

□ *Code*

Neural Networks

Joseph E. Gonzalez

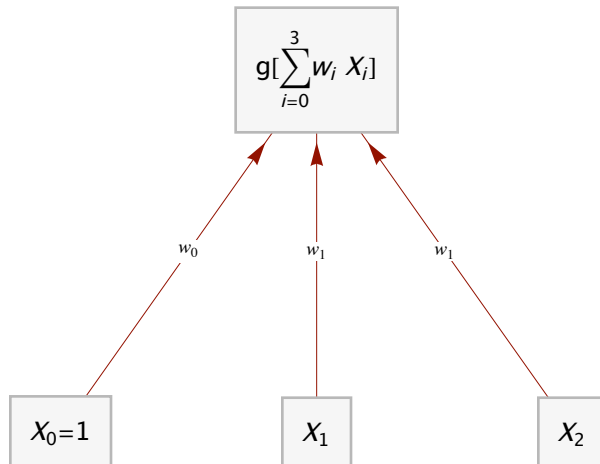
We are going to go through Neural Networks and review the process of back propagation.

Experimental *Mathematica* based presentation.

Single Perceptron

□ *The Perceptron*

perceptronPlot



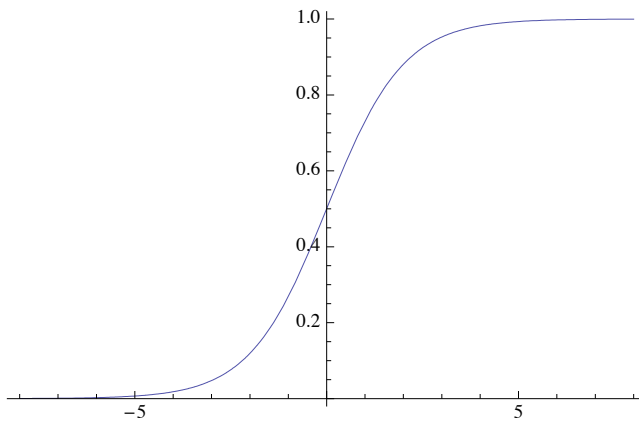
□ *There are several parts*

1. Link Function $g[u]$
2. Weights w_i
3. A bias term X_0

Link Function

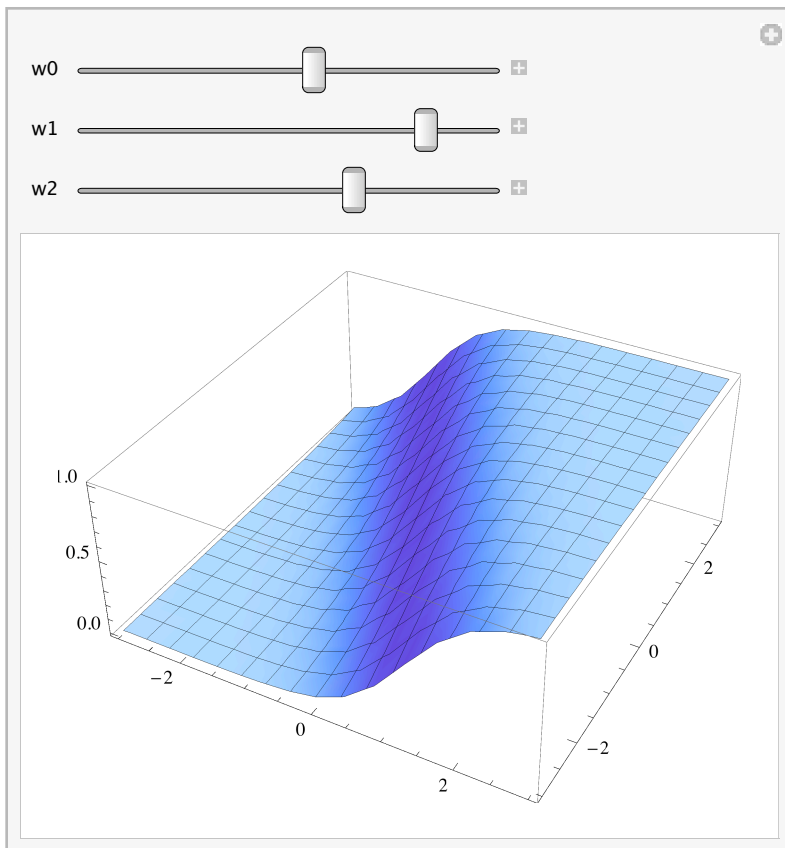
$$g[x] = \frac{1}{1 + e^{-x}} \quad (1)$$

```
g = Function[x,  $\frac{1}{1 + \text{Exp}[-\mathbf{x}]}$ ];  
Plot[g[x], {x, -8, 8}]
```



