

## Neural Network Training

Given training data  $\mathbf{D}=\{\mathbf{x}[m], \mathbf{y}[m]\}$ ,  $m=1,2,\dots, M$ , NN training is to learn the NN parameters  $\Theta$  by minimizing a loss function  $L(\mathbf{D}; \Theta)$ , i.e.,

$$L(\mathbf{D}; \Theta) = \frac{1}{M} \sum_{m=1}^N l(x[m], y[m], \Theta) + \lambda R(\Theta)$$

$$\Theta^* = \arg \min_{\Theta} L(\mathbf{D}; \Theta)$$

where

$$\Theta = [\mathbf{W}^1 \ \mathbf{W}_0^1 \ \mathbf{W}^2 \ \mathbf{W}_0^2 \ \dots \ \mathbf{W}^{L+1} \ \mathbf{W}_0^{L+1}]$$

## Loss Function for Binary Logistic Classification

For each training sample  $(\mathbf{x}[m], y[m])$ , where  $y[m]$  is a scalar, representing the probability of  $y[m]=1$ . Let  $\hat{y}[m]$  be the estimated output probability via sigmoid function through the forward propagation based on current parameters.

Commonly used loss functions:

- Squared loss

$$l(y[m], \hat{y}[m]) = \frac{1}{2} (y[m] - \hat{y}[m])^t (y[m] - \hat{y}[m])$$

## Loss function for multi-class logistic classification

For each training sample  $(\mathbf{x}[m], \mathbf{y}[m])$ , where  $\mathbf{y}[m]$  follows 1-of-K encoding and let  $\hat{\mathbf{y}}[m]$  be the estimated output vector via softmax through the forward propagation based on current parameters.

Commonly used loss functions:

- Squared loss

$$l(\mathbf{y}[m], \hat{\mathbf{y}}[m]) = \frac{1}{2} (\mathbf{y}[m] - \hat{\mathbf{y}}[m])^t (\mathbf{y}[m] - \hat{\mathbf{y}}[m])$$

- Cross-entropy loss:  $l(\mathbf{y}[m], \hat{\mathbf{y}}[m]) = -\sum_{k=1}^K y[m][k] * \log(\hat{\mathbf{y}}[m][k])$

## Loss Function for Regression

For each training sample  $(\mathbf{x}[m], y[m])$ , where  $y[m]$  is a real number and let  $\hat{y}[m]$  be the estimated output value through the forward propagation based on current parameters.

Squared loss :

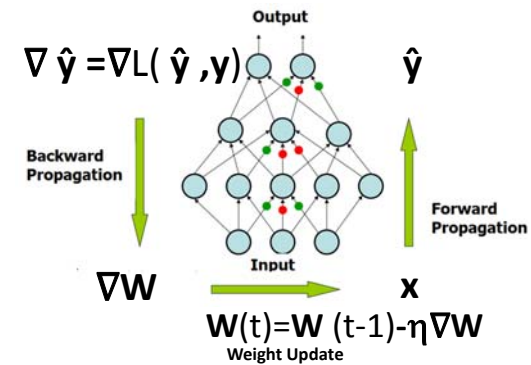
$$l(y[m], \hat{y}[m]) = (y[m] - \hat{y}[m])^2$$

## Backpropagation algorithm

- Three Steps of computation:
  - Forward pass:** run the NN and compute  $\hat{y}$
  - Backward pass:** start at the output layer, compute the output gradient  $\nabla \hat{y}$ , pass the output gradient backwards through the network, layer by layer, by recursively computing the weight gradients of each layer.
  - Weight updates:** update the weights for each layer based on the estimated gradients

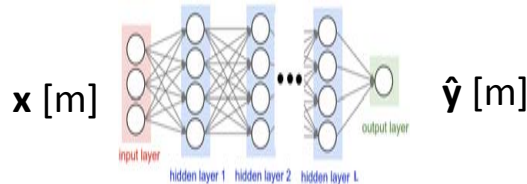
33

## Neural Network Learning



## Forward Propagation

- Given current  $W$ s and an input  $x[m]$ , compute  $\hat{y}[m]$



$$\begin{aligned}
 H^1[m] &= \phi((W^1)^T x[m] + W_0^1) \\
 H^2[m] &= \phi((W^2)^T H^1[m] + W_0^2) \\
 &\dots \\
 H^L[m] &= \phi((W^L)^T H^{L-1}[m] + W_0^L) \\
 \hat{y}[m] &= g((W^{L+1})^T H^L[m] + W_0^{L+1})
 \end{aligned}$$

- The output is computed through a series of recursive composition from input layer through the hidden layer until the output layer.
- $\phi()$  is the activation function and  $g()$  is the output function. Note they apply individually to each element of the vector.

## Derivatives with Matrices

Let  $A^{M \times N}$  be a matrix and  $z$  be a scalar function of  $A$

$$\frac{\partial z}{\partial A}^{M \times N} = \begin{bmatrix} \frac{\partial z}{\partial A[1][1]} & \frac{\partial z}{\partial A[1][2]} & \dots & \frac{\partial z}{\partial A[1][N]} \\ \frac{\partial z}{\partial A[2][1]} & \frac{\partial z}{\partial A[2][2]} & \dots & \frac{\partial z}{\partial A[2][N]} \\ \dots & \dots & \dots & \dots \\ \frac{\partial z}{\partial A[M][1]} & \frac{\partial z}{\partial A[M][2]} & \dots & \frac{\partial z}{\partial A[M][N]} \end{bmatrix} = \begin{pmatrix} \frac{\partial z}{\partial a_{11}} & \frac{\partial z}{\partial a_{12}} & \dots & \frac{\partial z}{\partial a_{1N}} \\ \frac{\partial z}{\partial a_{21}} & \frac{\partial z}{\partial a_{22}} & \dots & \frac{\partial z}{\partial a_{2N}} \\ \dots & \dots & \dots & \dots \\ \frac{\partial z}{\partial a_{M1}} & \frac{\partial z}{\partial a_{M2}} & \dots & \frac{\partial z}{\partial a_{MN}} \end{pmatrix}$$

## Derivatives with Matrices

- Vector by matrix . Let  $\mathbf{X}^{k \times 1}$  be a vector

$$\frac{\partial \mathbf{X}}{\partial \mathbf{A}} = \left[ \frac{\partial \mathbf{X}_1}{\partial \mathbf{A}} \quad \frac{\partial \mathbf{X}_2}{\partial \mathbf{A}} \quad \dots \quad \frac{\partial \mathbf{X}_K}{\partial \mathbf{A}} \right]^{M \times N \times K} \text{--Tensor}$$

