

Artificial intelligence in data science

Backpropagation

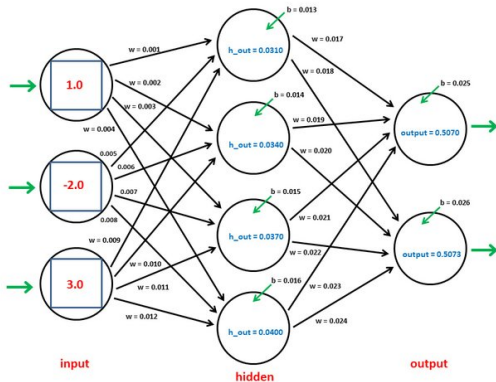
Janos Török

Department of Theoretical Physics

September 28, 2022

Fully connected neural networks

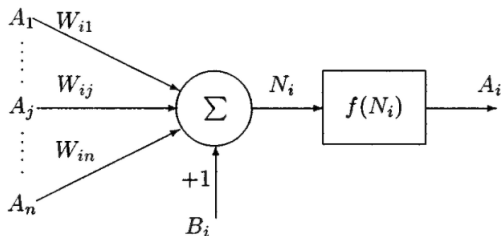
- ▶ Ideas from Piotr Skalski (practice), Pataki Bálint Ármin (lecture) and HMKCode (lecture)



Fully connected neural networks

► Model:

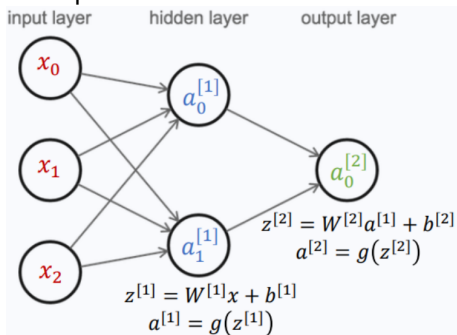
- Inputs (x_j) or for hidden layer l : A_j^{l-1}
- Weight w_{ij}^l
- Bias b_i^l
- Weighted sum of input and bias: $z_i^l = \sum_j A_j^{l-1} w_{ij}^l + b_i^l$
- Activation function (nonlinear) g : $A_i^l = g(z_i^l)$



Yang et al, 2000.

Feed forward

▶ Example



$$z_0^{[1]} = w_{0,0}^{[1]}x_0 + w_{0,1}^{[1]}x_1 + w_{0,2}^{[1]}x_2 + b_0^{[1]}$$
$$a_0^{[1]} = g(z_0^{[1]})$$

$$z_1^{[1]} = w_{1,0}^{[1]}x_0 + w_{1,1}^{[1]}x_1 + w_{1,2}^{[1]}x_2 + b_1^{[1]}$$
$$a_1^{[1]} = g(z_1^{[1]})$$

$$z_0^{[2]} = w_{0,0}^{[2]}a_0^{[1]} + w_{0,1}^{[2]}a_1^{[1]} + b_0^{[2]}$$
$$a_0^{[2]} = g(z_0^{[2]})$$

- ▶ We have an output, how to change weights and biases to achieve the desired output?
- ▶ Error L

Backpropagation



$$\Delta W = -\alpha \frac{\partial L}{\partial W}$$

- ▶ W is a large three dimensional matrix
- ▶ Chain rule!



Geoffrey Hinton

FOLLOW

Emeritus Prof. Comp Sci, U.Toronto & Engineering Fellow, Google
Verified email at cs.toronto.edu - [Homepage](#)

[machine learning](#) [neural networks](#) [artificial intelligence](#) [cognitive science](#)
[computer science](#)

TITLE	CITED BY	YEAR
Learning internal representations by error-propagation DE Rumelhart, GE Hinton, RJ Williams Parallel Distributed Processing: Explorations in the Microstructure of ...	63004 [*]	1986

Backpropagation

► Chain rule

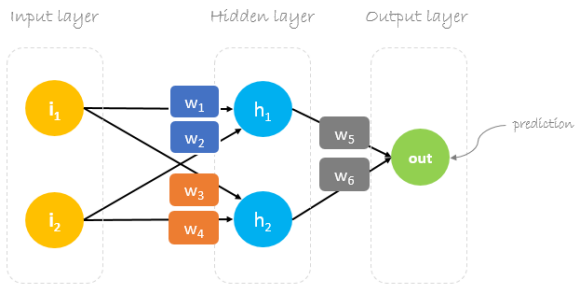
$$a_i = g(z_i) = g(w_{ij}a_j + b_i)$$

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial a_i} \frac{\partial a_i}{\partial z_i} \frac{\partial z_i}{\partial w_{ij}} = \frac{\partial L}{\partial a_i} g'(z_i) a_j$$

$$\frac{\partial L}{\partial a_i} = \sum_{l \in L} \frac{\partial L}{\partial a_l} \frac{\partial a_l}{\partial z_l} \frac{\partial z_l}{\partial a_i} = \sum_{l \in L} \frac{\partial L}{\partial a_l} g'(z_l) w_{li}$$

