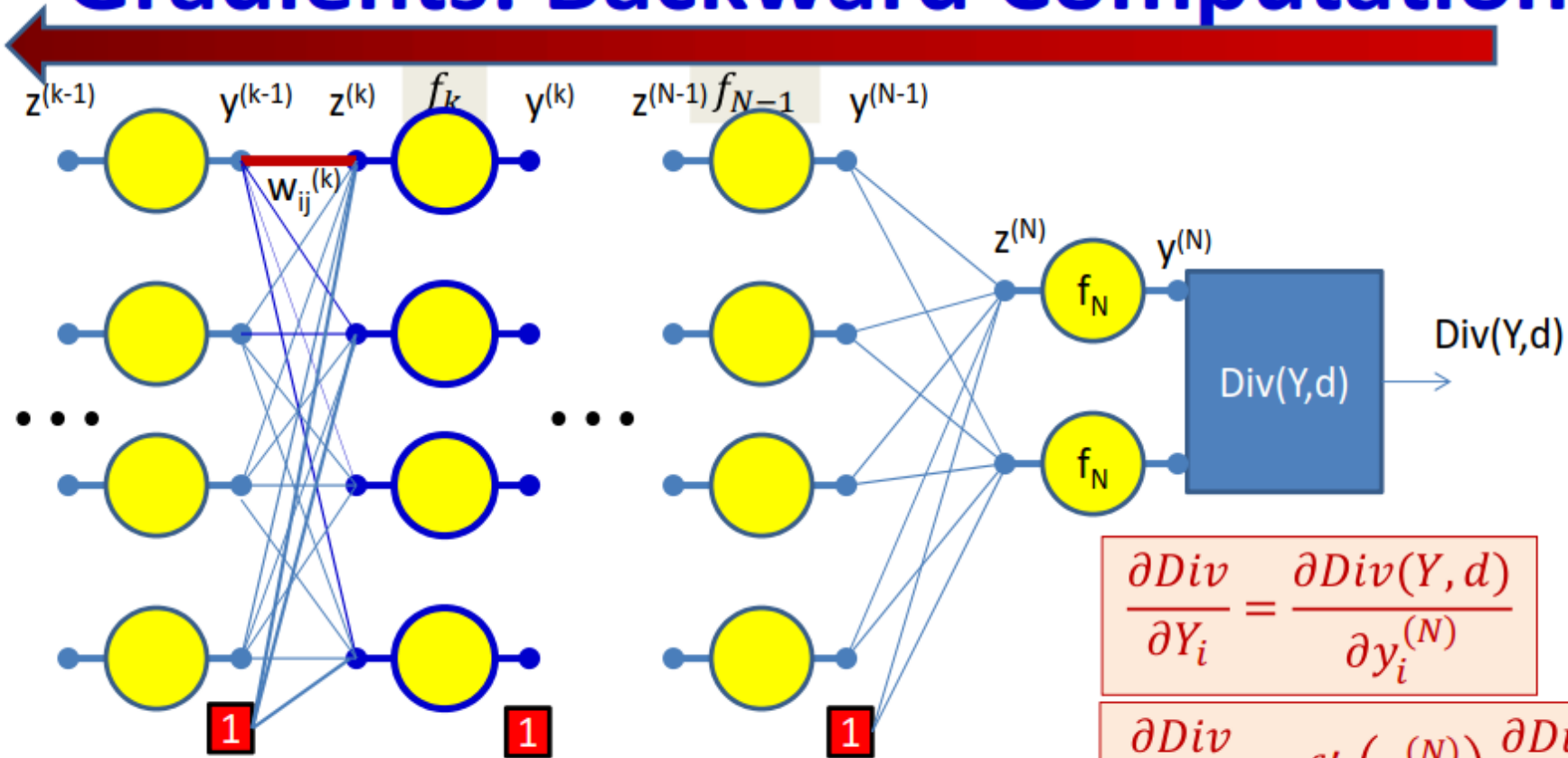


$$\frac{\partial Div}{\partial Y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

$$\frac{\partial Div}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j \frac{\partial z_j^{(k)}}{\partial y_i^{(k-1)}} \frac{\partial Div}{\partial z_j^{(k)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

# Gradients: Backward Computation



$$z_j^{[N]} = w^T y_i^{[N-1]}$$

$$\frac{\partial \text{Div}}{\partial w_{ij}^{(k)}} = \frac{\partial z_j^{(k)}}{\partial w_{ij}^{(k)}} \frac{\partial \text{Div}}{\partial z_j^{(k)}} = y_i^{(k-1)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$

$$\frac{\partial \text{Div}}{\partial Y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$



# Gradients: Backward Computation

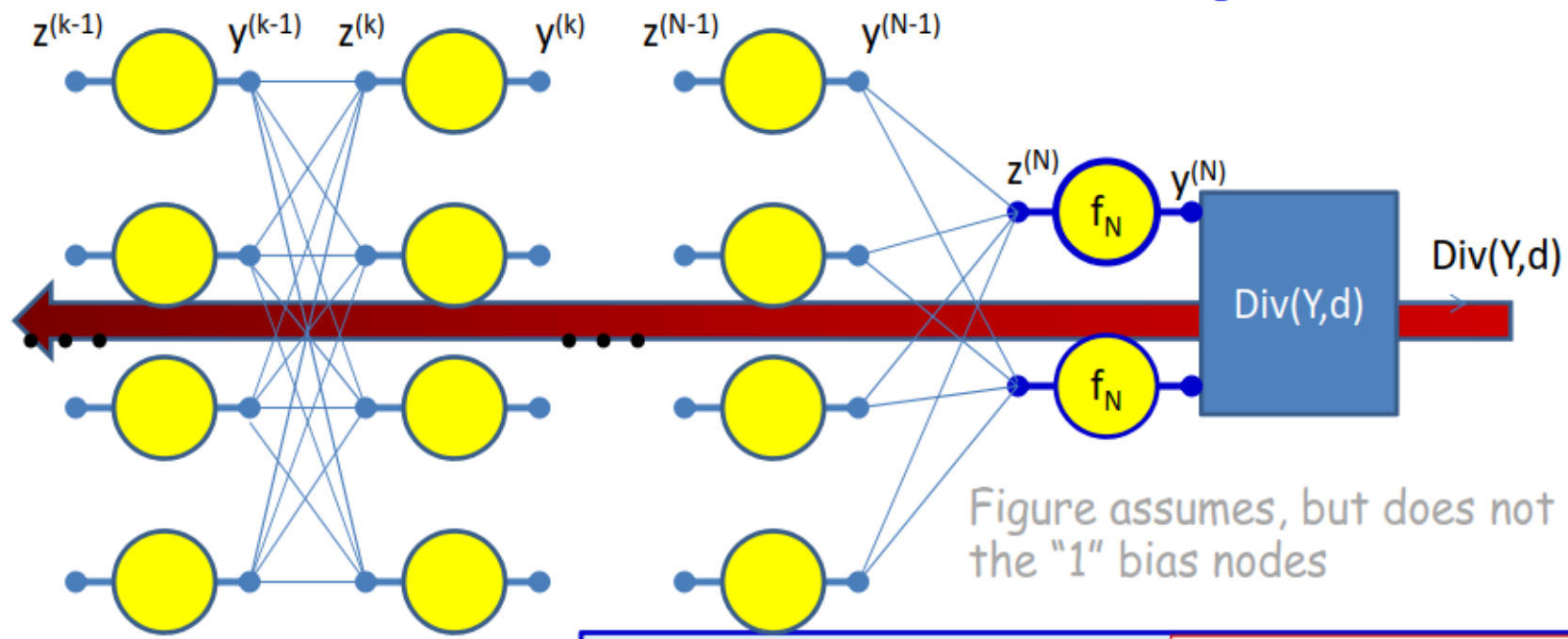


Figure assumes, but does not show the "1" bias nodes

Initialize: Gradient w.r.t network output

$$\frac{\partial Div}{\partial y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

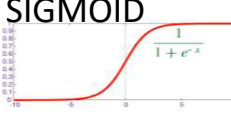
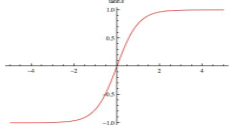
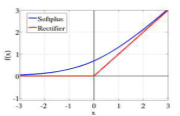


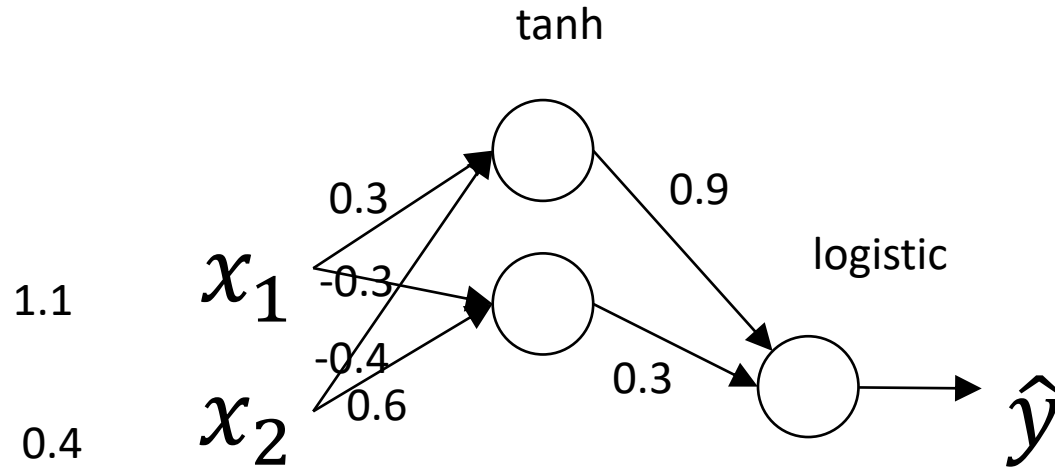
# Training by BackProp

- Initialize all weights  $(W^{(1)}, W^{(2)}, \dots, W^{(K)})$
- Do:
  - Initialize  $Err = 0$ ; For all  $i, j, k$ , initialize  $\frac{dErr}{dw_{i,j}^{(k)}} = 0$
  - For all  $t = 1:T$  (Loop over training instances)
    - **Forward pass:** Compute
      - Output  $Y_t$
      - $Err += Div(Y_t, d_t)$
    - **Backward pass:** For all  $i, j, k$ :
      - Compute  $\frac{dDiv(Y_t, d_t)}{dw_{i,j}^{(k)}}$
      - Compute  $\frac{dErr}{dw_{i,j}^{(k)}} += \frac{dDiv(Y_t, d_t)}{dw_{i,j}^{(k)}}$
  - For all  $i, j, k$ , update:
$$w_{i,j}^{(k)} = w_{i,j}^{(k)} - \frac{\eta}{T} \frac{dErr}{dw_{i,j}^{(k)}}$$
- Until  $Err$  has converged

# Exercise

## Activations and their derivatives

<p><b>SIGMOID</b></p>  <p><math>f(z) = \frac{1}{1 + \exp(-z)}</math></p>	<p><b>LOGISTIC FUNCTION</b></p> <p><math>f'(z) = f(z)(1 - f(z))</math></p>
 <p><math>f(z) = \tanh(z)</math></p>	<p><math>f'(z) = (1 - f^2(z))</math></p>
 <p><math>f(z) = \begin{cases} 0, &amp; z &lt; 0 \\ z, &amp; z \geq 0 \end{cases}</math></p> <p>softplus or SmoothReLU function</p> <p><math>f'(z) = \frac{1}{1 + \exp(-z)}</math></p>	<p>This space left intentionally (kind of) blank</p>



Expected output: 1

$$Div = \frac{(\hat{y}_i - y_i)^2}{2}$$

$$(y_i - \hat{y}_i)$$

?