Demo

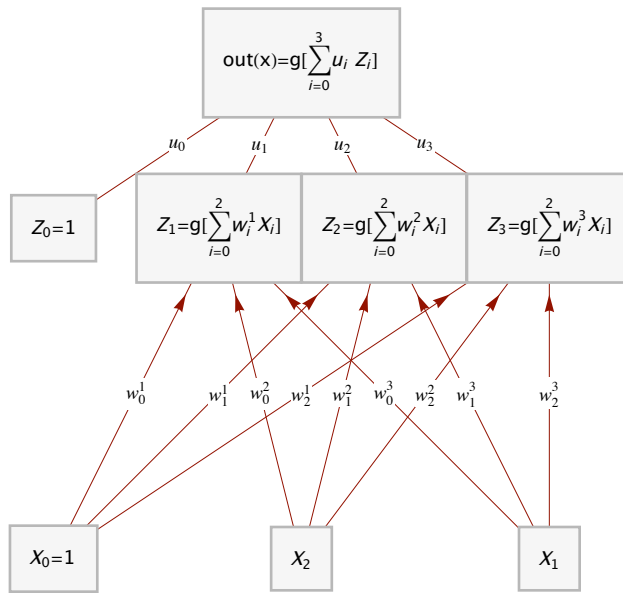
```
Manipulate[  
  g = Function[x,  $\frac{1}{1 + \text{Exp}[-x]}$ ];  
  Plot3D[g[w0 + w1 x1 + w2 x2], {x1, -3, 3}, {x2, -3, 3},  
    {{w0, 0}, -3, 3}, {{w1, 2}, -3, 3}, {{w2, -2}, -3, 3}]  
]
```



Neural Network with Multiple Hidden Layers

- Lets Consider what this network looks like

plt



- **Matlab Style Forward Propagation**

Lets define a matrix W as:

$$W = \begin{pmatrix} w_0^1 & w_1^1 & w_2^1 \\ w_0^2 & w_1^2 & w_2^2 \\ w_0^3 & w_1^3 & w_2^3 \end{pmatrix} \quad (2)$$

We can multiply this matrix by X where we have added a 1

$$W \cdot \begin{pmatrix} 1 \\ X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} w_0^1 & w_1^1 & w_2^1 \\ w_0^2 & w_1^2 & w_2^2 \\ w_0^3 & w_1^3 & w_2^3 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ X_2 \\ X_3 \end{pmatrix} = \begin{pmatrix} w_0^1 + w_1^1 X_1 + w_2^1 X_2 \\ w_0^2 + w_1^2 X_1 + w_2^2 X_2 \\ w_0^3 + w_1^3 X_1 + w_2^3 X_2 \end{pmatrix} \quad (3)$$

Lets define function application as element wise. Then we obtain:

$$g\left[W \cdot \begin{pmatrix} 1 \\ X_2 \\ X_3 \end{pmatrix}\right] = \begin{pmatrix} g[w_0^1 + w_1^1 X_1 + w_2^1 X_2] \\ g[w_0^2 + w_1^2 X_1 + w_2^2 X_2] \\ g[w_0^3 + w_1^3 X_1 + w_2^3 X_2] \end{pmatrix} = \begin{pmatrix} Z_1 \\ Z_2 \\ Z_3 \end{pmatrix} \quad (4)$$