- Tensor vector multiplication

Let  $\mathbf{Y}$  be a  $K \times 1$  vector

$$(\frac{\partial \mathbf{X}}{\partial \mathbf{A}})_{M \times N \times K} \mathbf{Y}^{K \times 1} = (\sum_{k=1}^K \frac{\partial \mathbf{X}_k}{\partial \mathbf{A}} \mathbf{Y}_k)_{M \times N}$$

- Tensor matrix multiplication

Let  $\mathbf{B}$  be a  $K \times D$  matrix

$$\frac{\partial \mathbf{X}}{\partial \mathbf{A}} \mathbf{B} = \left[ \frac{\partial \mathbf{X}}{\partial \mathbf{A}} \mathbf{B}_1, \frac{\partial \mathbf{X}}{\partial \mathbf{A}} \mathbf{B}_2, \dots, \frac{\partial \mathbf{X}}{\partial \mathbf{A}} \mathbf{B}_D \right]^{M \times N \times D}$$

where  $\mathbf{B}_i$  is the  $i$ th column of  $\mathbf{B}$

## Back Propagation

- The basic idea of backward propagation is to employ the gradient descent method to update the weight parameters for each layer.
- Let  $l(y, \hat{y})$  be the loss function for the predicted  $\hat{y}$  and given  $y$ , and  $W^l$  be the weight for the  $l$ th layer.
- Using the gradient descent method,  $W^l$  at  $t$ th iteration can be updated by

$$\mathbf{W}^l[t] = \mathbf{W}^l[t-1] - \eta \nabla_{\mathbf{W}^l} l(\mathbf{y}, \hat{\mathbf{y}})$$

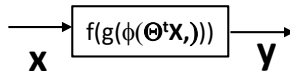
where

$\nabla_{\mathbf{w}} l(\mathbf{y}, \hat{\mathbf{y}})$  is the gradient of  $l(\mathbf{y}, \hat{\mathbf{y}})$  w.r.t.  $\mathbf{W}^1$  and by chain rule it can be computed as follows

$$\begin{aligned}\nabla_{\mathbf{w}^i} l(\mathbf{y}, \hat{\mathbf{y}}) &= \frac{\partial l(\mathbf{y}, \hat{\mathbf{y}})}{\mathbf{w}^i} \\ &= \frac{\partial \hat{\mathbf{y}}}{\mathbf{w}^i} \frac{\partial l(\mathbf{y}, \hat{\mathbf{y}})}{\hat{\mathbf{y}}} = \frac{\partial \hat{\mathbf{y}}}{\mathbf{w}^i} \nabla \hat{\mathbf{y}}\end{aligned}$$

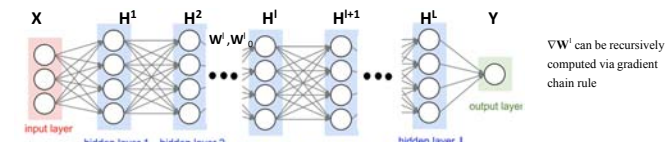
## Backward Propagation

## Gradient Chain rule



$$\nabla_{\Theta} = \frac{\partial \mathbf{y}}{\partial \Theta} = \frac{\partial f(g(\phi(\Theta' \mathbf{x})))}{\partial \Theta} = \frac{\partial f(g(\phi(\Theta' \mathbf{x})))}{\partial g(\phi(\mathbf{x}, \Theta))} \frac{\partial g(\phi(\Theta' \mathbf{x}))}{\partial \phi(\mathbf{x}, \Theta)} \frac{\partial \phi(\Theta' \mathbf{x})}{\partial (\Theta' \mathbf{x})} \frac{\partial \Theta' \mathbf{x}}{\partial \Theta}$$

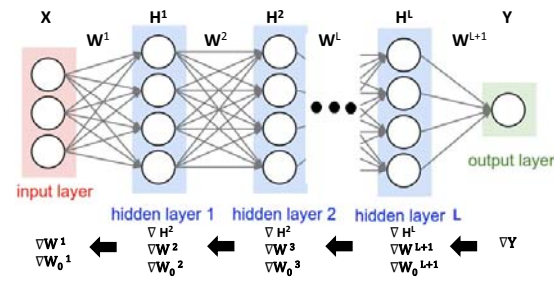
## Back Propagation-gradient chain rule



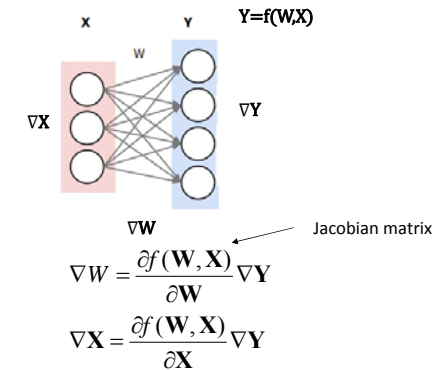
$$\begin{aligned} \nabla \mathbf{W}' &= \frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \mathbf{W}'} = \frac{N_{1-1} \times N_{\mathbf{N}_1} \times \mathbf{K}}{\partial \mathbf{W}'} \frac{\partial \hat{\mathbf{y}}[m]}{\partial \hat{\mathbf{y}}} \frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}} \\ &= \frac{\partial \mathbf{H}^{L-1}}{\partial \mathbf{W}'} \frac{\partial \mathbf{H}^{L-1}}{\partial \mathbf{H}^{L-1}} \frac{\partial \hat{\mathbf{y}}[m]}{\partial \hat{\mathbf{y}}} \frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}} = \frac{\partial \mathbf{H}'}{\partial \mathbf{W}'} \frac{\partial \mathbf{H}^{L-1}}{\partial \mathbf{H}'} \dots \frac{\partial \mathbf{H}^{L-1}}{\partial \mathbf{H}^{L-1}} \underbrace{\frac{\partial \hat{\mathbf{y}}[m]}{\partial \hat{\mathbf{y}}} \frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}}}_{\nabla \hat{\mathbf{y}}} \end{aligned}$$

Pay attention to matrix dimensions consistency!