## Gradients: Backward Computation

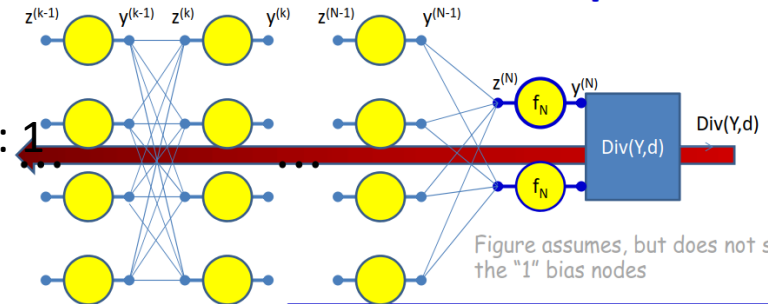


Figure assumes, but does not show the "1" bias nodes

Initialize: Gradient w.r.t network output

$$\frac{\partial Div}{\partial y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

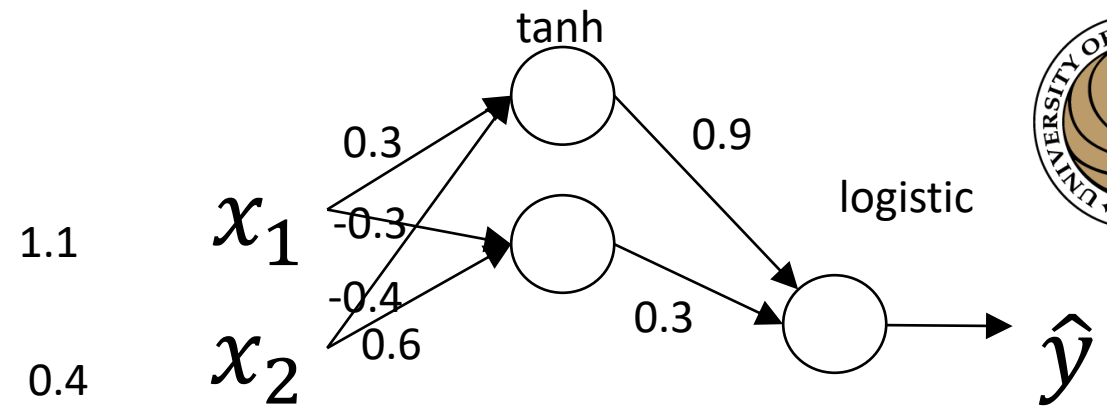
$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Forward



Expected output: 1  
 $Div = \frac{(\hat{y}_i - y_i)^2}{2}$

$$z_1^{[1]} = w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2$$

$$= 0.3 * 1.1 - 0.4 * 0.4$$

$$= 0.17$$

$$z_2^{[1]} = w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2$$

$$= -0.3 * 1.1 + 0.6 * 0.4$$

$$= -0.09$$

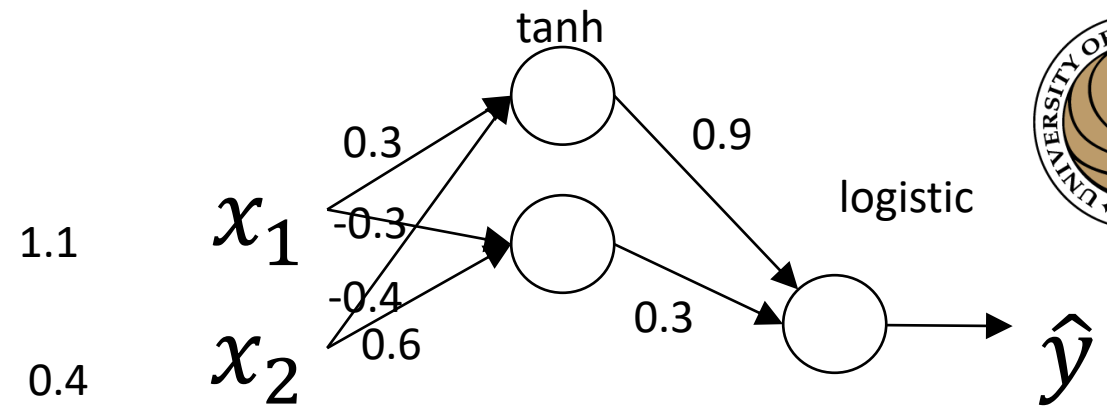
LAYER 1

$$y_1^{[1]} = \tanh(z_1^{[1]}) = \tanh(0.17) = 0.1683$$

$$y_2^{[1]} = \tanh(z_2^{[1]}) = \tanh(-0.09) = -0.0897$$



# Example: Forward



$$y_1^{[1]} = 0.1683$$

$$y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = w_{11}^{[2]} y_1^{[1]} + w_{21}^{[2]} y_2^{[1]}$$

$$= 0.9 * 0.1683 - 0.3 * 0.0897$$

$$= -0.124615$$

Expected output: 1

$$Div = \frac{(\hat{y}_i - y_i)^2}{2}$$

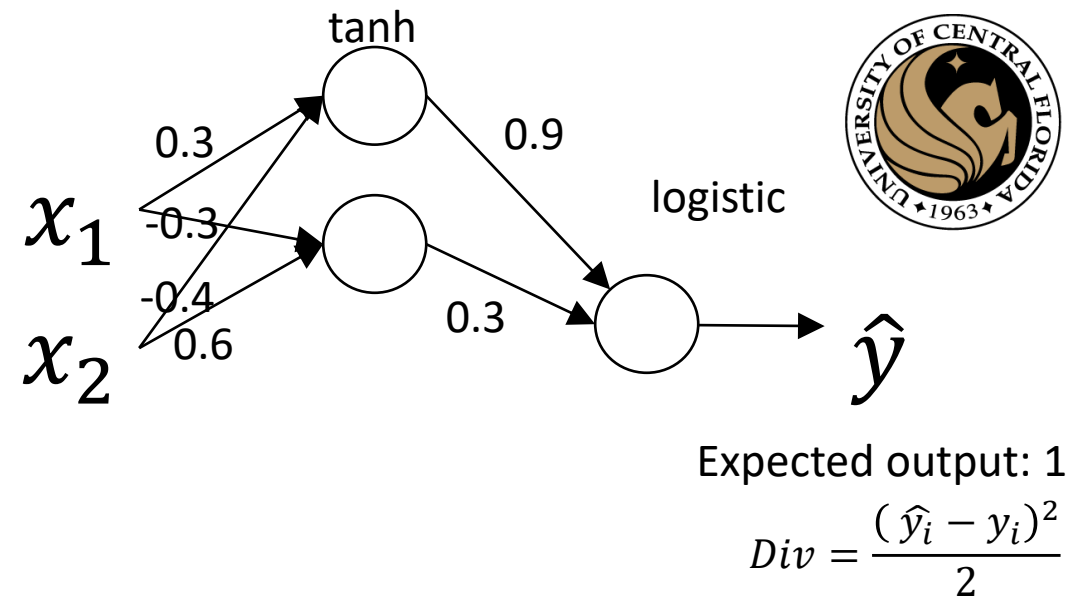
LAYER 2

$$y_1^{[2]} = \text{logistic}(z_1^{[2]}) = \frac{1}{1 + e^{-z_1^{[2]}}} = \frac{1}{1 + e^{0.1246}} = 0.531113$$

# Example: Backward

$z_1^{[1]} = 0.17$	$y_1^{[1]} = 0.1683$
$z_2^{[1]} = -0.019$	$y_2^{[1]} = -0.0897$
$z_1^{[2]} = 0.1246$	$y_1^{[2]} = 0.5311$

1.1  
0.4



INITIALIZE

$$Div = \frac{(\hat{y}_i - y_i)^2}{2}$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = (y_1^{[2]} - \hat{y})$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = (0.5311 - 1)$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



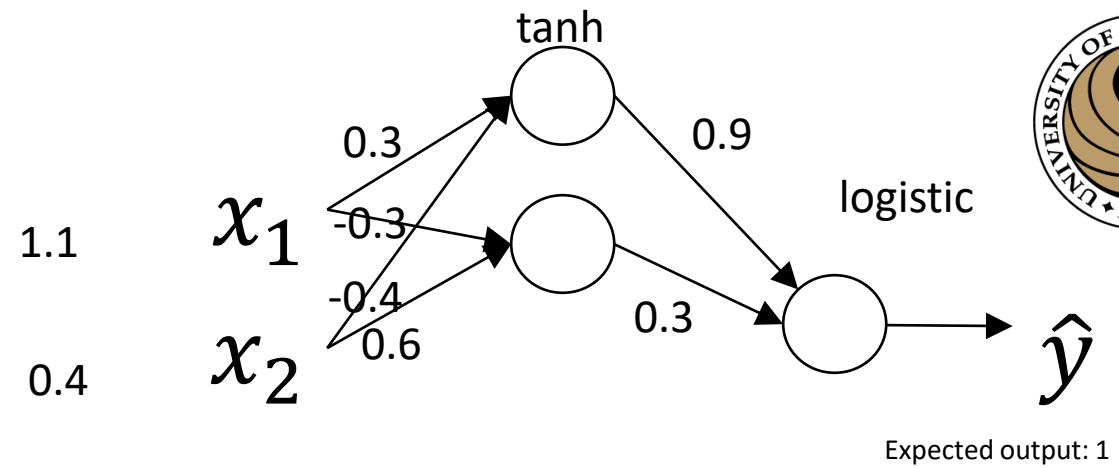
# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2'(z_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

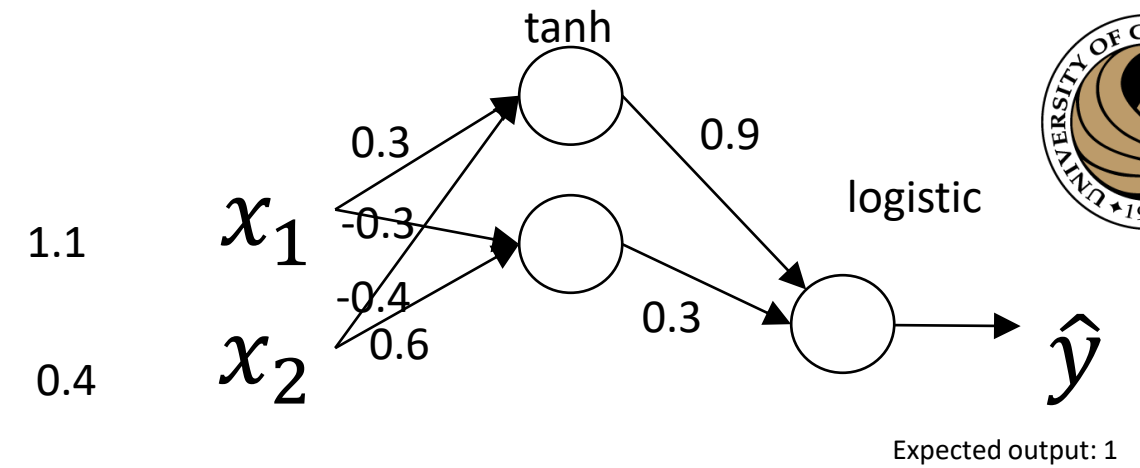
$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



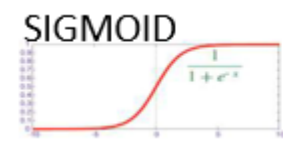
$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2'(z_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$


LOGISTIC FUNCTION

$$f(z) = \frac{1}{1 + \exp(-z)}$$

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

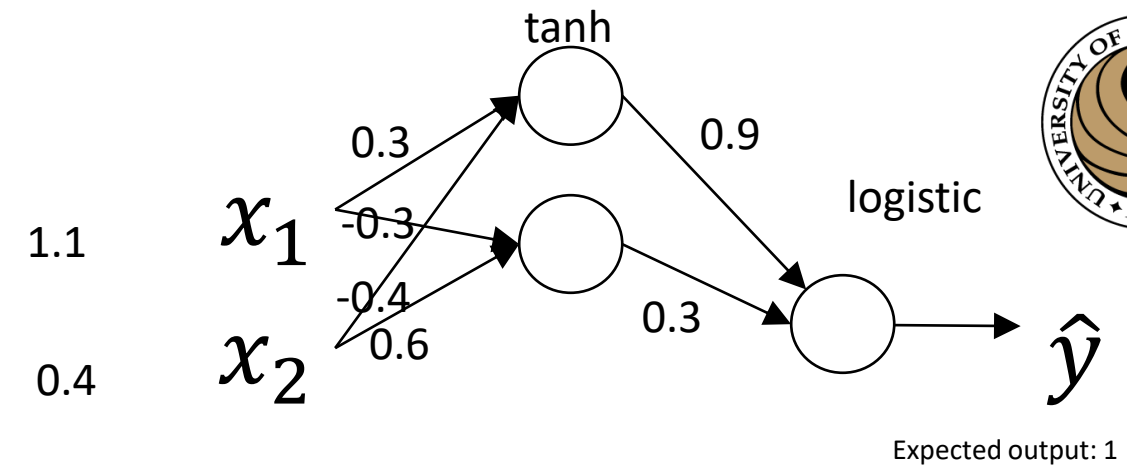




# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2'(z_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