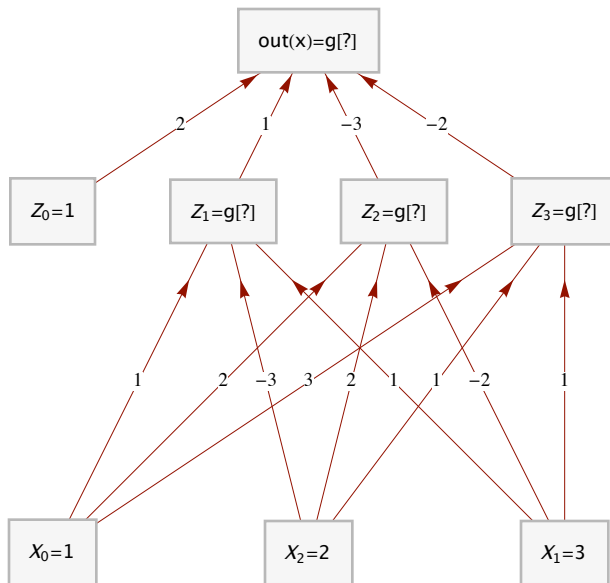
We can then prepend a 1 to the result to obtain:

$$\text{out}(X) = g\left[\begin{pmatrix} u_0 & u_1 & u_2 & u_3 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ Z_1 \\ Z_2 \\ Z_3 \end{pmatrix} \right] = g\left[u_0 + \sum_{i=1}^3 u_i Z_i \right] \quad (5)$$

Forward Propagation (Example) #1

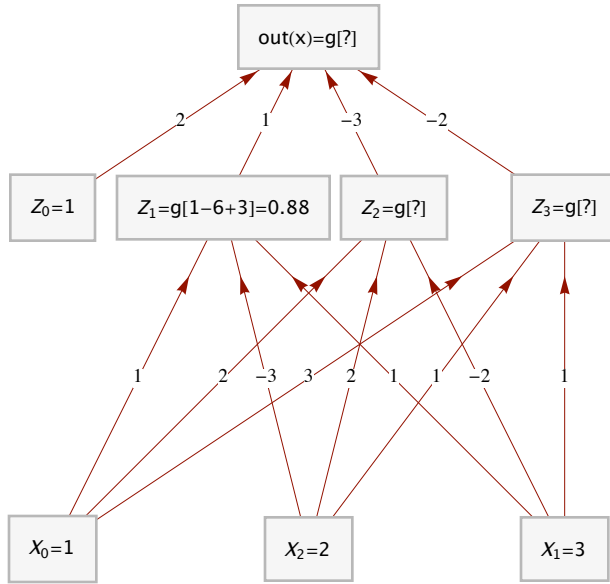
- What is the value of Z_1

```
plotTree[{"out(x)=g[?]", "z0=1", "z1=g[?]", "z2=g[?]", "z3=g[?]",
  "x0=1", "x2=2", "x1=3"}]
```



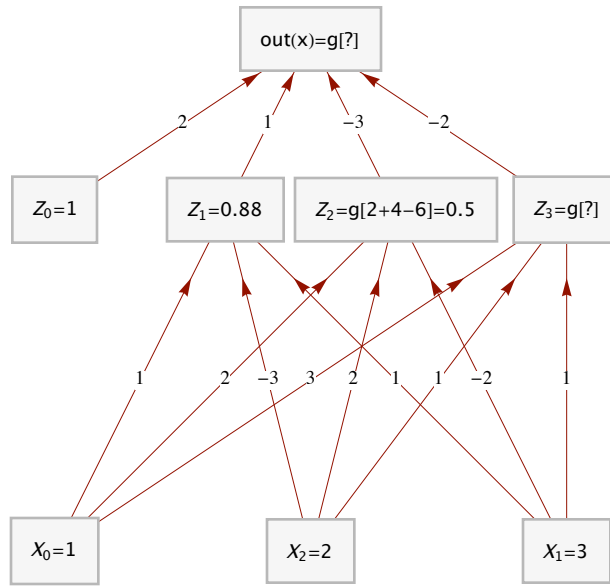
What is the value of Z_2 ?

```
plotTree[{"out(x)=g[?]", "z0=1", "z1=g[1-6+3]=0.88", "z2=g[?]",
"z3=g[?]", "x0=1", "x2=2", "x1=3"}]
```