## Back Propagation



## Back Propagation for Each Layer



## Compute output gradient

For the  $m$ th training sample, given groundtruth  $\mathbf{y}[m]$  and the computed  $\hat{\mathbf{y}}[m]$  by forward propagation, and let  $l(\mathbf{y}[m], \hat{\mathbf{y}}[m])$  be the loss function

$$\nabla \hat{\mathbf{y}}[m] = \nabla_{\hat{\mathbf{y}}} l(\mathbf{y}[m], \hat{\mathbf{y}}[m]) = \frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}}$$

## Compute output gradient (cont'd)

- 1) For squared loss function

$$l(\mathbf{y}[m], \hat{\mathbf{y}}[m]) = \frac{1}{2} (\mathbf{y}[m] - \hat{\mathbf{y}}[m])^t (\mathbf{y}[m] - \hat{\mathbf{y}}[m])$$

$$\begin{aligned} \nabla \hat{\mathbf{y}}[m] &= \frac{1}{2} \frac{\partial (\mathbf{y}[m] - \hat{\mathbf{y}}[m])^t (\mathbf{y}[m] - \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}} \\ &= -(\mathbf{y}[m] - \hat{\mathbf{y}}[m]) \end{aligned}$$

## Compute output gradient (cont'd)

2) For cross-entropy loss function

$$l(\mathbf{y}[m], \hat{\mathbf{y}}[m]) = -\sum_{k=1}^K \mathbf{y}[m][k] * \log(\hat{\mathbf{y}}[m][k])$$

$$\nabla \hat{\mathbf{y}}[m] = -\frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}}$$

$$= -\begin{bmatrix} \frac{\partial \sum_{k=1}^K \mathbf{y}[m][k] \log \hat{\mathbf{y}}[m][k]}{\partial \hat{\mathbf{y}}[m][1]} \\ \frac{\partial \sum_{k=1}^K \mathbf{y}[m][k] \log \hat{\mathbf{y}}[m][k]}{\partial \hat{\mathbf{y}}[m][2]} \\ \dots \\ \frac{\partial \sum_{k=1}^K \mathbf{y}[m][k] \log \hat{\mathbf{y}}[m][K]}{\partial \hat{\mathbf{y}}[m][K]} \end{bmatrix} = \begin{bmatrix} \frac{\mathbf{y}[m][1]}{\hat{\mathbf{y}}[m][1]} \\ \frac{\mathbf{y}[m][2]}{\hat{\mathbf{y}}[m][2]} \\ \dots \\ \frac{\mathbf{y}[m][K]}{\hat{\mathbf{y}}[m][K]} \end{bmatrix}$$

## Backward Propagation-Output Layer

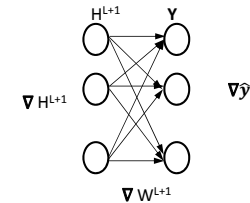
Given  $\nabla \hat{\mathbf{y}}[m]$ , compute  $\nabla \mathbf{H}^L[m]$  and  $\nabla \mathbf{W}^{L+1}[m]$

$$\hat{\mathbf{y}}[m] = g((\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1})$$

$$\nabla \mathbf{W}^{L+1}[m] = \frac{\partial \hat{\mathbf{y}}[m]}{\partial \mathbf{W}^{L+1}} \nabla \hat{\mathbf{y}}[m]$$

$$\nabla \mathbf{W}_0^{L+1}[m] = \frac{\partial \hat{\mathbf{y}}[m]}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{\mathbf{y}}[m]$$

$$\nabla \mathbf{H}^L[m] = \frac{\partial \hat{\mathbf{y}}[m]}{\partial \mathbf{H}^L[m]} \nabla \hat{\mathbf{y}}[m]$$



where  $g()$  is the output function

53

## Backward Propagation-Output Layer

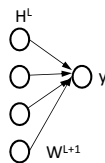
For regression, the output function is a scalar node and  $g()$  is

$$\hat{\mathbf{y}}[m] = g((\mathbf{W}^{L+1})^T \mathbf{H}^L[m]) = (\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}$$

$$\nabla \mathbf{W}^{L+1}[m] = \frac{\partial \hat{\mathbf{y}}[m]}{\partial \mathbf{W}^{L+1}} \nabla \hat{\mathbf{y}}[m] = \mathbf{H}^L[m] \nabla \hat{\mathbf{y}}[m]$$

$$\nabla \mathbf{W}_0^{L+1}[m] = \frac{\partial \hat{\mathbf{y}}[m]}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{\mathbf{y}}[m] = \nabla \hat{\mathbf{y}}[m]$$

$$\nabla \mathbf{H}^L[m] = \frac{\partial \hat{\mathbf{y}}[m]}{\partial \mathbf{H}^L} \nabla \hat{\mathbf{y}}[m] = \mathbf{W}^{L+1} \nabla \hat{\mathbf{y}}[m]$$



Note  $\mathbf{W}^{L+1}$  is a vector and  $\nabla \hat{\mathbf{y}}$  is a scalar.

54

## Backward Propagation-Output Layer

For binary classification  $y \in \{+1, -1\}$ , the output is a probability scalar node and output function is  $\sigma()$

$$\hat{\mathbf{y}}[m] = \sigma(\mathbf{z}[m]), \text{ where } \mathbf{z}[m] = (\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}$$

$$\nabla \mathbf{W}^{L+1}[m] = \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{W}^{L+1}} \nabla \hat{\mathbf{y}}[m]$$

$$= \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{z}[m]} \frac{\partial \mathbf{z}[m]}{\partial \mathbf{W}^{L+1}} \nabla \hat{\mathbf{y}}[m]$$

$$= \sigma(\mathbf{z}[m])(1 - \sigma(\mathbf{z}[m])) \mathbf{H}^L[m] \nabla \hat{\mathbf{y}}[m]$$

$$\nabla \mathbf{W}_0^{L+1}[m] = \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{\mathbf{y}}[m]$$

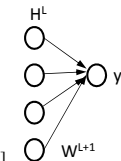
$$= \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{z}[m]} \frac{\partial \mathbf{z}[m]}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{\mathbf{y}}[m]$$

$$= \sigma(\mathbf{z}[m])(1 - \sigma(\mathbf{z}[m])) \nabla \hat{\mathbf{y}}[m]$$

$$\nabla \mathbf{H}^L[m] = \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{H}^L} \nabla \hat{\mathbf{y}}[m]$$

$$= \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{z}[m]} \frac{\partial \mathbf{z}[m]}{\partial \mathbf{H}^L} \mathbf{W}^{L+1} \nabla \hat{\mathbf{y}}[m]$$

$$= \sigma(\mathbf{z}[m])(1 - \sigma(\mathbf{z}[m])) \mathbf{W}^{L+1} \nabla \hat{\mathbf{y}}[m]$$



Note  $\mathbf{W}^{L+1}$  is a vector and  $\nabla \hat{\mathbf{y}}$  is a scalar.