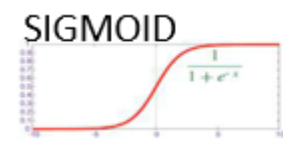
$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2(z_1^{[2]}) (1 - f_2(z_1^{[2]})) \frac{\partial Div}{\partial y_1^{[2]}}$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$


LOGISTIC FUNCTION

$$f(z) = \frac{1}{1 + \exp(-z)}$$

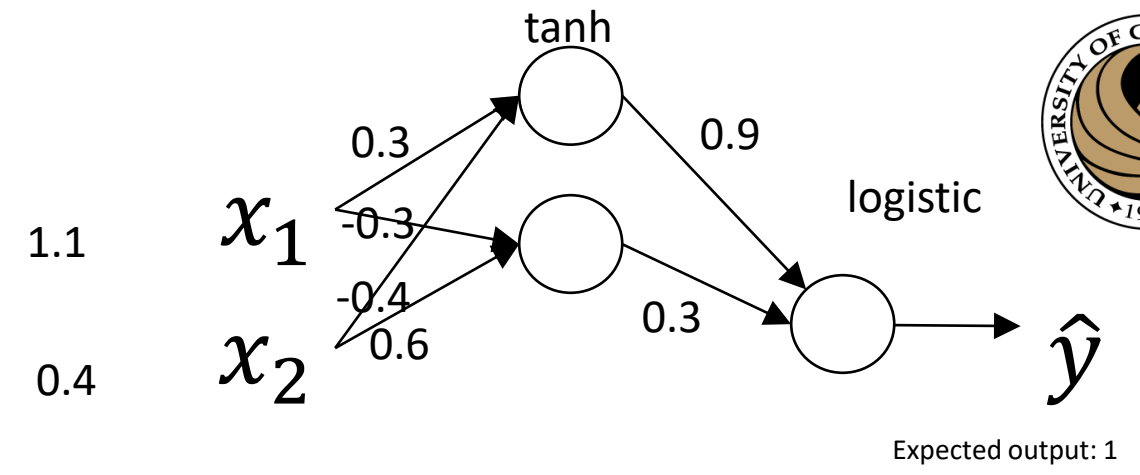
$$f'(z) = f(z)(1 - f(z))$$



# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2'(z_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2(z_1^{[2]}) (1 - f_2(z_1^{[2]})) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = y_1^{[2]} (1 - y_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

LAYER 2

K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

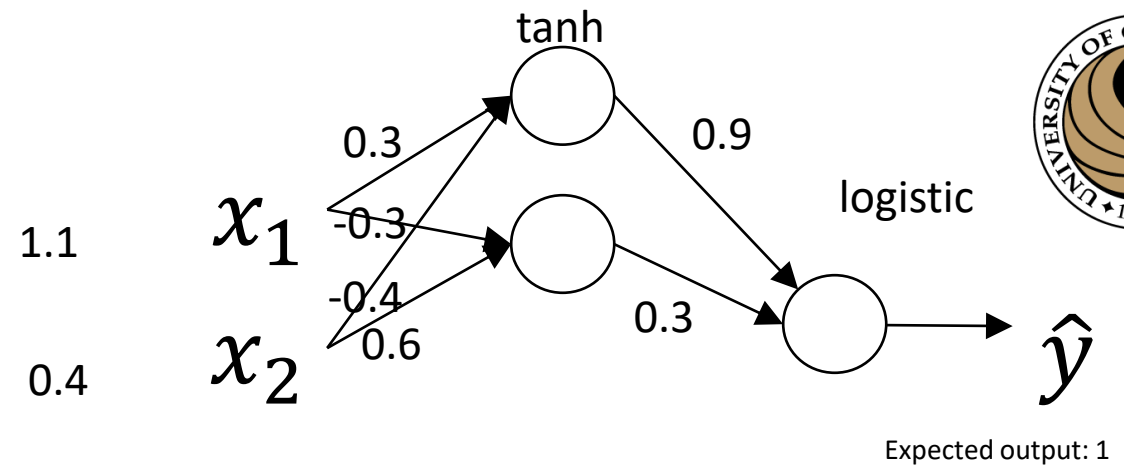
$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2'(z_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2(z_1^{[2]}) (1 - f_2(z_1^{[2]})) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = y_1^{[2]} (1 - y_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = 0.5311(1 - 0.5311)(-0.4688)$$

LAYER 2

K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

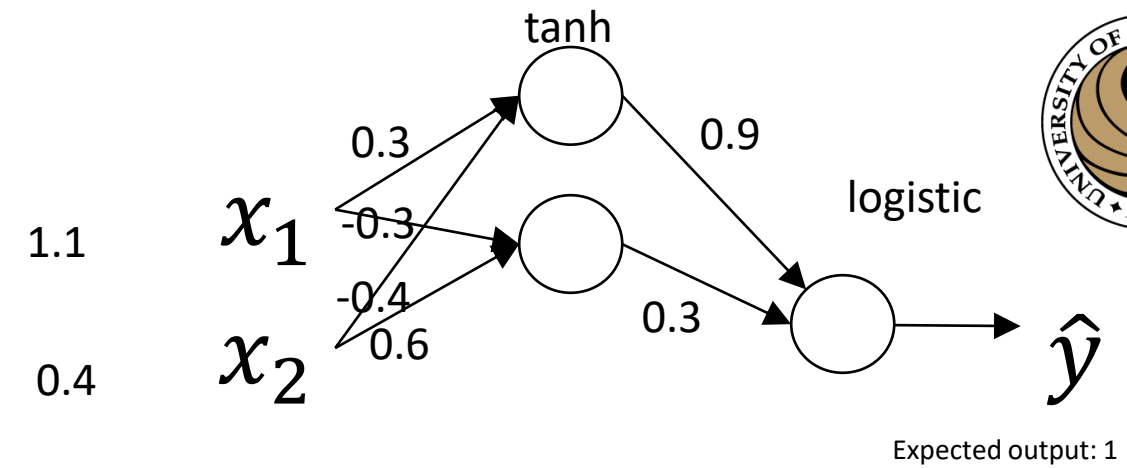
$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$



$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2'(z_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = f_2(z_1^{[2]}) (1 - f_2(z_1^{[2]})) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = y_1^{[2]} (1 - y_1^{[2]}) \frac{\partial Div}{\partial y_1^{[2]}}$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

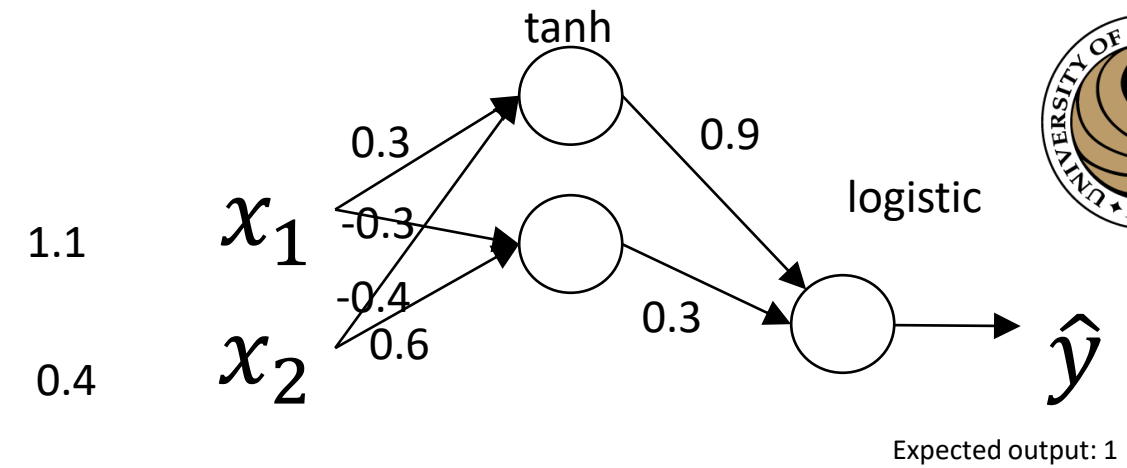
$$\begin{aligned}
 \frac{\partial Div}{\partial y_1^{[2]}} &= -0.4688 \\
 \frac{\partial Div}{\partial z_1^{[2]}} &= (-0.1167)
 \end{aligned}$$

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



$$\frac{\partial Div}{\partial y_1^{[1]}} = \sum_{j=1}^1 w_{1j}^{[2]} \frac{\partial Div}{\partial z_j^{[2]}}$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = w_{11}^{[2]} \frac{\partial Div}{\partial z_1^{[2]}}$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = 0.9(-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509$$

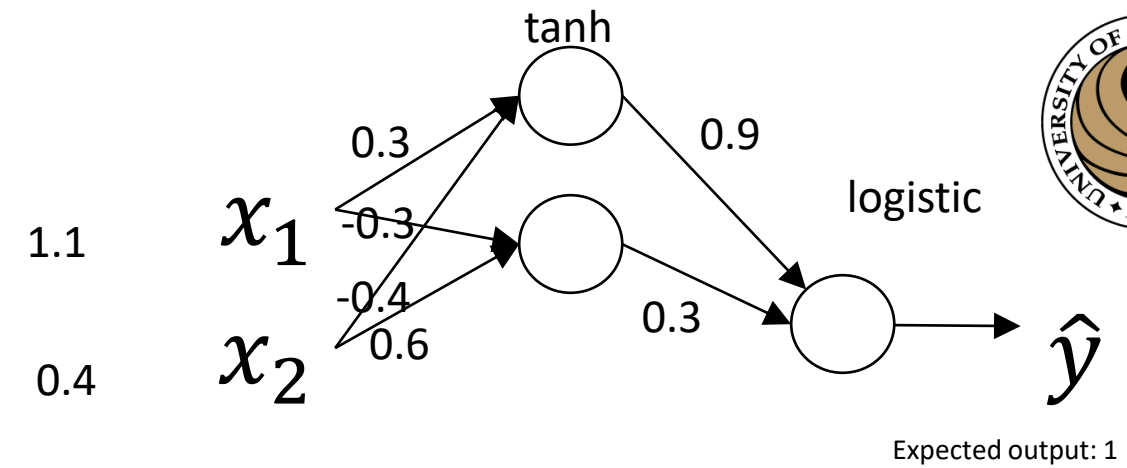
LAYER 2  
K=2



# Example: Backward

$$\begin{aligned}
 z_1^{[1]} &= 0.17 & y_1^{[1]} &= 0.1683 \\
 z_2^{[1]} &= -0.019 & y_2^{[1]} &= -0.0897 \\
 z_1^{[2]} &= 0.1246 & y_1^{[2]} &= 0.5311
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial Div}{\partial y_1^{[2]}} &= -0.4688 \\
 \frac{\partial Div}{\partial z_1^{[2]}} &= (-0.1167) \\
 \frac{\partial Div}{\partial y_1^{[1]}} &= -0.10509
 \end{aligned}$$



$$\frac{\partial Div}{\partial y_2^{[1]}} = \sum_{j=1}^1 w_{2j}^{[2]} \frac{\partial Div}{\partial z_j^{[2]}}$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = w_{21}^{[2]} \frac{\partial Div}{\partial z_1^{[2]}}$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = 0.3(-0.1167)$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