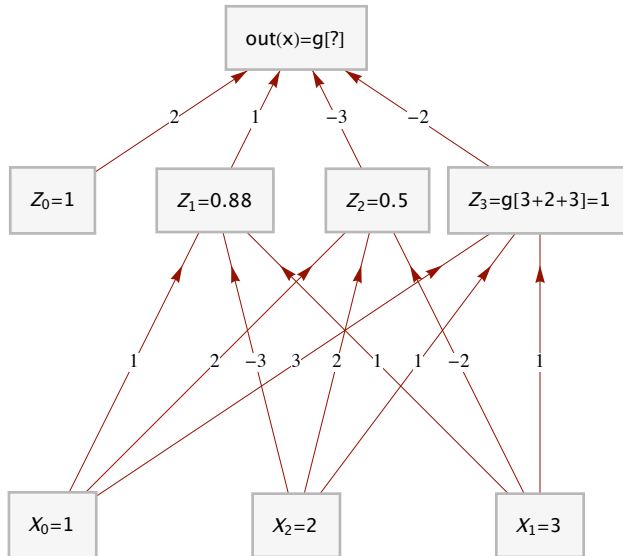
What is the value of Z_3 ?

```
plotTree[{"out(x)=g[?]", "z0=1", "z1=0.88", "z2=g[2+4-6]=0.5", "z3=g[?]",
"X0=1", "X2=2", "X1=3"}]
```



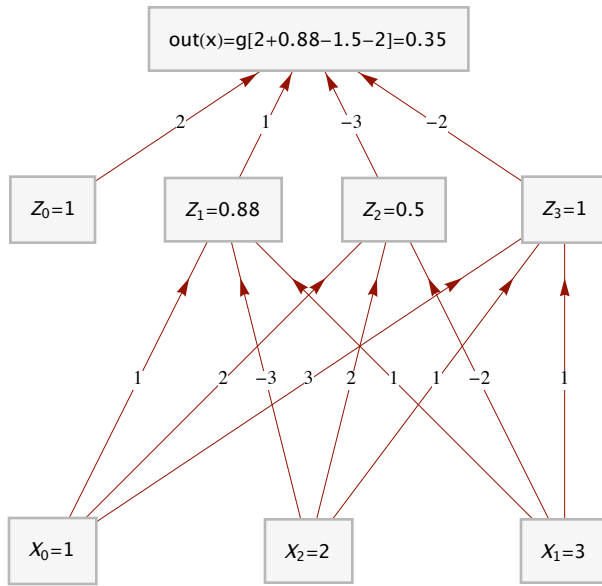
What is the value of $\text{out}(X)$?

```
plotTree[{"out(x)=g[?]", "z0=1", "z1=0.88", "z2=0.5", "z3=g[3+2+3]=1",
  "x0=1", "x2=2", "x1=3"}]
```