55

## Backward Propagation-Output Layer

For binary classification  $y \in \{1,0\}$ , the output is a scalar node and output function is  $\sigma()$

$$\hat{y}[m] = \frac{e^{z[m]}}{1 + e^{z[m]}}, \text{ where } z[m] = (\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}$$

$$\nabla \mathbf{W}^{L+1}[m] = \frac{\partial \hat{y}[m]}{\partial \mathbf{W}^{L+1}} \nabla \hat{y}[m]$$

$$= \frac{\partial \hat{y}[m]}{\partial z[m]} \frac{\partial z[m]}{\partial \mathbf{W}^{L+1}} \nabla \hat{y}[m]$$

$$= (\hat{y}[m])^2 \mathbf{H}^L[m] \nabla \hat{y}[m]$$

$$\nabla \mathbf{W}_0^{L+1}[m] = \frac{\partial \hat{y}[m]}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{y}[m]$$

$$= \frac{\partial \hat{y}[m]}{\partial z[m]} \frac{\partial (z[m])}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{y}[m]$$

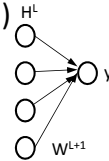
$$= (\hat{y}[m])^2 \nabla \hat{y}[m]$$

$$\text{Note } \frac{\partial \hat{y}[m]}{\partial z[m]} = (\hat{y}[m])^2$$

$$\nabla \mathbf{H}^L[m] = \frac{\partial \hat{y}[m]}{\partial \mathbf{H}^L} \nabla \hat{y}[m]$$

$$= \frac{\partial \hat{y}[m]}{\partial z[m]} \frac{\partial z[m]}{\partial \mathbf{H}^L} \mathbf{W}^{L+1} \nabla \hat{y}[m]$$

$$= (\hat{y}[m])^2 \mathbf{W}^{L+1} \nabla \hat{y}[m]$$

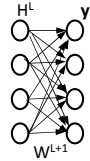


Note  $\mathbf{W}^{L+1}$  is a vector and  $\nabla \hat{y}$  is a scalar.

57

## Backward Propagation-Output Layer

For multi-class classification, the output has multiple nodes and  $g()$  is multi-class Sigmoid function  $\sigma_M()$



$$\text{Let } z[m] = (\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}$$

$$\hat{y}[m] = \sigma_M(z[m]) = \begin{bmatrix} \sigma_M(z[m][1]) \\ \sigma_M(z[m][2]) \\ \vdots \\ \sigma_M(z[m][K]) \end{bmatrix}$$

$$\text{where } z[m][k] = (\mathbf{W}_k^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[k]$$

$$\sigma_M(z[m][k]) = \frac{\exp((\mathbf{W}_k^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[k])}{\sum_{j=1}^K \exp((\mathbf{W}_j^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[j])}$$

$$\nabla \mathbf{W}^{L+1}[m] = \frac{\partial \sigma_M(z[m])}{\partial \mathbf{W}^{L+1}} \nabla \hat{y}[m]$$

$$= \sum_{i=1}^K \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}^{L+1}} \nabla \hat{y}[m][i]$$

$$= \sum_{i=1}^K \left[ \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_1^{L+1}} \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_2^{L+1}} \dots \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_K^{L+1}} \right] \nabla \hat{y}[m][i]$$

where

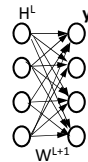
$$\frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_k^{L+1}} = \begin{cases} \sigma_M(z[m][k])(1 - \sigma_M(z[m][k])) \mathbf{H}^L[m] & \text{if } i = k \\ -\sigma_M(z[m][k]) \sigma_M(z[m][i]) \mathbf{H}^L[m] & \text{else} \end{cases}$$

Note  $\mathbf{W}^{L+1}$  is a matrix and  $\nabla \hat{y}$  is a vector

59

## Backward Propagation-Output Layer

For multi-class classification, the output has multiple nodes and  $g()$  is multi-class Sigmoid function  $\sigma_M()$



$$\text{Let } z[m] = (\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}$$

$$\hat{y}[m] = \sigma_M(z[m]) = \begin{bmatrix} \sigma_M(z[m][1]) \\ \sigma_M(z[m][2]) \\ \vdots \\ \sigma_M(z[m][K]) \end{bmatrix}$$

$$\text{where } z[m][k] = (\mathbf{W}_k^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[k]$$

$$\nabla \mathbf{W}_0^{L+1}[m] = \frac{\partial \sigma_M(z[m])}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{y}[m]$$

$$= \sum_{i=1}^K \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_0^{L+1}} \nabla \hat{y}[m][i] = \sum_{i=1}^K \left[ \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_0^{L+1}[1]} \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_0^{L+1}[2]} \dots \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{W}_0^{L+1}[K]} \right] \nabla \hat{y}[m][i]$$

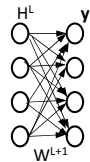
$$\sigma_M(z[m][k]) = \frac{\exp((\mathbf{W}_k^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[k])}{\sum_{j=1}^K \exp((\mathbf{W}_j^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[j])}$$

Note  $\mathbf{W}^{L+1}$  is a matrix and  $\nabla \hat{y}$  is a vector

61

## Backward Propagation-Output Layer

For multi-class classification, the output has multiple nodes and  $g()$  is multi-class Sigmoid function  $\sigma_M()$



$$\text{Let } z[m] = (\mathbf{W}^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}$$

$$\hat{y}[m] = \sigma_M(z[m]) = \begin{bmatrix} \sigma_M(z[m][1]) \\ \sigma_M(z[m][2]) \\ \vdots \\ \sigma_M(z[m][K]) \end{bmatrix}$$

$$\text{where } z[m][k] = (\mathbf{W}_k^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[k]$$

$$\sigma_M(z[m][k]) = \frac{\exp((\mathbf{W}_k^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[k])}{\sum_{j=1}^K \exp((\mathbf{W}_j^{L+1})^T \mathbf{H}^L[m] + \mathbf{W}_0^{L+1}[j])}$$

$$\nabla \mathbf{H}^L[m] = \frac{\partial \sigma_M(z[m])}{\partial \mathbf{H}^L} \nabla \hat{y}[m]$$

$$= \sum_{i=1}^K \frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{H}^L} \nabla \hat{y}[m][i]$$

where

$$\frac{\partial \sigma_M(z[m][i])}{\partial \mathbf{H}^L} = \sigma_M(z[m][i]) (\mathbf{W}_i^{L+1} - \sum_{j=1}^K \sigma_M(z[m][j]) \mathbf{W}_j^{L+1})$$