# Example: Backward



$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

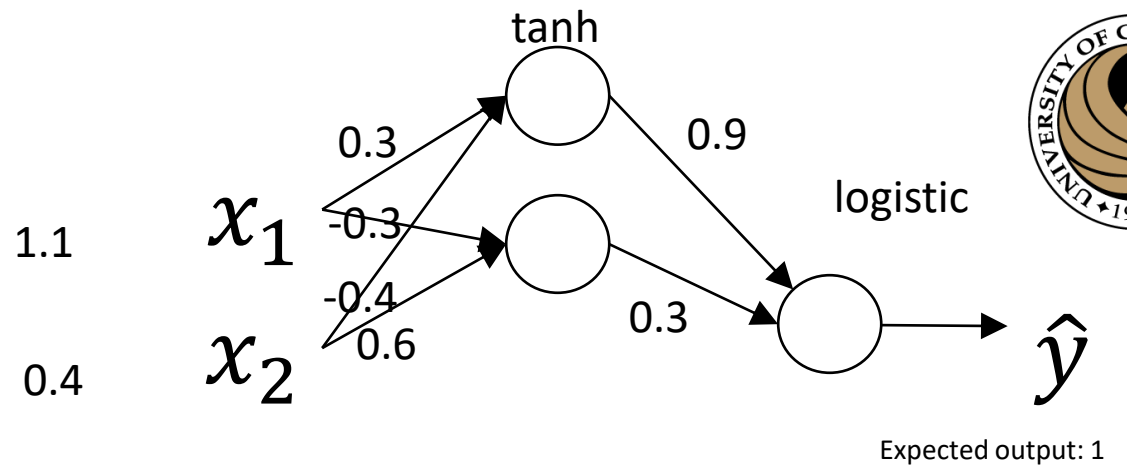
$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$



$$\frac{\partial Div}{\partial w_{11}^{[2]}} = y_1^{[1]} \frac{\partial Div}{\partial z_j^{[2]}}$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = 0.1683(-0.1167)$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

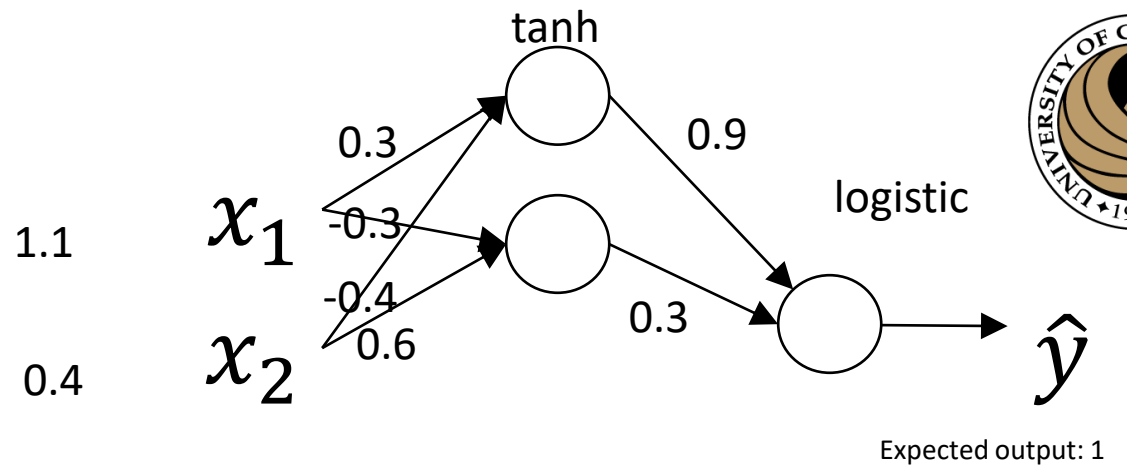
$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$



$$\frac{\partial Div}{\partial w_{21}^{[2]}} = y_2^{[1]} \frac{\partial Div}{\partial z_1^{[2]}}$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = (-0.0897)(-0.1167)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

LAYER 2  
K=2

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$





# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

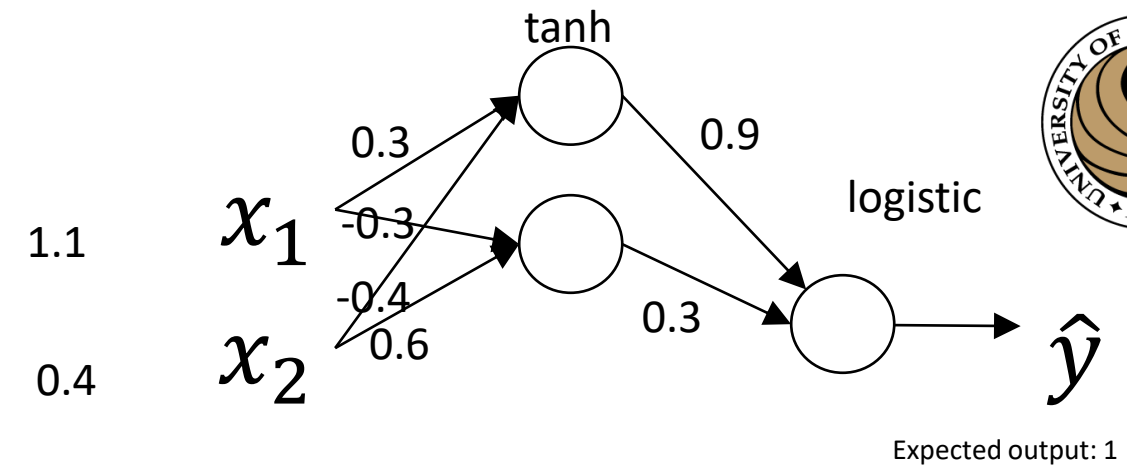
$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$



$$\frac{\partial Div}{\partial z_1^{[1]}} = f_1'(z_1^{[1]}) \frac{\partial Div}{\partial y_1^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = \left(1 - f_1^2(z_1^{[1]})\right) \frac{\partial Div}{\partial y_1^{[1]}}$$

LAYER 1  
K=1

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

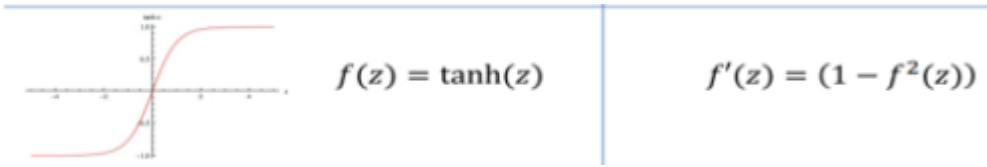
$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$





# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

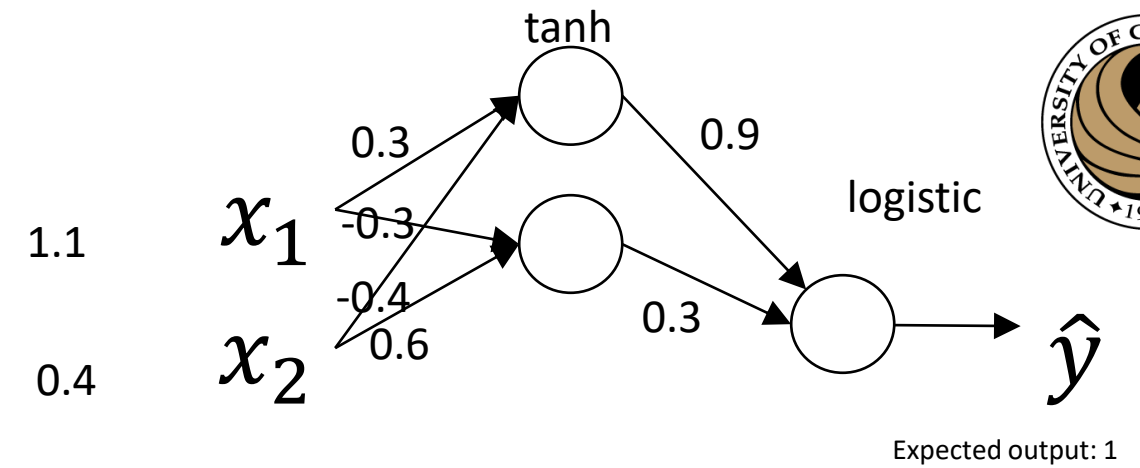
$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$



For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = f_1'(z_1^{[1]}) \frac{\partial Div}{\partial y_1^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = (1 - f_1^2(z_1^{[1]})) \frac{\partial Div}{\partial y_1^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = (1 - (y_1^{[1]})^2) \frac{\partial Div}{\partial y_1^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = (1 - 0.1683^2)(-0.10509)$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

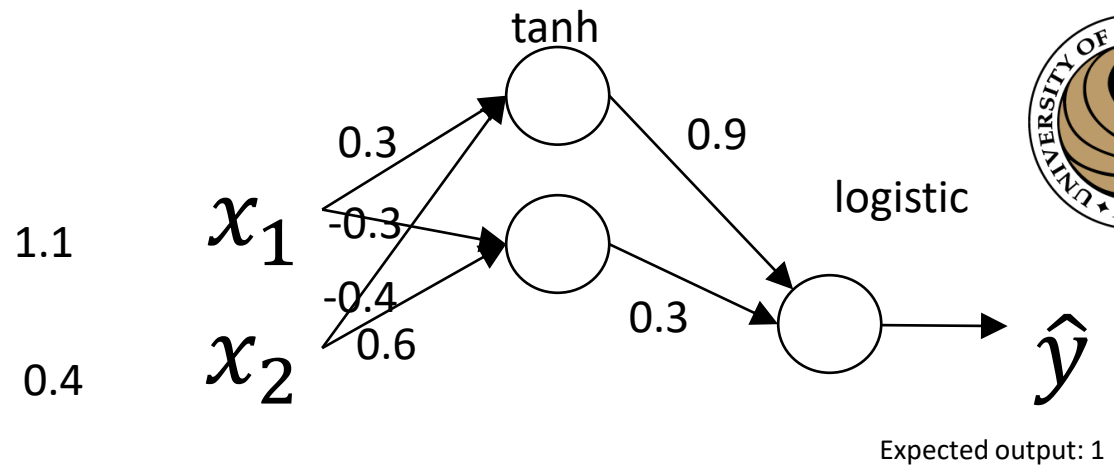
$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688$$

$$\frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509$$

$$\frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$



For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = f_1'(z_1^{[1]}) \frac{\partial Div}{\partial y_1^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = (1 - f_1^2(z_1^{[1]})) \frac{\partial Div}{\partial y_1^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = (1 - (y_1^{[1]})^2) \frac{\partial Div}{\partial y_1^{[1]}}$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

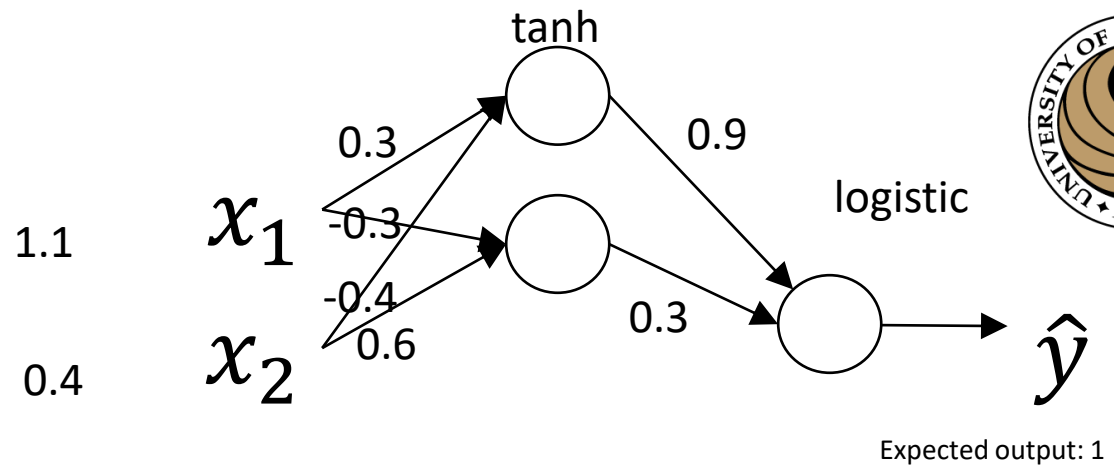
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021$$



$$\frac{\partial Div}{\partial z_2^{[1]}} = f_1'(z_2^{[1]}) \frac{\partial Div}{\partial y_2^{[1]}}$$

$$\frac{\partial Div}{\partial z_2^{[1]}} = (1 - f_1^2(z_2^{[1]})) \frac{\partial Div}{\partial y_2^{[1]}}$$

$$\frac{\partial Div}{\partial z_2^{[1]}} = (1 - (y_2^{[1]})^2) \frac{\partial Div}{\partial y_2^{[1]}}$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = (1 - 0.0897^2)(-0.03503)$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