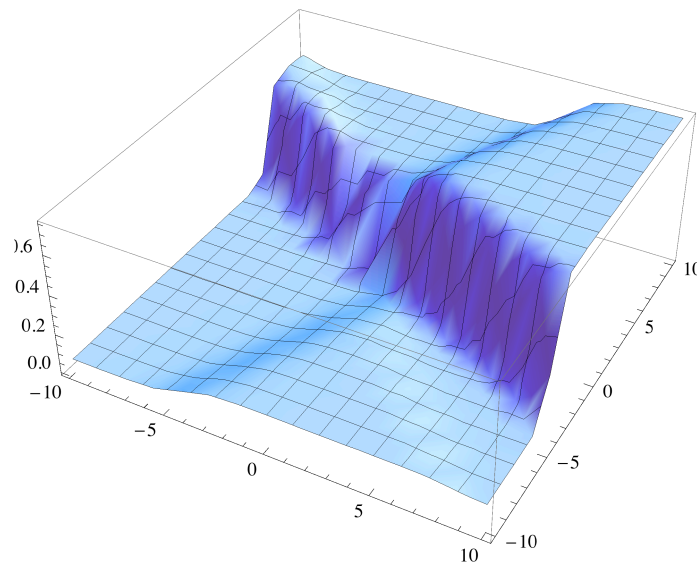
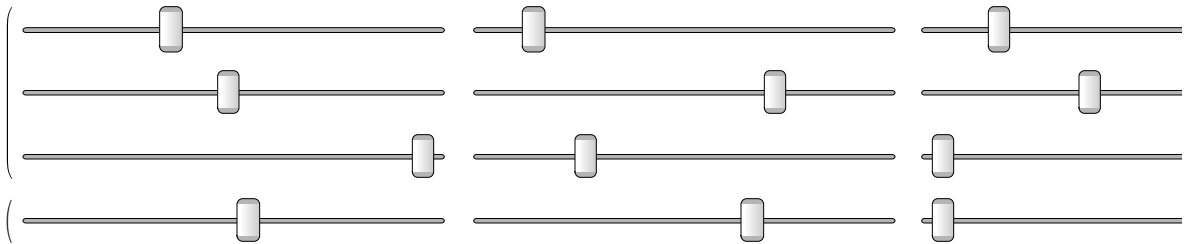
Done!

```
plotTree[{"out(x)=g[2+0.88-1.5-2]=0.35", "z_0=1", "z_1=0.88", "z_2=0.5",
"z_3=1", "x_0=1", "x_2=2", "x_1=3"}]
```



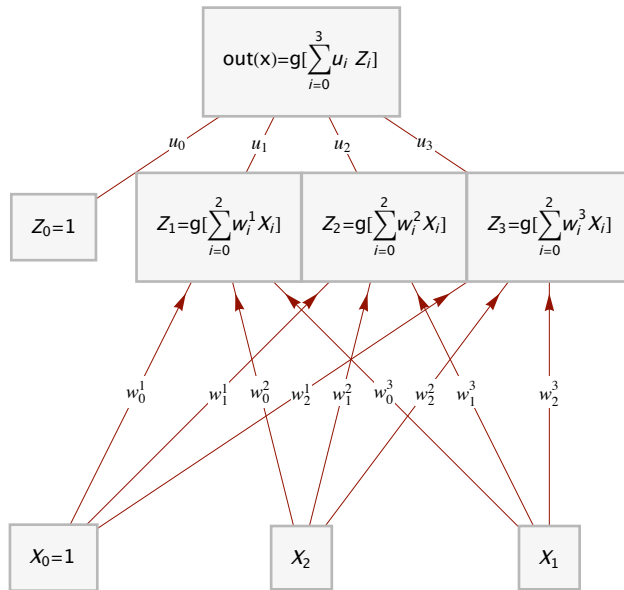
Demo

dynamicDemo



Generalized Back Propagation

plt



Suppose we want to find the best model $out(x; U, W)$ with respect to the parameters W and U . How can we quantify best? Lets considered mean squared error.