Note  $\mathbf{W}^{L+1}$  is a matrix and  $\nabla \hat{y}$  is a vector

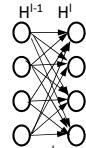
63

## Backward Propagation-Hidden Layers

For each hidden layer  $l$ , let  $\mathbf{H}^l$  be its output,  $\mathbf{H}^{l-1}$  be its inputs, and  $\mathbf{W}^l$  be its current parameters. Given  $\nabla \mathbf{H}^l[m]$ , compute  $\nabla \mathbf{H}^{l-1}[m]$ ,  $\nabla \mathbf{W}^l[m]$ , and  $\nabla \mathbf{W}_0^l[m]$

$$\mathbf{H}^l = \phi((\mathbf{W}^l)^t \mathbf{H}^{l-1} + \mathbf{W}_0^l)$$

$$\nabla \mathbf{H}^{l-1}[m] = \frac{\partial \mathbf{H}^l[m]}{\partial \mathbf{H}^{l-1}[m]} \nabla \mathbf{H}^l[m]$$

$$\nabla \mathbf{W}^l[m] = \frac{\partial \mathbf{H}^l[m]}{\partial \mathbf{W}^l} \nabla \mathbf{H}^l[m] \quad \nabla \mathbf{W}_0^l[m] = \frac{\partial \mathbf{H}^l[m]}{\partial \mathbf{W}_0^l} \nabla \mathbf{H}^l[m]$$


66

## Backward Propagation-Hidden Layers

Let  $\mathbf{z}[m] = (\mathbf{W}^l)^t \mathbf{H}^{l-1}[m] + \mathbf{W}_0^l$

$$\mathbf{H}^l[m] = \phi(\mathbf{z}[m]) = \begin{bmatrix} \phi(\mathbf{z}[m][1]) \\ \phi(\mathbf{z}[m][2]) \\ \vdots \\ \phi(\mathbf{z}[m][N_l]) \end{bmatrix}, \text{ where } \mathbf{z}[m][i] = (\mathbf{W}_i^l)^t \mathbf{H}^{l-1}[m] + \mathbf{W}_0^l[i]$$

$$\nabla \mathbf{W}^l[m] = \frac{\partial \phi(\mathbf{z}[m])}{\partial \mathbf{W}^l} \nabla \mathbf{H}^l[m]$$

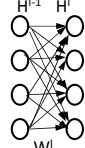
$$= \sum_{i=1}^{N_l} \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}^l} \nabla \mathbf{H}^l[m][i]$$

$$= \sum_{i=1}^{N_l} \left[ \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_1^l} \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_2^l} \dots \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_{N_l}^l} \right] \nabla \mathbf{H}_i^l[m][i], \text{ where } \mathbf{W}_j \text{ is the } j\text{th column of } \mathbf{W}$$

For sigmoid transfer function

$$\frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_k^l} = \begin{cases} \sigma(\mathbf{z}[m][k])(1 - \sigma(\mathbf{z}[m][k])) \mathbf{H}^{l-1}[m] & \text{if } i = k \\ 0 & \text{else} \end{cases}$$

For ReLU transfer function

$$\frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_k^l} = \frac{\partial \text{ReLU}(\mathbf{z}[m][i])}{\partial \mathbf{W}_k^l} = \begin{cases} \mathbf{H}^{l-1}[m] & \text{if } i = k \text{ and } \mathbf{z}[m][i] > 0 \\ 0 & \text{else} \end{cases}$$


67

## Backward Propagation-Hidden Layers

Let  $\mathbf{z}[m] = (\mathbf{W}^l)^t \mathbf{H}^{l-1}[m] + \mathbf{W}_0^l$

$$\mathbf{H}^l[m] = \phi(\mathbf{z}[m]) = \begin{bmatrix} \phi(\mathbf{z}[m][1]) \\ \phi(\mathbf{z}[m][2]) \\ \vdots \\ \phi(\mathbf{z}[m][N_l]) \end{bmatrix}, \text{ where } \mathbf{z}[m][i] = (\mathbf{W}_i^l)^t \mathbf{H}^{l-1}[m] + \mathbf{W}_0^l[i]$$

$$\nabla \mathbf{W}_0^l[m] = \frac{\partial \phi(\mathbf{z}[m])}{\partial \mathbf{W}_0^l} \nabla \mathbf{H}^l[m]$$

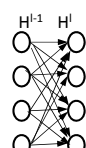
$$= \sum_{i=1}^{N_l} \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_0^l} \nabla \mathbf{H}^l[m][i]$$

$$= \sum_{i=1}^{N_l} \left[ \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_0^l[1]} \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_0^l[2]} \dots \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_0^l[N_l]} \right] \nabla \mathbf{H}^l[m][i]$$

For sigmoid transfer function

$$\frac{\partial \phi(\mathbf{z}[m])}{\partial \mathbf{W}_0^l[k]} = \begin{cases} \sigma(\mathbf{z}[m])(1 - \sigma(\mathbf{z}[m])) & \text{if } i = k \\ 0 & \text{else} \end{cases}$$

For ReLU transfer function

$$\frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{W}_0^l[k]} = \frac{\partial \text{ReLU}(\mathbf{z}[m][i])}{\partial \mathbf{W}_0^l[k]} = \begin{cases} 1 & \text{if } i = k \text{ and } \mathbf{z}[m][i] > 0 \\ 0 & \text{else} \end{cases}$$


68

## Backward Propagation-Hidden Layers

Let  $\mathbf{z}[m] = (\mathbf{W}^l)^t \mathbf{H}^{l-1}[m] + \mathbf{W}_0^l$

$$\mathbf{H}^l[m] = \phi(\mathbf{z}[m]) = \begin{bmatrix} \phi(\mathbf{z}[m][1]) \\ \phi(\mathbf{z}[m][2]) \\ \vdots \\ \phi(\mathbf{z}[m][N_l]) \end{bmatrix}, \text{ where } \mathbf{z}[m][i] = (\mathbf{W}_i^l)^t \mathbf{H}^{l-1}[m] + \mathbf{W}_0^l[i]$$