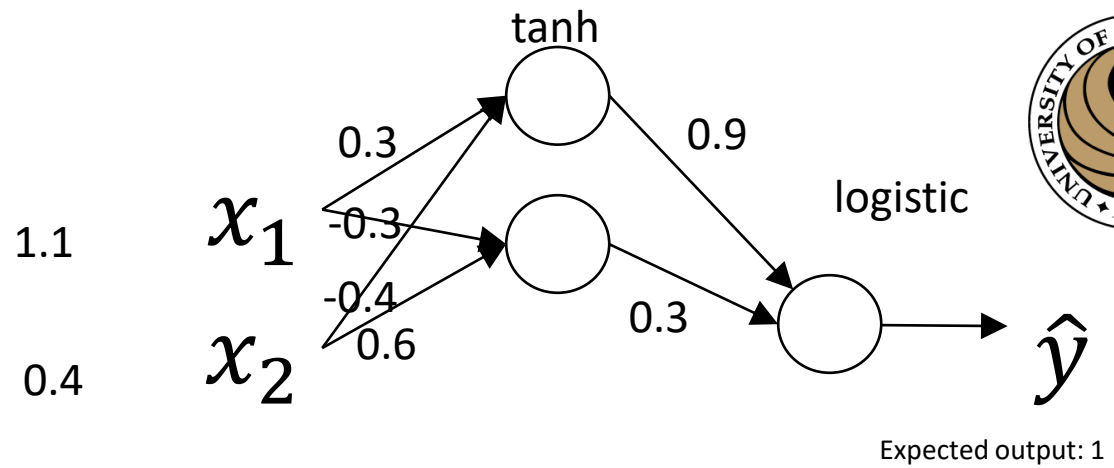
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021$$



$$\frac{\partial Div}{\partial z_2^{[1]}} = f_1'(z_2^{[1]}) \frac{\partial Div}{\partial y_2^{[1]}}$$

$$\frac{\partial Div}{\partial z_2^{[1]}} = \left(1 - f_1^2(z_2^{[1]})\right) \frac{\partial Div}{\partial y_2^{[1]}}$$

$$\frac{\partial Div}{\partial z_2^{[1]}} = \left(1 - (y_2^{[1]})^2\right) \frac{\partial Div}{\partial y_2^{[1]}}$$

$$\frac{\partial Div}{\partial z_2^{[1]}} = -0.0347$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f_k'(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

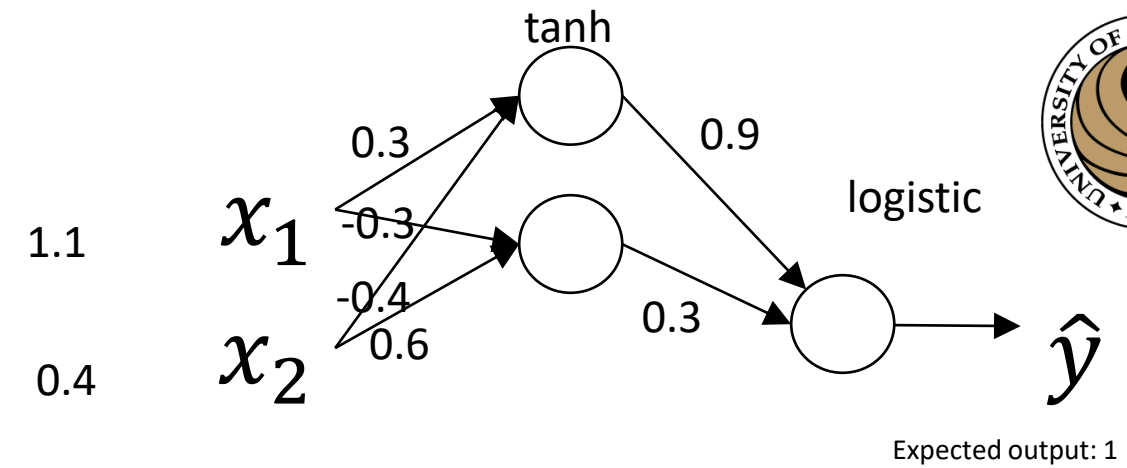
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021 \quad \frac{\partial Div}{\partial z_2^{[1]}} = -0.0347$$



$$\frac{\partial Div}{\partial w_{11}^{[1]}} = y_1^{[0]} \frac{\partial Div}{\partial z_1^{[1]}}$$

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = x_1 \frac{\partial Div}{\partial z_1^{[1]}}$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = 1.1(-0.1021)$$

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = -0.1123$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

# Example: Backward



$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

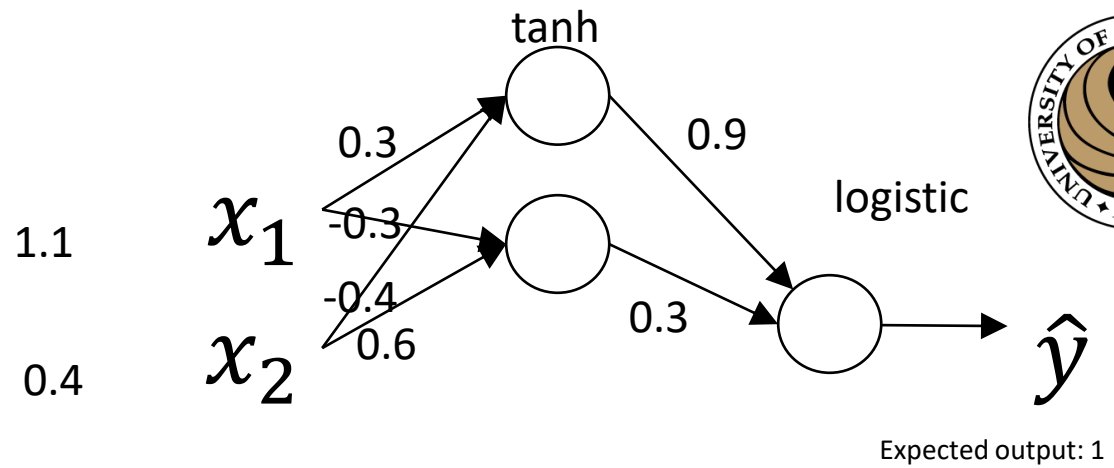
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021 \quad \frac{\partial Div}{\partial z_2^{[1]}} = -0.0347$$



$$\frac{\partial Div}{\partial w_{12}^{[1]}} = y_1^{[0]} \frac{\partial Div}{\partial z_2^{[1]}}$$

$$\frac{\partial Div}{\partial w_{12}^{[1]}} = x_1 \frac{\partial Div}{\partial z_2^{[1]}}$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{12}^{[1]}} = 1.1(-0.0347)$$

$$\frac{\partial Div}{\partial w_{12}^{[1]}} = -0.03822$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966)$$

$$\frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = -0.1123$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

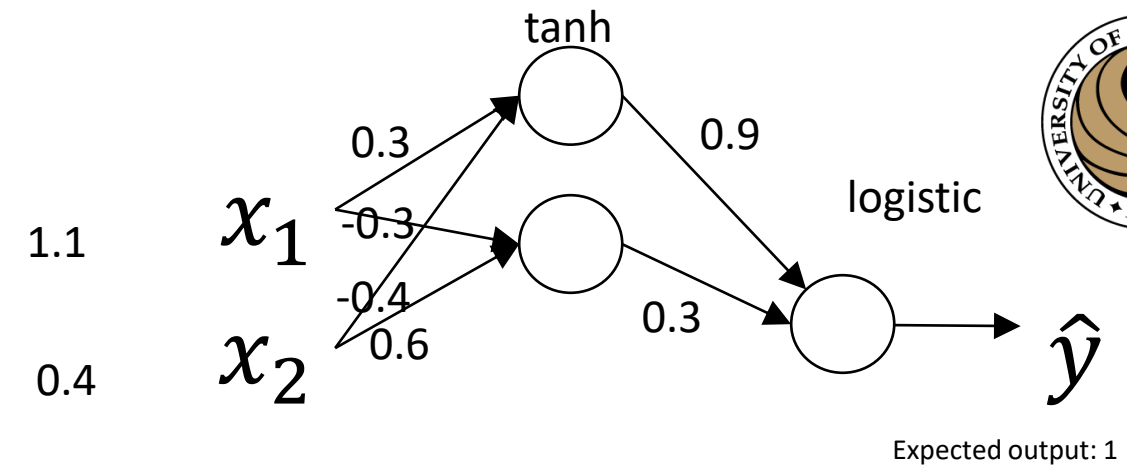
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021 \quad \frac{\partial Div}{\partial z_2^{[1]}} = -0.0347$$



$$\frac{\partial Div}{\partial w_{21}^{[1]}} = y_2^{[0]} \frac{\partial Div}{\partial z_1^{[1]}}$$

$$\frac{\partial Div}{\partial w_{21}^{[1]}} = x_2 \frac{\partial Div}{\partial z_1^{[1]}}$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{21}^{[1]}} = 0.4(-0.1021)$$

$$\frac{\partial Div}{\partial w_{21}^{[1]}} = -0.04084$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966) \quad \frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = -0.1123 \quad \frac{\partial Div}{\partial w_{12}^{[1]}} = -0.03822$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$





# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

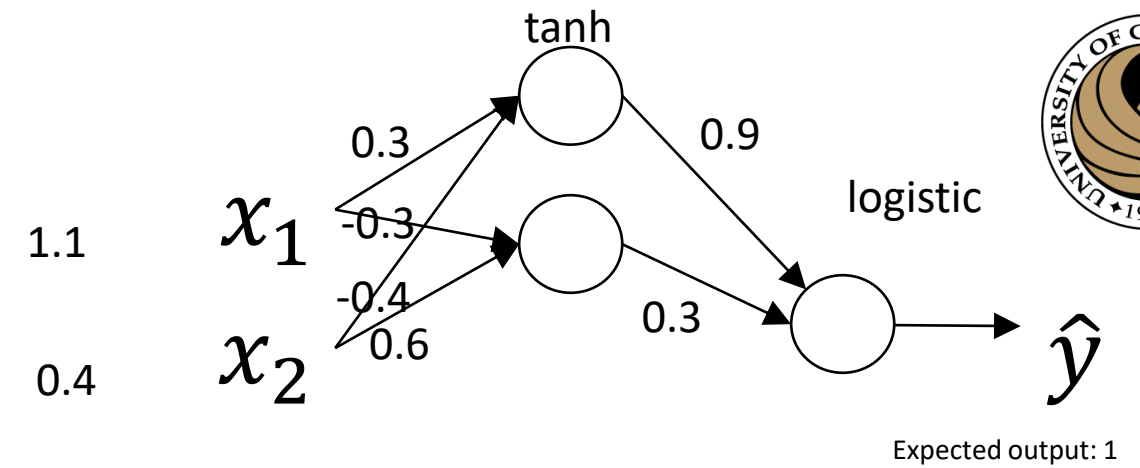
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021 \quad \frac{\partial Div}{\partial z_2^{[1]}} = -0.0347$$



$$\frac{\partial Div}{\partial w_{22}^{[1]}} = y_2^{[0]} \frac{\partial Div}{\partial z_2^{[1]}}$$

$$\frac{\partial Div}{\partial w_{22}^{[1]}} = x_2 \frac{\partial Div}{\partial z_2^{[1]}}$$