$$E = \sum_{i=1}^n (out(X_i) - Y_i)^2 \tag{6}$$

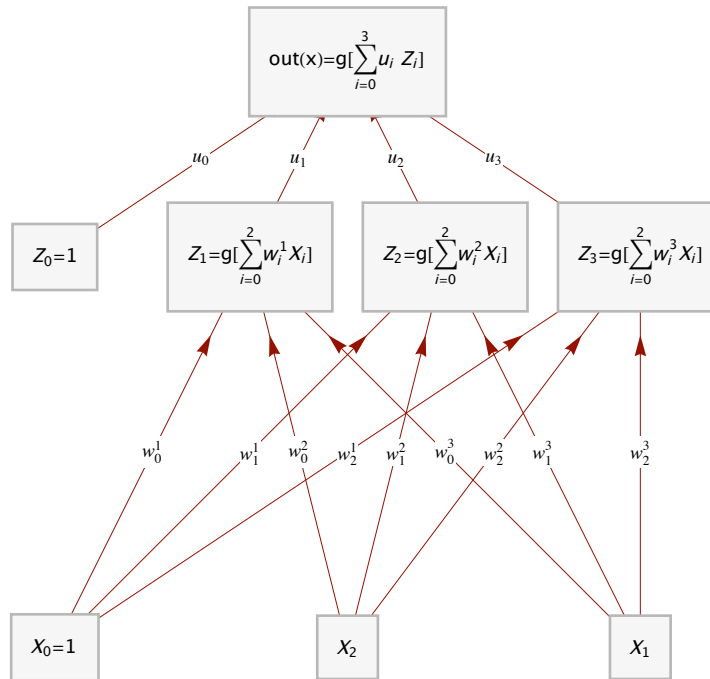
There are many ways to do this. One of the most common (and least effective) methods is to use gradient descent. This corresponds to the update rule:

$$u_i^{(t+1)} \leftarrow u_i^{(t)} - \eta \frac{\partial E}{\partial u_i} \Bigg|_{u_i^{(t)}} \tag{7}$$

$$w_{ij}^{(t+1)} \leftarrow w_{ij}^{(t)} - \eta \frac{\partial E}{\partial w_{ij}} \Bigg|_{w_{ij}^{(t)}} \tag{8}$$

Recall we have the following graph:

plt



Lets first derive the update rule for U :

$$E = (\text{out}(X) - Y)^2 \quad (9)$$

Taking the derivative we get (stuck?):

$$\frac{\partial E}{\partial u_k} = \frac{\partial}{\partial u_k} (\text{out}(X) - Y)^2 \quad (10)$$

Applying the infamous chain rule:

$$\frac{\partial}{\partial x} f(g(x)) = \left(\frac{\partial}{\partial u} f(u) \Big|_{u=g(x)} \right) \left(\frac{\partial}{\partial x} g(x) \right) \quad (11)$$

$$\frac{\partial E}{\partial u_k} = 2(\text{out}(X) - Y) \left(\frac{\partial}{\partial u_k} \text{out}(X) \right) \quad (12)$$

$\frac{\partial}{\partial x} x^2 = 2x$

$$\frac{\partial}{\partial x} f(g(x)) = f'(g(x)) g'(x)$$

Now we need to take the derivative of the neural network. Lets first replace out with the function from the top perceptron

$$\frac{\partial E}{\partial u_k} = 2(\text{out}(X) - Y) \left(\frac{\partial}{\partial u_k} g \left[\sum_{i=0}^3 u_i Z_i \right] \right) \quad (13)$$

Chain rule again

$$\frac{\partial E}{\partial u_k} = 2 (\text{out}(X) - Y) g' \left[\sum_{i=0}^3 u_i Z_i \right] \left(\sum_{i=0}^3 \frac{\partial}{\partial u_k} u_i Z_i \right) \quad (14)$$

We know that only one term in the Z_i sum will remain and that is $Z_{i=k}$

$$\frac{\partial E}{\partial u_k} = 2 (\text{out}(X) - Y) g' \left[\sum_{i=0}^3 u_i Z_i \right] Z_k \quad (15)$$