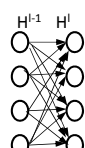
$$\nabla \mathbf{H}^{l-1}[m] = \frac{\partial \phi(\mathbf{z}[m])}{\partial \mathbf{H}^{l-1}} \nabla \mathbf{H}^l[m]$$

$$= \sum_{i=1}^{N_l} \frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{H}^{l-1}} \nabla \mathbf{H}^l[m][i]$$

For sigmoid transfer function

$$\frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{H}^{l-1}} = \sigma(\mathbf{z}[m][i])(1 - \sigma(\mathbf{z}[m][i])) \mathbf{W}_i^l$$

For ReLU transfer function

$$\frac{\partial \phi(\mathbf{z}[m][i])}{\partial \mathbf{H}^{l-1}} = \frac{\partial \text{ReLU}(\mathbf{z}[m][i])}{\partial \mathbf{H}^{l-1}} = \begin{cases} \mathbf{W}_i^l & \text{if } \mathbf{z}[m][i] > 0 \\ 0 & \text{else} \end{cases}$$


69

## Backward Propagation-First Hidden Layer

For  $l=1$ , given  $\nabla \mathbf{H}^1[m]$ , compute  $\nabla \mathbf{W}^1[m]$

$$\mathbf{z}[m] = (\mathbf{W}^1)^T \mathbf{X}[m] + \mathbf{W}_0^1, \quad \mathbf{H}^1[m] = \phi(\mathbf{z}[m])$$

$$\nabla \mathbf{W}^1[m] = \frac{\partial \mathbf{H}^1[m]}{\partial \mathbf{W}^1} \nabla \mathbf{H}^1[m] = \frac{\partial \sigma(\mathbf{z}[m])}{\partial \mathbf{W}^1} \nabla \mathbf{H}^1[m]$$

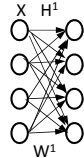
$$= \sum_{i=1}^{N_1} \frac{\partial \sigma(\mathbf{z}[m][i])}{\partial \mathbf{W}^1} \nabla \mathbf{H}^1[m][i] = \sum_{i=1}^{N_1} \left[ \frac{\partial \sigma(\mathbf{z}[m][i])}{\partial \mathbf{W}_1^1} \frac{\partial \sigma(\mathbf{z}[m][i])}{\partial \mathbf{W}_2^1} \dots \frac{\partial \sigma(\mathbf{z}[m][i])}{\partial \mathbf{W}_{N_1}^1} \right] \nabla \mathbf{H}^1[m][i]$$

For sigmoid activation function

$$\frac{\partial \sigma(\mathbf{z}[m][i])}{\partial \mathbf{W}_k^1} = \begin{cases} \sigma(\mathbf{z}[m][k])(1 - \sigma(\mathbf{z}[m][k])) \mathbf{X}[m] & \text{if } i = k \\ 0 & \text{else} \end{cases}$$

For ReLU activation function

$$\frac{\partial \text{ReLU}(\mathbf{z}[m][i])}{\partial \mathbf{W}_k^1} = \begin{cases} \mathbf{X}[m] & \text{if } i = k \text{ and } \mathbf{z}[m][i] > 0 \\ 0 & \text{else} \end{cases}$$



71

## Add Regularization

For squared of L2-norm, its gradient is

$$\|\mathbf{W}^l\|_2^2 = \left[ \sum_{i=1}^{N_{l-1}} \sum_{j=1}^{N_l} (\mathbf{W}^l[i][j])^2 \right]^{\frac{1}{2}}$$

$$\begin{aligned} \nabla R_{L2}(\mathbf{W}^l) &= \frac{\partial (\|\mathbf{W}^l\|_2^2)}{\partial \mathbf{W}^l} = \frac{\partial \left( \sum_{i=1}^{N_{l-1}} \sum_{j=1}^{N_l} (\mathbf{W}^l[i][j])^2 \right)}{\partial \mathbf{W}^l} \\ &= \sum_{i=1}^{N_{l-1}} \sum_{j=1}^{N_l} \frac{\partial ((\mathbf{W}^l[i][j])^2)}{\partial \mathbf{W}^l} = \sum_{i=1}^{N_{l-1}} \sum_{j=1}^{N_l} \frac{\partial ((\mathbf{W}^l[i][j])^2)}{\partial \mathbf{W}^l[i][j]} \frac{\partial \mathbf{W}^l[i][j]}{\partial \mathbf{W}^l} \\ &= \sum_{i=1}^{N_{l-1}} \sum_{j=1}^{N_l} 2\mathbf{W}^l[i][j] \frac{\partial \mathbf{W}^l[i][j]}{\partial \mathbf{W}^l} \end{aligned}$$

=  $\{2\mathbf{W}^l[i][j]\}$   
where  $\{\mathbf{W}^l[i][j]\}$  is a matrix, whose  $i$ th row and  $j$ th column entry is  $2\mathbf{W}^l[i][j]$

73

## Add Regularization

For L1 norm, its gradient is

$$\begin{aligned} \|\mathbf{W}^l\|_1 &= \max_{1 \leq j \leq N_l} \sum_{i=1}^{N_{l-1}} |\mathbf{W}^l[i][j]| = \max_{i=1}^{N_{l-1}} \left( \sum_{j=1}^{N_l} |\mathbf{W}^l[i][j]| \right) \\ \nabla R_{L1}(\mathbf{W}^l) &= \frac{\partial \|\mathbf{W}^l\|_1}{\partial \mathbf{W}^l} = \left[ \frac{\partial \max_{1 \leq j \leq N_l} \sum_{i=1}^{N_{l-1}} |\mathbf{W}^l[i][j]|}{\partial \mathbf{W}_1^l} \quad \frac{\partial \max_{1 \leq j \leq N_l} \sum_{i=1}^{N_{l-1}} |\mathbf{W}^l[i][j]|}{\partial \mathbf{W}_2^l} \quad \dots \quad \frac{\partial \max_{1 \leq j \leq N_l} \sum_{i=1}^{N_{l-1}} |\mathbf{W}^l[i][j]|}{\partial \mathbf{W}_{N_l}^l} \right] \\ &= \begin{bmatrix} 0 & 0 & \text{sign}(\mathbf{W}^l[1][j^*]) & \dots & 0 \\ 0 & 0 & \text{sign}(\mathbf{W}^l[2][j^*]) & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \text{sign}(\mathbf{W}^l[N_{l-1}][j^*]) & \dots & 0 \end{bmatrix} \quad j^* = \arg \max_j \left( \sum_{i=1}^{N_{l-1}} |\mathbf{W}^l[i][j]| \right) \end{aligned}$$