LAYER 1  
K=1

$$\frac{\partial Div}{\partial w_{22}^{[1]}} = 0.4(-0.0347)$$

$$\frac{\partial Div}{\partial w_{22}^{[1]}} = -0.013899$$

$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966) \quad \frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = -0.1123 \quad \frac{\partial Div}{\partial w_{12}^{[1]}} = -0.03822$$

$$\frac{\partial Div}{\partial w_{21}^{[1]}} = -0.04084$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Example: Backward

$$z_1^{[1]} = 0.17 \quad y_1^{[1]} = 0.1683$$

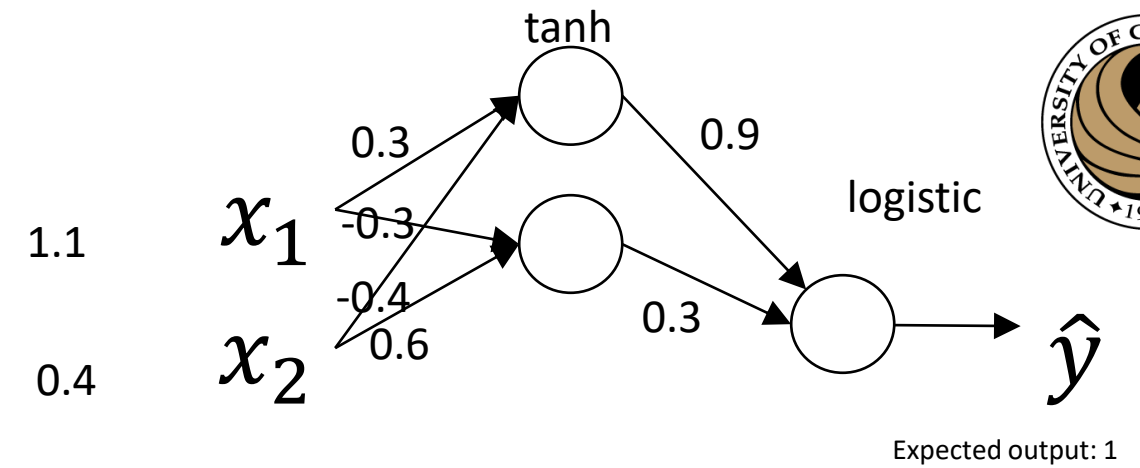
$$z_2^{[1]} = -0.019 \quad y_2^{[1]} = -0.0897$$

$$z_1^{[2]} = 0.1246 \quad y_1^{[2]} = 0.5311$$

$$\frac{\partial Div}{\partial y_1^{[2]}} = -0.4688 \quad \frac{\partial Div}{\partial z_1^{[2]}} = (-0.1167)$$

$$\frac{\partial Div}{\partial y_1^{[1]}} = -0.10509 \quad \frac{\partial Div}{\partial y_2^{[1]}} = -0.03503$$

$$\frac{\partial Div}{\partial z_1^{[1]}} = -0.1021 \quad \frac{\partial Div}{\partial z_2^{[1]}} = -0.0347$$



$$\frac{\partial Div}{\partial w_{11}^{[2]}} = (-0.01966) \quad \frac{\partial Div}{\partial w_{21}^{[2]}} = 0.010481$$

$$\frac{\partial Div}{\partial w_{11}^{[1]}} = -0.1123 \quad \frac{\partial Div}{\partial w_{12}^{[1]}} = -0.03822$$

$$\frac{\partial Div}{\partial w_{21}^{[1]}} = -0.04084 \quad \frac{\partial Div}{\partial w_{22}^{[1]}} = -0.013899$$

LAYER 1  
K=1

For  $k = N..1$   
 For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$



# Softmax

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

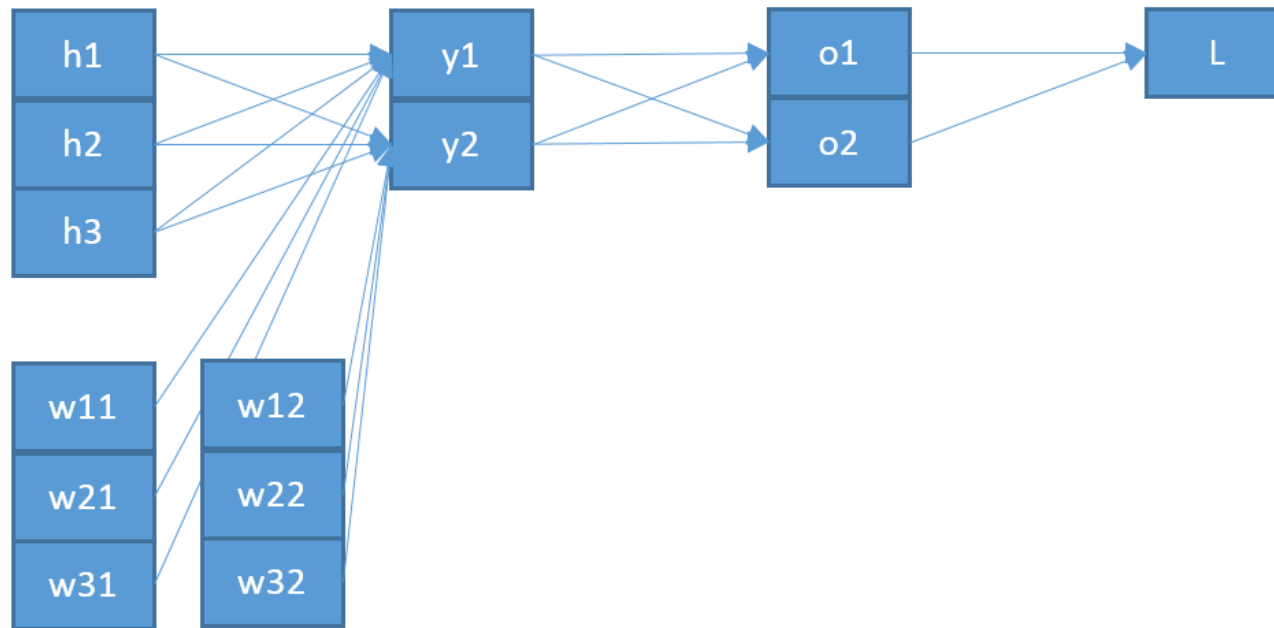
Used to interpret outputs as probabilities

$\vec{z}$	The input vector to the softmax function, made up of (z0, ... zK)
$z_i$	All the $z_i$ values are the elements of the input vector to the softmax function, and they can take any real value, positive, zero or negative. For example a neural network could have output a vector such as (-0.62, 8.12, 2.53), which is not a valid probability distribution, hence why the softmax would be necessary.
$e^{z_i}$	The standard exponential function is applied to each element of the input vector. This gives a positive value above 0, which will be very small if the input was negative, and very large if the input was large. However, it is still not fixed in the range (0, 1) which is what is required of a probability.
$\sum_{j=1}^K e^{z_j}$	The term on the bottom of the formula is the normalization term. It ensures that all the output values of the function will sum to 1 and each be in the range (0, 1), thus constituting a valid probability distribution.
$K$	The number of classes in the multi-class classifier.

$$\begin{aligned} \begin{bmatrix} P(\text{cat}) \\ P(\text{dog}) \end{bmatrix} &= \sigma\left(\begin{bmatrix} 1.2 \\ 0.3 \end{bmatrix}\right) \\ &= \begin{bmatrix} \frac{e^{1.2}}{e^{1.2} + e^{0.3}} \\ \frac{e^{0.3}}{e^{1.2} + e^{0.3}} \end{bmatrix} \\ &= \begin{bmatrix} 0.71 \\ 0.29 \end{bmatrix} \end{aligned}$$

[The Softmax function and its derivative - Eli Bendersky's website \(thegreenplace.net\)](http://thegreenplace.net)

# Backpropagation with Softmax / Cross Entropy



# Backpropagation with Softmax / Cross Entropy

$$L = -t_1 \log o_1 - t_2 \log o_2$$

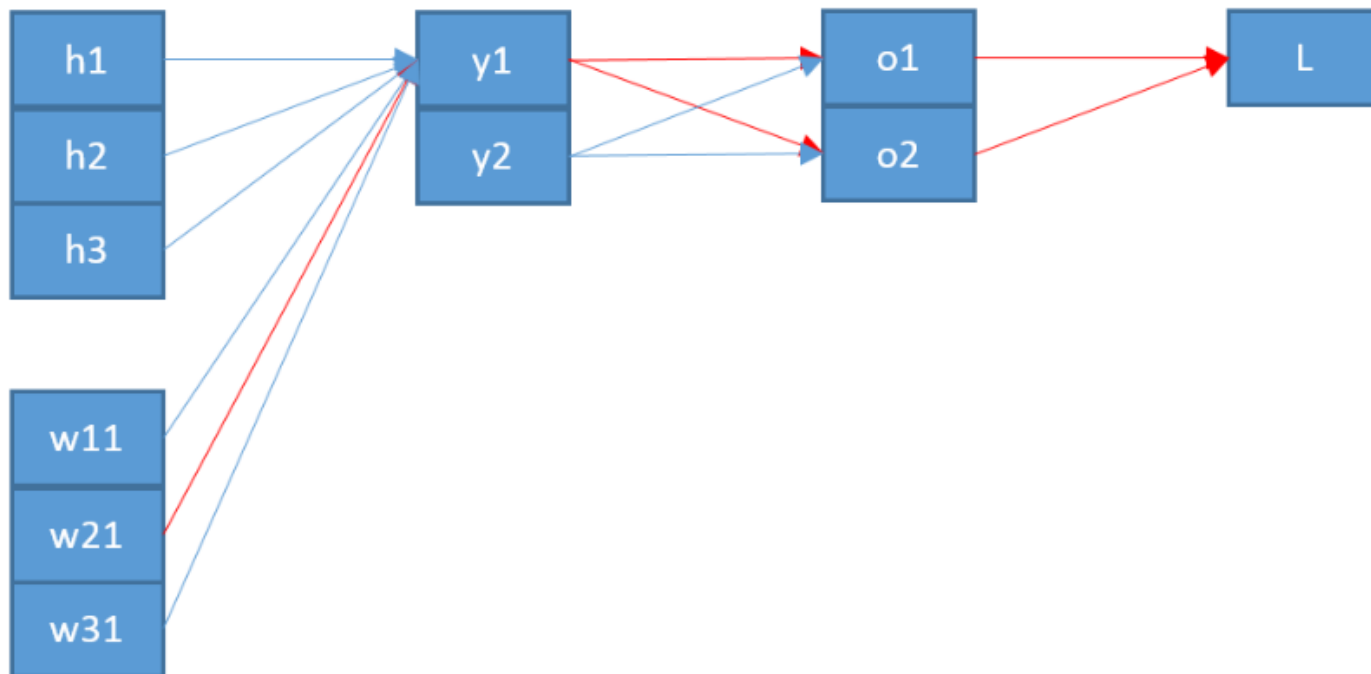
$$o_1 = \frac{\exp(y_1)}{\exp(y_1) + \exp(y_2)}$$

$$o_2 = \frac{\exp(y_2)}{\exp(y_1) + \exp(y_2)}$$

$$y_1 = w_{11}h_1 + w_{21}h_2 + w_{31}h_3$$

$$y_2 = w_{12}h_1 + w_{22}h_2 + w_{32}h_3$$

Say I want to calculate the derivative of the loss with respect to  $w_{21}$ . I can just use my picture to trace back the path from the loss to the weight I'm interested in (removed the second column of  $w$ 's for clarity):