Done that's it!!! Sort of. Lets look at the derivative of $g[x] = \frac{1}{1 + \text{Exp}[-x]}$

$$g'[x] = \frac{\partial}{\partial x} (1 + \text{Exp}[-x])^{-1} \quad (16)$$

$$g'[x] = -(1 + \text{Exp}[-x])^{-2} \frac{\partial}{\partial x} (1 + \text{Exp}[-x]) \quad (17)$$

$$g'[x] = -(1 + \text{Exp}[-x])^{-2} \left(\frac{\partial}{\partial x} 1 + \frac{\partial}{\partial x} \text{Exp}[-x] \right) \quad (18)$$

$$g'[x] = -(1 + \text{Exp}[-x])^{-2} \left(0 + \text{Exp}[-x] \frac{\partial}{\partial x} (-x) \right) \quad (19)$$

$$g'[x] = -(1 + \text{Exp}[-x])^{-2} (0 - \text{Exp}[-x]) \quad (20)$$

$$g'[x] = (1 + \text{Exp}[-x])^{-2} \text{Exp}[-x] \quad (21)$$

With some manipulation we get:

$$g'[x] = \frac{\text{Exp}[-x]}{(1 + \text{Exp}[-x])} \frac{1}{(1 + \text{Exp}[-x])} \quad (22)$$

$$g'[x] = \frac{\text{Exp}[-x]}{1 + \text{Exp}[-x]} g[x] \quad (23)$$

$$g'[x] = \frac{1 + \text{Exp}[-x] - 1}{1 + \text{Exp}[-x]} g[x] \quad (24)$$

$$g'[x] = \left(\frac{1 + \text{Exp}[-x]}{1 + \text{Exp}[-x]} + \frac{-1}{1 + \text{Exp}[-x]} \right) g[x] \quad (25)$$

$$g'[x] = \left(1 - \frac{1}{1 + \text{Exp}[-x]} \right) g[x] \quad (26)$$

$$g'[x] = (1 - g[x]) g[x] \quad (27)$$

Recall that we earlier had:

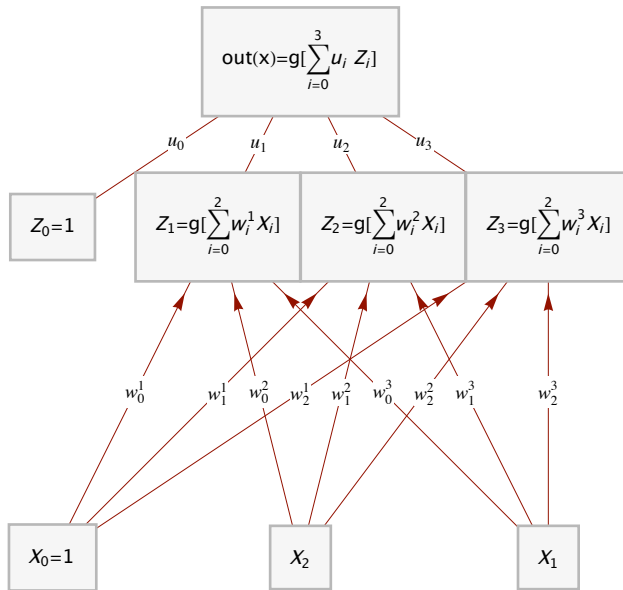
$$\frac{\partial E}{\partial u_k} = 2 (\text{out}(X) - Y) g' \left[\sum_{i=0}^3 u_i Z_i \right] Z_k \quad (28)$$

we can make a simple substitution to get:

$$\frac{\partial E}{\partial u_k} = 2(\text{out}(X) - Y) \left(1 - g \left[\sum_{i=0}^3 u_i Z_i \right] \right) g \left[\sum_{i=0}^3 u_i Z_i \right] Z_k \quad (29)$$

$$\frac{\partial E}{\partial u_k} = 2(\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) Z_k \quad (30)$$

plt



Gradient of W

That wasn't too bad. How about the next layer. We again start with:

$$E = (\text{out}(X) - Y)^2 \quad (31)$$

Taking the derivative with respect to w_k^r (and applying the chain rule)

$$\frac{\partial E}{\partial w_k^r} = \frac{\partial E}{\partial \text{out}(X)} \frac{\partial \text{out}(X)}{\partial w_k^r} = 2(\text{out}(X) - Y) \frac{\partial}{\partial w_k^r} \text{out}(X) \quad (32)$$