74

## Add Regularization

For regularization on  $\mathbf{W}^l_0$ , we can use the vector regularization equations to derive its L1 and L2-squared norms

75



## Parameter Updating

- Update the parameters  $\mathbf{W}^l$  for each layer  
For  $l=1$  to  $L+1$  do

$$\mathbf{W}^l[t] = \mathbf{W}^l[t-1] - \eta_l (\nabla \mathbf{W}^l[m] + \lambda \frac{\partial R(\mathbf{W}^l)}{\partial \mathbf{W}^l})$$

$$\mathbf{W}_0^l[t] = \mathbf{W}_0^l[t-1] - \eta_{l0} (\nabla \mathbf{W}_0^l[m] + \lambda_0 \frac{\partial R(\mathbf{W}_0^l)}{\partial \mathbf{W}_0^l})$$

76

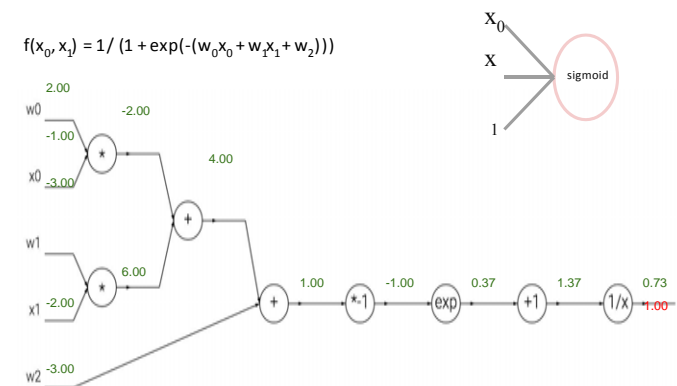
## Back Propagation Steps by sample

- Initialize the weights
- For each training sample  $\mathbf{x}[m]$ , perform forward propagation to compute  $\nabla \hat{\mathbf{y}}[m]$
- For hidden layers  $l=L+1$  to 1 do (note  $\mathbf{H}^{L+1}=\mathbf{y}$ ,  $\mathbf{H}^0=\mathbf{x}$ )  
Compute  $\nabla \mathbf{W}^l(m)$ ,  $\nabla \mathbf{W}_0^l(m)$ , and  $\nabla \mathbf{H}^{l-1}(m)$   
For  $l=1$  to  $N_{l+1}$   
Update weights for each layer  
 $\mathbf{W}^l[t] = \mathbf{W}^l[t-1] - \eta_l (\nabla \mathbf{W}^l(m) + \lambda \frac{\partial R(\mathbf{W}^l)}{\partial \mathbf{W}^l})$   
Update bias weights for each layer  
 $\mathbf{W}_0^l[t] = \mathbf{W}_0^l[t-1] - \eta_{l0} (\nabla \mathbf{W}_0^l(m) + \lambda_0 \frac{\partial R(\mathbf{W}_0^l)}{\partial \mathbf{W}_0^l})$
- Repeat for all training samples until convergence

## Back Propagation Steps by Batch

- Initialize the weights
- Randomly select from all training samples a mini-batch  $\mathbf{D}_i$  of size  $M_i$ .
- For each training sample  $\mathbf{x}[m]$  in  $\mathbf{D}_i$ , perform forward propagation to compute  $\nabla \hat{\mathbf{y}}[m]$
- For hidden layers  $l=L+1$  to 1 do (note  $\mathbf{H}^{L+1}=\mathbf{y}$ ,  $\mathbf{H}^0=\mathbf{x}$ )  
Compute  $\nabla \mathbf{W}^l(m)$ ,  $\nabla \mathbf{W}_0^l(m)$ , and  $\nabla \mathbf{H}^{l-1}(m)$   
Compute the average weight gradients for all samples in  $\mathbf{D}_i$   
 $\nabla \mathbf{W}^l = \frac{1}{M_i} \sum_{m=1}^{M_i} \nabla \mathbf{W}^l[m]$     $\nabla \mathbf{W}_0^l = \frac{1}{M_i} \sum_{m=1}^{M_i} \nabla \mathbf{W}_0^l[m]$   
For  $l=1$  to  $N_{l+1}$   
Update weights for each layer    $\mathbf{W}^l[t] = \mathbf{W}^l[t-1] - \eta_l (\nabla \mathbf{W}^l + \lambda \frac{\partial R(\mathbf{W}^l)}{\partial \mathbf{W}^l})$   
Update bias weights for each layer    $\mathbf{W}_0^l[t] = \mathbf{W}_0^l[t-1] - \eta_{l0} (\nabla \mathbf{W}_0^l + \lambda_0 \frac{\partial R(\mathbf{W}_0^l)}{\partial \mathbf{W}_0^l})$
- Repeat by selecting another mini-batch until convergence

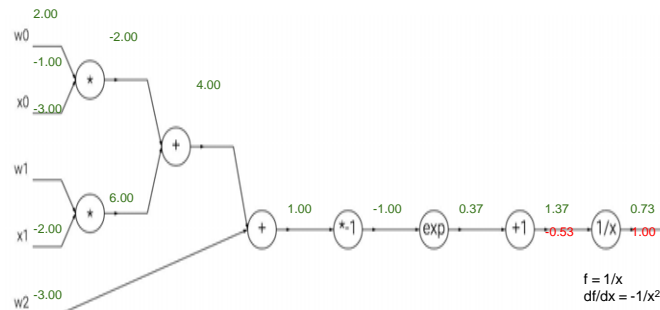
## Backward Computation



<http://cs231n.github.io/optimization-2/>

## Backward Computation

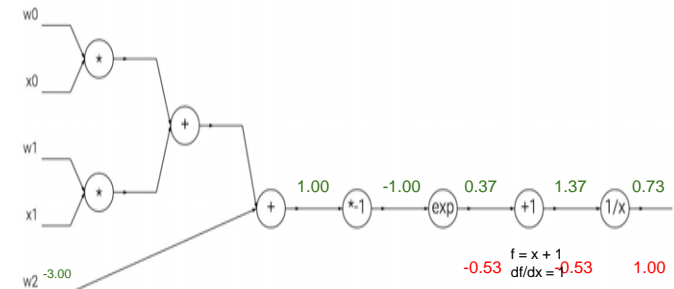
$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Backward Computation