# Backpropagation with Softmax / Cross Entropy

$$\frac{\partial L}{\partial o_1} = -\frac{t_1}{o_1}$$

$$\frac{\partial L}{\partial o_2} = -\frac{t_2}{o_2}$$

$$\frac{\partial o_1}{\partial y_1} = \frac{\exp(y_1)}{\exp(y_1) + \exp(y_2)} - \left( \frac{\exp(y_1)}{\exp(y_1) + \exp(y_2)} \right)^2 = o_1(1 - o_1)$$

$$\frac{\partial o_2}{\partial y_1} = \frac{-\exp(y_2) \exp(y_1)}{(\exp(y_1) + \exp(y_2))^2} = -o_2 o_1$$

$$\frac{\partial y_1}{\partial w_{21}} = h_2$$

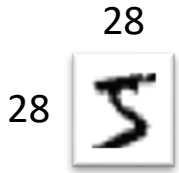
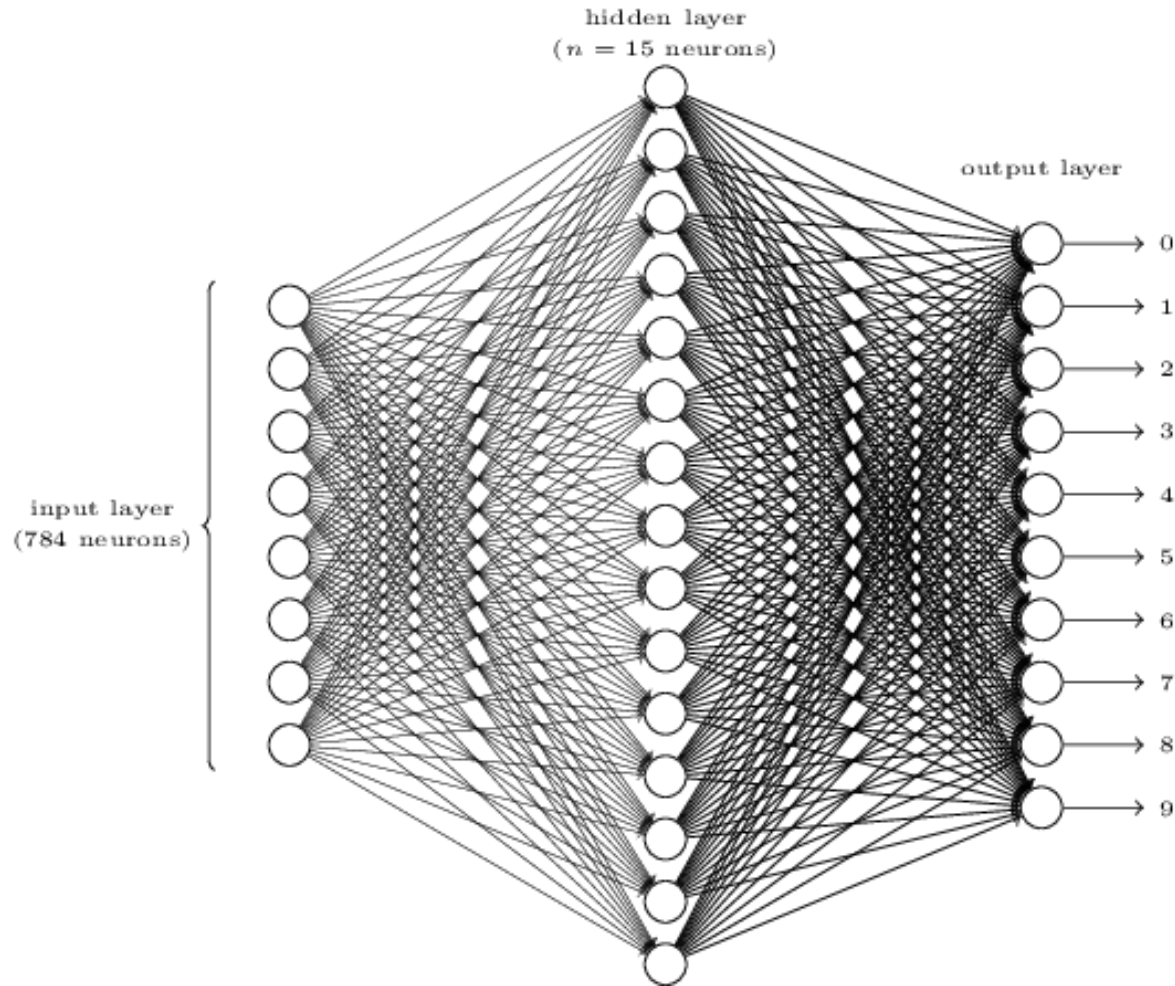
Finally, putting the chain rule together:

$$\begin{aligned} \frac{\partial L}{\partial w_{21}} &= \frac{\partial L}{\partial o_1} \frac{\partial o_1}{\partial y_1} \frac{\partial y_1}{\partial w_{21}} + \frac{\partial L}{\partial o_2} \frac{\partial o_2}{\partial y_1} \frac{\partial y_1}{\partial w_{21}} \\ &= \frac{-t_1}{o_1} [o_1(1 - o_1)] h_2 + \frac{-t_2}{o_2} (-o_2 o_1) h_2 \\ &= h_2(t_2 o_1 - t_1 + t_1 o_1) \\ &= h_2(o_1(t_1 + t_2) - t_1) \\ &= h_2(o_1 - t_1) \end{aligned}$$

Note that in the last step,  $t_1 + t_2 = 1$  because the vector  $\mathbf{t}$  is a one-hot vector.

# A real example

# Digit classification



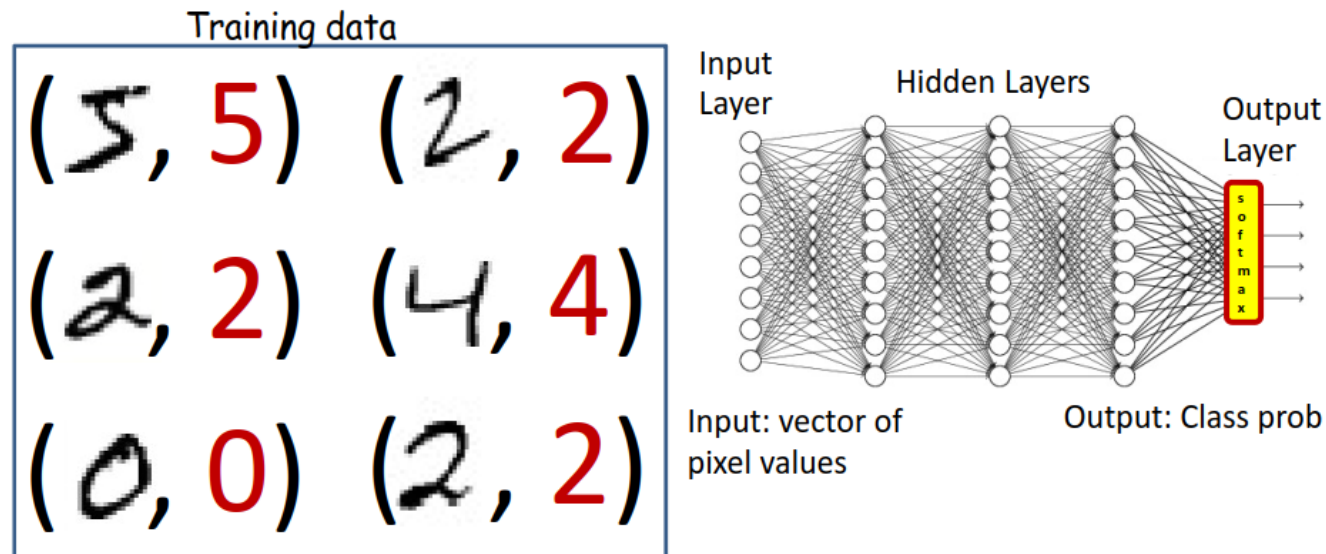
- MNIST dataset:
  - 70000 grayscale images of digits scanned.
  - 60000 for training
  - 10000 for testing
- Loss function

$$J_2(w) = \frac{1}{m} \sum_{train} (\hat{y}_i - y_i)^2$$



# Digit classification

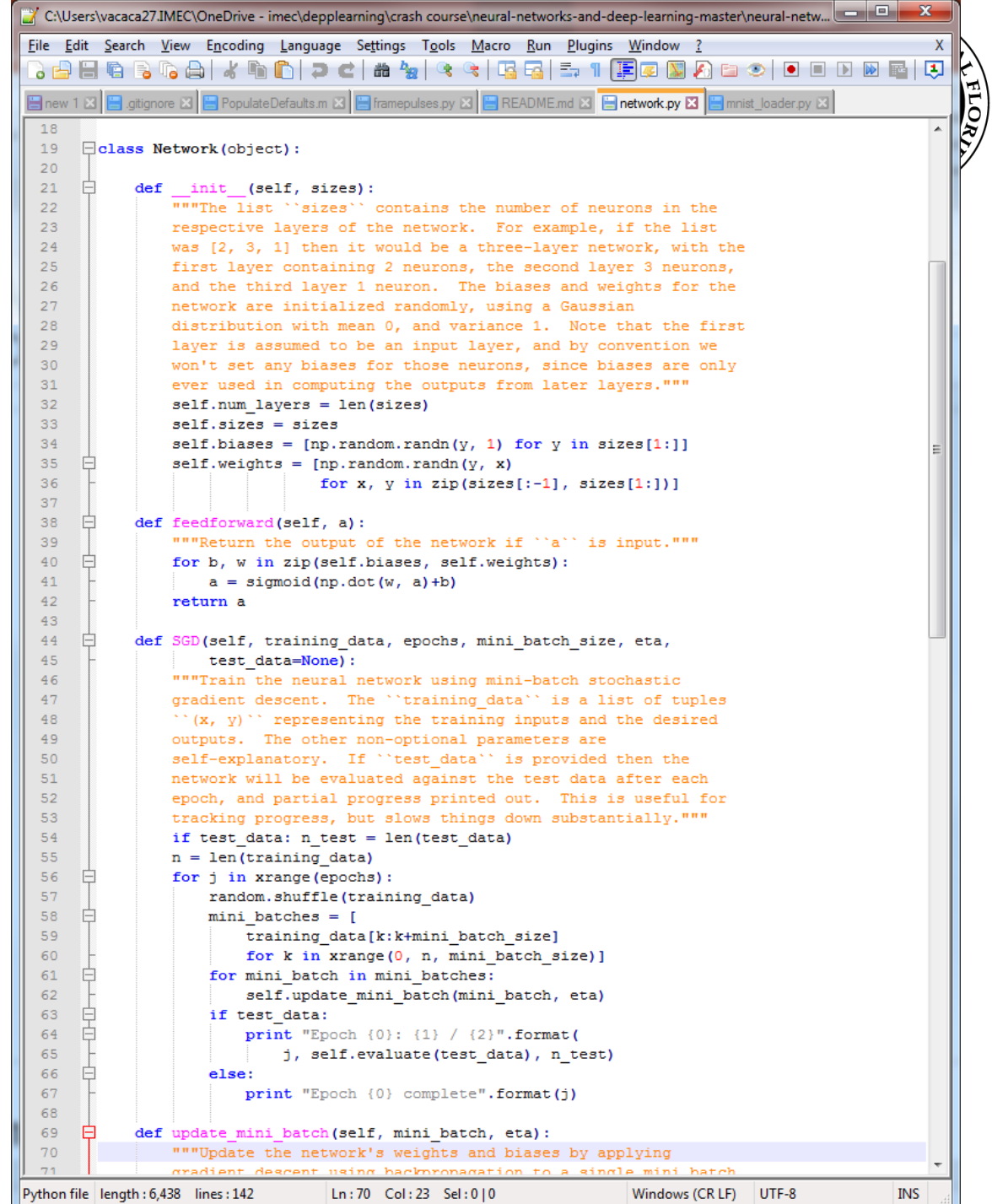
## Typical Problem statement: multiclass classification



- Given, many positive and negative examples (training data),
  - learn all weights such that the network does the desired job