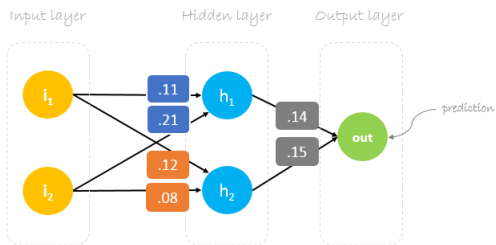
Backpropagation: Example

- ▶ From HMKCode
- ▶ Note that there is no activation function (it would just add one more step in the chain rule)



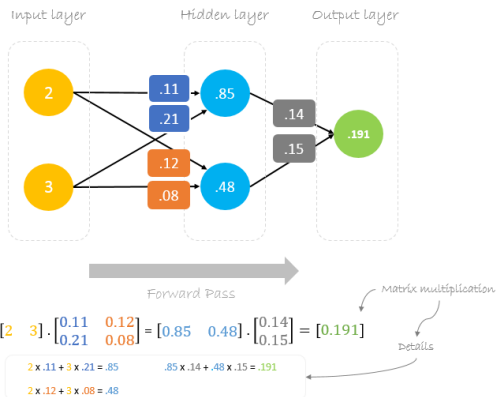
Backpropagation: Example

► Weights



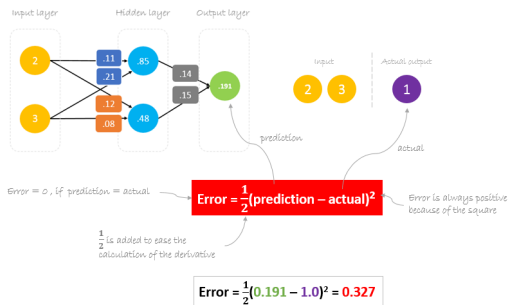
Backpropagation: Example

► Feedforward



Backpropagation: Example

► Error from the desired target



Backpropagation: Example

► Prediction function

prediction = out

$$\text{prediction} = (h_1) w_5 + (h_2) w_6$$

$$\begin{aligned} h_1 &= i_1 w_1 + i_2 w_2 \\ h_2 &= i_1 w_3 + i_2 w_4 \end{aligned}$$

$$\text{prediction} = (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$$

to change **prediction** value,
we need to change **weights**

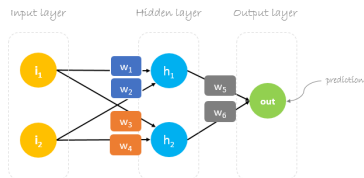
Backpropagation: Example

► Gradient descent

$$*W_x = W_x - a \left(\frac{\partial \text{Error}}{\partial W_x} \right)$$

Annotations for the equation above:
- W_x (old): Old weight
- a : Learning rate
- $\left(\frac{\partial \text{Error}}{\partial W_x} \right)$: Derivative of Error with respect to weight
- $*W_x$ (new): New weight

$$*W_6 = W_6 - a \left(\frac{\partial \text{Error}}{\partial W_6} \right)$$



Backpropagation: Example

► Chain rule

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial W_6} \quad \leftarrow \text{chain rule}$$

Error = $\frac{1}{2}(\text{prediction} - \text{actual})^2$

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{1}{2}(\text{prediction} - \text{actual})^2 * \frac{\partial (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6}{\partial W_6}$$

prediction = $(i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$

$$\frac{\partial \text{Error}}{\partial W_6} = 2 * \frac{1}{2}(\text{prediction} - \text{actual}) * \frac{\partial (\text{prediction} - \text{actual})}{\partial \text{prediction}} * (i_1 w_3 + i_2 w_4)$$

$h_2 = i_1 w_3 + i_2 w_4$

$$\frac{\partial \text{Error}}{\partial W_6} = (\text{prediction} - \text{actual}) * (h_2)$$

$\Delta = \text{prediction} - \text{actual}$ ← delta

$$\frac{\partial \text{Error}}{\partial W_6} = \Delta h_2$$

Backpropagation: Example

► Chain rule

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial W_6} \quad \leftarrow \text{chain rule}$$

Error = $\frac{1}{2}(\text{prediction} - \text{actual})^2$

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{1}{2}(\text{prediction} - \text{actual})^2 * \frac{\partial ((i_1 W_1 + i_2 W_2) W_5 + (i_1 W_3 + i_2 W_4) W_6)}{\partial W_6}$$

prediction = $(i_1 W_1 + i_2 W_2) W_5 + (i_1 W_3 + i_2 W_4) W_6$

$$\frac{\partial \text{Error}}{\partial W_6} = 2 * \frac{1}{2}(\text{prediction} - \text{actual}) \frac{\partial ((i_1 W_3 + i_2 W_4) W_6)}{\partial W_6}$$