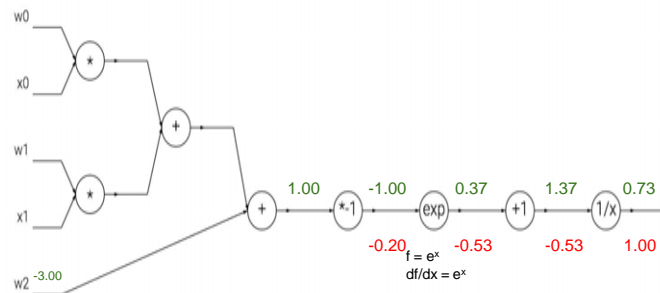
$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Backward Computation

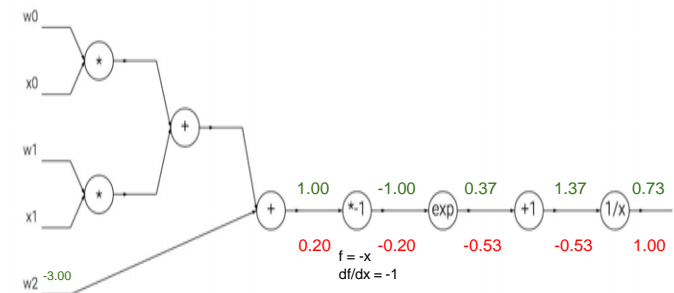
$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Backward Computation

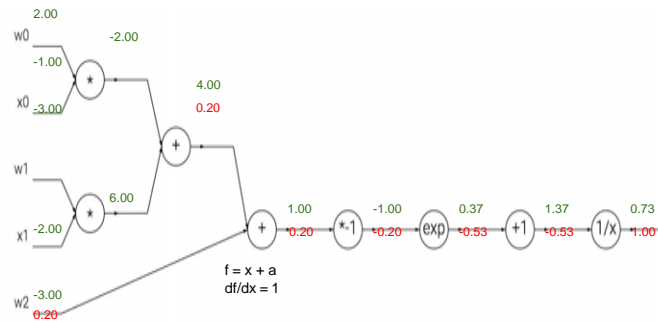
$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Backward Computation

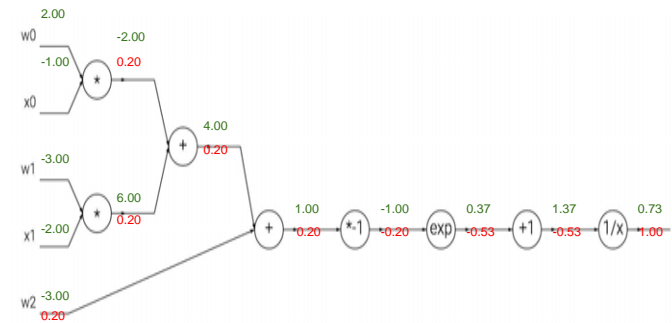
$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0x_0 + w_1x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Backward Computation

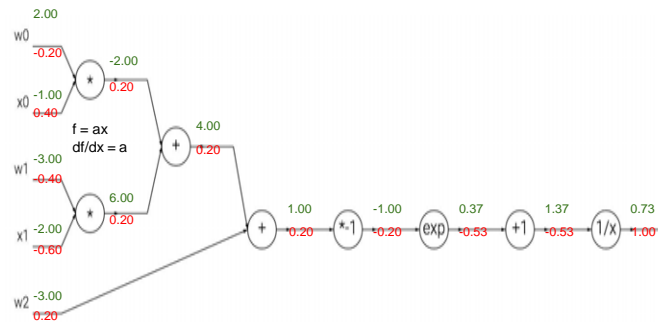
$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0x_0 + w_1x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Backward Computation

$$f(x_0, x_1) = 1 / (1 + \exp(-(w_0x_0 + w_1x_1 + w_2)))$$



<http://cs231n.github.io/optimization-2/>

## Training

- The training continues epoch by epoch. Each epoch goes through all training samples (all mini-batches)
- For each epoch, training samples order may be **randomized** to avoid any correlations between epochs
- The learning continues until the stopping criterion is satisfied.

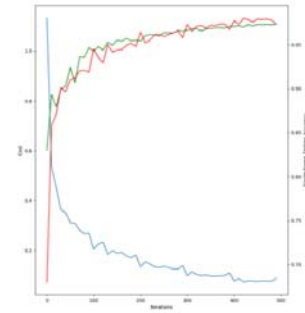
98

## Stopping criteria

- Average loss change:
  - Back-prop is considered to have converged when the absolute rate of change in the average loss function per epoch is sufficiently small
- Classification performance
  - After each epoch, the NN is tested on
    - the training data to evaluate classification performance
    - testing data to evaluate the generalization performance
    - If the performance does not improve further, stop.

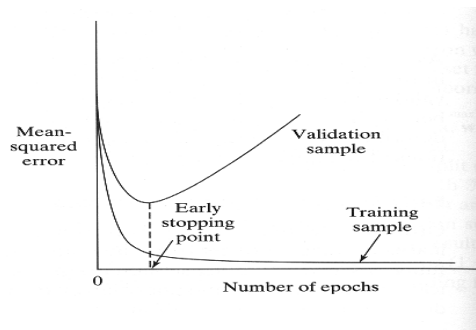
99

## NN Training



NOTE: Red is training accuracy, green is testing accuracy, and blue is training loss. (It's not labeled on the graph.)

## Early stopping



PMR5406 Redes Neurais e  
Lógica Fuzzy

Multilayer Perceptrons

101

## Initializations

- Initialize all biases to zero
- Initialize the weights to small non-zero random values. They should not be the same or all equal to zero. For example,  $w = N(0,1)$

102

## Hyper-parameter Tuning

---

On the validation dataset (different from testing data),

- Use the grid method to greedily tune the hyper-parameters
  - The grid search method first estimates the range for each parameter and then quantizes each parameter into a finite set of intervals.
  - Given initial values of hyper- parameters, it then tunes each hyper-parameter individually, while fixing other parameters.
  - The process repeats until convergence.
- 

## Hyper-parameter Tuning

---

Adaptive learning rate

- Learning rate should vary with the iteration. It can be set as a function of the gradient magnitude. Two common choice of choosing learning rates are  $O(1/t)$  and  $O(1/\sqrt{t})$  for strong convexity. Polynomial decay or AdaGad are also widely used.
  - Learning rates for different parameters should be different
-