# A look in the code

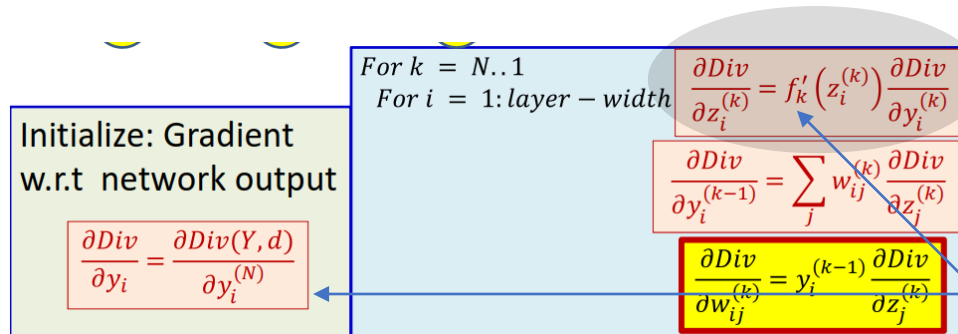
- To run this code do:
  - import network
  - net = network.Network([784, 30, 10])
  - net.SGD(training\_data, 30, 10, 3.0, test\_data=test\_data)



```
18
19 class Network(object):
20
21     def __init__(self, sizes):
22         """The list ``sizes`` contains the number of neurons in the
23         respective layers of the network. For example, if the list
24         was [2, 3, 1] then it would be a three-layer network, with the
25         first layer containing 2 neurons, the second layer 3 neurons,
26         and the third layer 1 neuron. The biases and weights for the
27         network are initialized randomly, using a Gaussian
28         distribution with mean 0, and variance 1. Note that the first
29         layer is assumed to be an input layer, and by convention we
30         won't set any biases for those neurons, since biases are only
31         ever used in computing the outputs from later layers."""
32         self.num_layers = len(sizes)
33         self.sizes = sizes
34         self.biases = [np.random.randn(y, 1) for y in sizes[1:]]
35         self.weights = [np.random.randn(y, x)
36                         for x, y in zip(sizes[:-1], sizes[1:])]
37
38     def feedforward(self, a):
39         """Return the output of the network if ``a`` is input."""
40         for b, w in zip(self.biases, self.weights):
41             a = sigmoid(np.dot(w, a)+b)
42         return a
43
44     def SGD(self, training_data, epochs, mini_batch_size, eta,
45            test_data=None):
46         """Train the neural network using mini-batch stochastic
47         gradient descent. The ``training_data`` is a list of tuples
48         ``(x, y)`` representing the training inputs and the desired
49         outputs. The other non-optional parameters are
50         self-explanatory. If ``test_data`` is provided then the
51         network will be evaluated against the test data after each
52         epoch, and partial progress printed out. This is useful for
53         tracking progress, but slows things down substantially."""
54         if test_data: n_test = len(test_data)
55         n = len(training_data)
56         for j in xrange(epochs):
57             random.shuffle(training_data)
58             mini_batches = [
59                 training_data[k:k+mini_batch_size]
60                 for k in xrange(0, n, mini_batch_size)]
61             for mini_batch in mini_batches:
62                 self.update_mini_batch(mini_batch, eta)
63             if test_data:
64                 print "Epoch {0}: {1} / {2}".format(
65                     j, self.evaluate(test_data), n_test)
66             else:
67                 print "Epoch {0} complete".format(j)
68
69     def update_mini_batch(self, mini_batch, eta):
70         """Update the network's weights and biases by applying
71         gradient descent using backpropagation to a single mini batch
```

Python file length: 6,438 lines: 142 Ln: 70 Col: 23 Sel: 0 | 0 Windows (CR LF) UTF-8 INS

# A look in code



```

93
94
95 def backprop(self, x, y):
96     """Return a tuple ``(nabla_b, nabla_w)`` representing the
97     gradient for the cost function C_x. ``nabla_b`` and
98     ``nabla_w`` are layer-by-layer lists of numpy arrays, similar
99     to ``self.biases`` and ``self.weights``."""
100     nabla_b = [np.zeros(b.shape) for b in self.biases]
101     nabla_w = [np.zeros(w.shape) for w in self.weights]
102     # feedforward
103     activation = x
104     activations = [x] # list to store all the activations, layer by layer
105     zs = [] # list to store all the z vectors, layer by layer
106     for b, w in zip(self.biases, self.weights):
107         z = np.dot(w, activation)+b
108         zs.append(z)
109         activation = sigmoid(z)
110         activations.append(activation)
111     # backward pass
112     delta = self.cost_derivative(activations[-1], y) * \
113             sigmoid_prime(zs[-1])
114     nabla_b[-1] = delta
115     nabla_w[-1] = np.dot(delta, activations[-2].transpose())
116     # Note that the variable l in the loop below is used a little
117     # differently to the notation in Chapter 2 of the book. Here,
118     # l = 1 means the last layer of neurons, l = 2 is the
119     # second-last layer, and so on. It's a renumbering of the
120     # scheme in the book, used here to take advantage of the fact
121     # that Python can use negative indices in lists.
122     for l in xrange(2, self.num_layers):
123         z = zs[-l]
124         sp = sigmoid_prime(z)
125         delta = np.dot(self.weights[-l+1].transpose(), delta) * sp
126         nabla_b[-l] = delta
127         nabla_w[-l] = np.dot(delta, activations[-l-1].transpose())
128     return (nabla_b, nabla_w)
129
130 def cost_derivative(self, output_activations, y):
131     """Return the vector of partial derivatives \partial C_x /
132     \partial a for the output activations."""
133     return (output_activations-y)
134
135 ##### Miscellaneous functions
136 def sigmoid(z):
137     """The sigmoid function."""
138     return 1.0/(1.0+np.exp(-z))
139
140 def sigmoid_prime(z):
141     """Derivative of the sigmoid function."""
142     return sigmoid(z)*(1-sigmoid(z))

```

Python file length: 6,439 lines: 142 Ln: 140 Col: 5 Sel: 0 | 0 Windows (CR LF) UTF-8 INS

# A look in code

Initialize: Gradient  
w.r.t network output

$$\frac{\partial Div}{\partial y_i} = \frac{\partial Div(Y, d)}{\partial y_i^{(N)}}$$

For  $k = N..1$   
For  $i = 1: \text{layer} - \text{width}$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\frac{\partial Div}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial Div}{\partial z_j^{(k)}}$$

$$\frac{\partial Div}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial Div}{\partial z_j^{(k)}}$$

```
*C:\Users\vacaca27\IMEC\OneDrive - imec\deplearning\crash course\neural-networks-and-deep-learning-master\neural-net...
File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?
new 1 x | gitignore x | PopulateDefaults.m x | framepulses.py x | README.md x | network.py x | mnist_loader.py x
93
94 def backprop(self, x, y):
95     """Return a tuple ``(nabla_b, nabla_w)`` representing the
96     gradient for the cost function C_x. ``nabla_b`` and
97     ``nabla_w`` are layer-by-layer lists of numpy arrays, similar
98     to ``self.biases`` and ``self.weights``."""
99     nabla_b = [np.zeros(b.shape) for b in self.biases]
100    nabla_w = [np.zeros(w.shape) for w in self.weights]
101    # feedforward
102    activation = x
103    activations = [x] # list to store all the activations, layer by layer
104    zs = [] # list to store all the z vectors, layer by layer
105    for b, w in zip(self.biases, self.weights):
106        z = np.dot(w, activation)+b
107        zs.append(z)
108        activation = sigmoid(z)
109        activations.append(activation)
110    # backward pass
111    delta = self.cost_derivative(activations[-1], y) * \
112            sigmoid_prime(zs[-1])
113    nabla_b[-1] = delta
114    nabla_w[-1] = np.dot(delta, activations[-2].transpose())
115    # Note that the variable l in the loop below is used a little
116    # differently to the notation in Chapter 2 of the book. Here,
117    # l = 1 means the last layer of neurons, l = 2 is the
118    # second-last layer, and so on. It's a renumbering of the
119    # scheme in the book, used here to take advantage of the fact
120    # that Python can use negative indices in lists.
121    for l in xrange(2, self.num_layers):
122        z = zs[-l]
123        sp = sigmoid_prime(z)
124        delta = np.dot(self.weights[-l+1].transpose(), delta) * sp
125        nabla_b[-l] = delta
126        nabla_w[-l] = np.dot(delta, activations[-l-1].transpose())
127    return (nabla_b, nabla_w)
128
129 def cost_derivative(self, output_activations, y):
130     """Return the vector of partial derivatives \partial C_x /
131     \partial a for the output activations."""
132     return (output_activations-y)
133
134 ##### Miscellaneous functions
135 def sigmoid(z):
136     """The sigmoid function."""
137     return 1.0/(1.0+np.exp(-z))
138
139 def sigmoid_prime(z):
140     """Derivative of the sigmoid function."""
141     return sigmoid(z)*(1-sigmoid(z))
142
Python file length: 6,439 lines: 142 Ln: 140 Col: 5 Sel: 0 | 0 Windows (CR LF) UTF-8 INS
```

# A look in the code

random Initialization

Feed forward 'a' thru all the layers

A Epoch is when all the training data has been used to update weights

A minibatch is a subset of all the data used to obtain a 'quick' weight updates

If there is test data perform evaluation

```
File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?
new 1 x gitignore x PopulateDefaults.m x framepulses.py x README.md x network.py x mnist_loader.py x
18
19 class Network(object):
20
21     def __init__(self, sizes):
22         """The list ``sizes`` contains the number of neurons in the
23         respective layers of the network. For example, if the list
24         was [2, 3, 1] then it would be a three-layer network, with the
25         first layer containing 2 neurons, the second layer 3 neurons,
26         and the third layer 1 neuron. The biases and weights for the
27         network are initialized randomly, using a Gaussian
28         distribution with mean 0, and variance 1. Note that the first
29         layer is assumed to be an input layer, and by convention we
30         won't set any biases for those neurons, since biases are only
31         ever used in computing the outputs from later layers."""
32         self.num_layers = len(sizes)
33         self.sizes = sizes
34         self.biases = [np.random.randn(y, 1) for y in sizes[1:]]
35         self.weights = [np.random.randn(y, x)
36                         for x, y in zip(sizes[:-1], sizes[1:])]
37
38     def feedforward(self, a):
39         """Return the output of the network if ``a`` is input."""
40         for b, w in zip(self.biases, self.weights):
41             a = sigmoid(np.dot(w, a)+b)
42         return a
43
44     def SGD(self, training_data, epochs, mini_batch_size, eta,
45            test_data=None):
46         """Train the neural network using mini-batch stochastic
47         gradient descent. The ``training_data`` is a list of tuples
48         ``(x, y)`` representing the training inputs and the desired
49         outputs. The other non-optional parameters are
50         self-explanatory. If ``test_data`` is provided then the
51         network will be evaluated against the test data after each
52         epoch, and partial progress printed out. This is useful for
53         tracking progress, but slows things down substantially."""
54         if test_data: n_test = len(test_data)
55         n = len(training_data)
56         for j in xrange(epochs):
57             random.shuffle(training_data)
58             mini_batches = [
59                 training_data[k:k+mini_batch_size]
60                 for k in xrange(0, n, mini_batch_size)]
61             for mini_batch in mini_batches:
62                 self.update_mini_batch(mini_batch, eta)
63             if test_data:
64                 print "Epoch {0}: {1} / {2}".format(
65                     j, self.evaluate(test_data), n_test)
66             else:
67                 print "Epoch {0} complete".format(j)
68
69     def update_mini_batch(self, mini_batch, eta):
70         """Update the network's weights and biases by applying
71         gradient descent using backpropagation to a single mini batch
```



# A look in the code

Add errors from all the training data from the mini-batch

Update the weights

```
*C:\Users\vacaca27\IMEC\OneDrive - imec\depplearning\crash course\neural-networks-and-deep-learning-master\neural-net...
File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?
new 1 x .gitignore x PopulateDefaults.m x framepulses.py x README.md x network.py x mnist_loader.py x
107
108     def update_mini_batch(self, mini_batch, eta):
109         """Update the network's weights and biases by applying
110         gradient descent using backpropagation to a single mini batch.
111         The ``mini_batch`` is a list of tuples ``(x, y)``, and ``eta``
112         is the learning rate."""
113         nabla_b = [np.zeros(b.shape) for b in self.biases]
114         nabla_w = [np.zeros(w.shape) for w in self.weights]
115         for x, y in mini_batch:
116             delta_nabla_b, delta_nabla_w = self.backprop(x, y)
117             nabla_b = [nb+dnb for nb, dnb in zip(nabla_b, delta_nabla_b)]
118             nabla_w = [nw+dnw for nw, dnw in zip(nabla_w, delta_nabla_w)]
119         self.weights = [w-(eta/len(mini_batch))*nw
120                        for w, nw in zip(self.weights, nabla_w)]
121         self.biases = [b-(eta/len(mini_batch))*nb
122                       for b, nb in zip(self.biases, nabla_b)]
123
124     def evaluate(self, test_data):
125         """Return the number of test inputs for which the neural
126         network outputs the correct result. Note that the neural
127         network's output is assumed to be the index of whichever
128         neuron in the final layer has the highest activation."""
129         test_results = [(np.argmax(self.feedforward(x)), y)
130                        for (x, y) in test_data]
131         return sum(int(x == y) for (x, y) in test_results)
132
```



# references

- <http://neuralnetworksanddeeplearning.com/chap1.html>
- <https://www.cs.cmu.edu/~bhiksha/courses/deeplearning/Fall.2015/>



# Questions?