Expanding out we get:

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) \frac{\partial}{\partial w_k^r} g\left[\sum_{i=0}^3 u_i Z_i\right] \quad (33)$$

Chain rule:

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) \left(g'\left[\sum_{i=0}^3 u_i Z_i\right] \right) \left(\sum_{i=0}^3 \frac{\partial}{\partial w_k^r} u_i Z_i \right) \quad (34)$$

Recall that $g'[x] = (1 - g[x])g[x]$

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) \left(\sum_{i=0}^3 \frac{\partial}{\partial w_k^r} u_i Z_i \right) \quad (35)$$

Remember that each of the Z_i is connected to all the perceptrons from the lower level so we must take the derivative of each.

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) \left(\sum_{i=0}^3 u_i \frac{\partial}{\partial w_k^r} Z_i \right) \quad (36)$$

Which becomes:

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) \left(\sum_{i=0}^3 u_i \frac{\partial}{\partial w_k^r} g\left[\sum_{s=0}^3 w_s^i X_s\right] \right) \quad (37)$$

Another application of the chain rule:

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) \sum_{i=0}^3 u_i g'\left[\sum_{s=0}^3 w_s^i X_s\right] \frac{\partial}{\partial w_k^r} \sum_{s=0}^3 w_s^i X_s \quad (38)$$

Taking the final derivative we have:

$$\frac{\partial E}{\partial w_k^r} = 2(\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) \left(\sum_{i=0}^3 u_i (1 - Z_i) (Z_i) \right) X_k \quad (39)$$

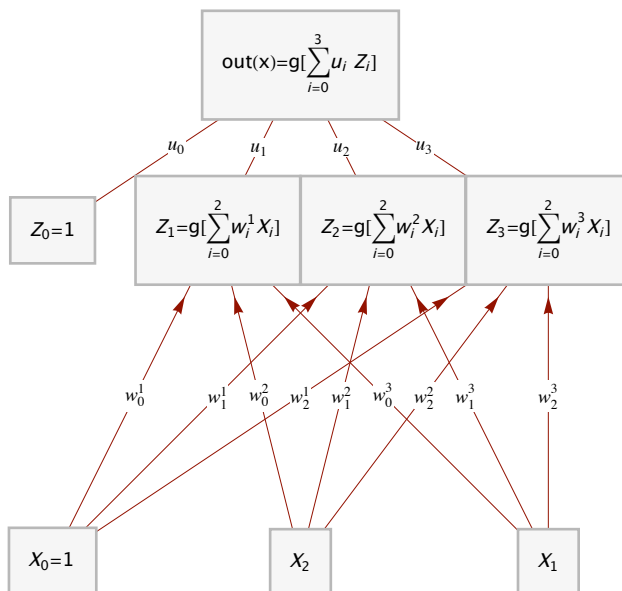
Why is it called back propagation?

Lets look at the following two equations:

$$\frac{\partial E}{\partial u_k} = 2 (\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) Z_k \quad (40)$$

$$\frac{\partial E}{\partial w_k^r} = 2 (\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X) \sum_{i=0}^3 u_i (1 - Z_i) (Z_i) X_k \quad (41)$$

plt



We propagate the derivative information backwards to the inputs:

$$\frac{\partial E}{\partial u_k} = (2 (\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X)) Z_k \quad (42)$$

$$\frac{\partial E}{\partial w_k^r} = (2 (\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X)) \sum_{i=0}^3 u_i (1 - Z_i) (Z_i) X_k \quad (43)$$

$$\frac{\partial E}{\partial w_k^r} = (2 (\text{out}(X) - Y) (1 - \text{out}(X)) \text{out}(X)) \sum_{i=0}^3 u_i (1 - Z_i) (Z_i) X_k \quad (44)$$

Questions?



- *Any feedback about Mathematica style presentations is welcome.;*