$h_2 = i_1 W_3 + i_2 W_4$

$$\frac{\partial \text{Error}}{\partial W_6} = (\text{prediction} - \text{actual}) * (h_2)$$

$\Delta = \text{prediction} - \text{actual}$ ← delta

$$\frac{\partial \text{Error}}{\partial W_6} = \Delta h_2$$

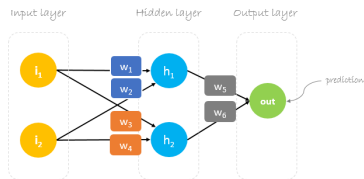
$$*W_6 = W_6 - a \Delta h_2$$

Backpropagation: Example

► Chain rule

$$*W_6 = W_6 - a \Delta h_2$$

$$*W_5 = W_5 - a \Delta h_1$$



Backpropagation: Example

► Chain rule

$$\frac{\partial \text{Error}}{\partial w_1} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial h_1} * \frac{\partial h_1}{\partial w_1} \leftarrow \text{chain rule}$$

$$\text{Error} = \frac{1}{2}(\text{prediction} - \text{actual})^2$$

$$\text{prediction} = (h_1) w_5 + (h_2) w_6$$

$$h_1 = i_1 w_1 + i_2 w_2$$

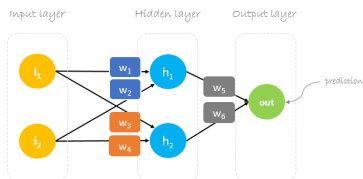
$$\frac{\partial \text{Error}}{\partial w_1} = \frac{\partial \frac{1}{2}(\text{prediction} - \text{actual})^2}{\partial \text{prediction}} * \frac{\partial (h_1) w_5 + (h_2) w_6}{\partial h_1} * \frac{\partial i_1 w_1 + i_2 w_2}{\partial w_1}$$

$$\frac{\partial \text{Error}}{\partial w_1} = 2 * \frac{1}{2}(\text{prediction} - \text{actual}) \frac{\partial (\text{prediction} - \text{actual})}{\partial \text{prediction}} * (w_5) * (i_1)$$

$$\frac{\partial \text{Error}}{\partial w_1} = (\text{prediction} - \text{actual}) * (w_5 i_1)$$

$$\Delta = \text{prediction} - \text{actual} \leftarrow \text{delta}$$

$$\frac{\partial \text{Error}}{\partial w_1} = \Delta w_5 i_1$$



Backpropagation: Example

► Summarized

updated weights

$$\begin{aligned} *w_6 &= w_6 - a (h_2 \cdot \Delta) \\ *w_5 &= w_5 - a (h_1 \cdot \Delta) \\ *w_4 &= w_4 - a (i_2 \cdot \Delta w_6) \\ *w_3 &= w_3 - a (i_1 \cdot \Delta w_6) \\ *w_2 &= w_2 - a (i_2 \cdot \Delta w_5) \\ *w_1 &= w_1 - a (i_1 \cdot \Delta w_5) \end{aligned}$$

Backpropagation: Example

- ▶ Summarized in matrix form

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \mathbf{a} \Delta \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \begin{bmatrix} a h_1 \Delta \\ a h_2 \Delta \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \mathbf{a} \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \cdot [w_5 \quad w_6] = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \begin{bmatrix} a i_1 \Delta w_5 & a i_1 \Delta w_6 \\ a i_2 \Delta w_5 & a i_2 \Delta w_6 \end{bmatrix}$$

Backpropagation: Multiple data points

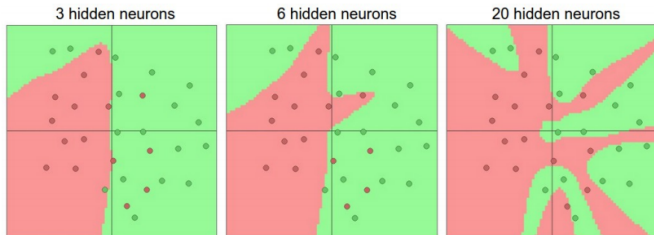
- ▶ Generally Δ is a vector, with the dimension of the number of training data points.
- ▶ The error can be the average of the error, so repeat the equations below for all training points and average the changes (the part after a)
- ▶ Fortunately numpy does not care about the number of dimensions, so instead of the multiplication in the right matrices we can use dot product.

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - a \Delta \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \begin{bmatrix} a h_1 \Delta \\ a h_2 \Delta \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - a \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \cdot [w_5 \quad w_6] = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \begin{bmatrix} a i_1 \Delta w_5 & a i_1 \Delta w_6 \\ a i_2 \Delta w_5 & a i_2 \Delta w_6 \end{bmatrix}$$

How many layers?

- ▶ Neural network with at least one hidden layer is a universal approximator (can represent any function).



Do Deep Nets Really Need to be Deep? Jimmy Ba, Rich Caruana,