

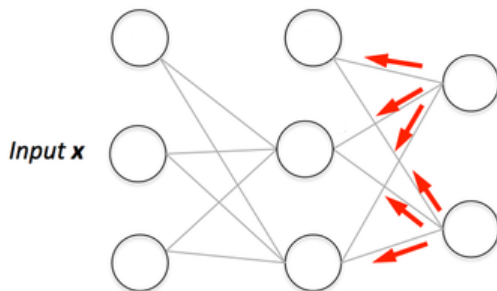
Artificial Intelligence

ANN – Back propagation

Jonathan Mwaura

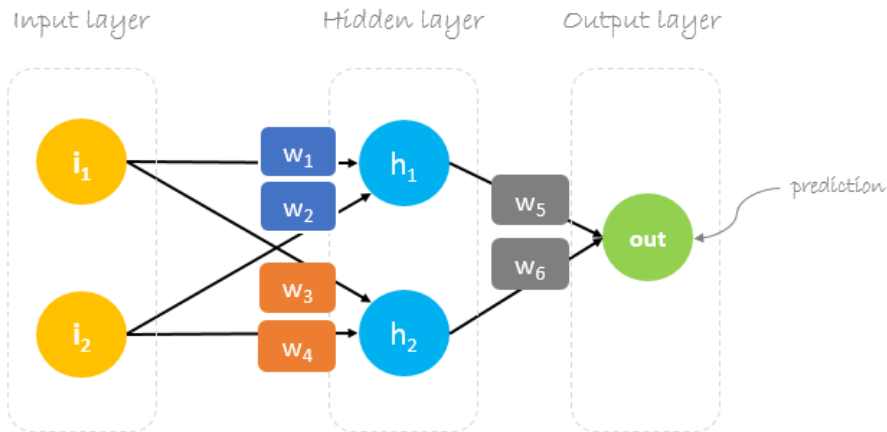
December 6, 2018

Backpropagation Step by Step

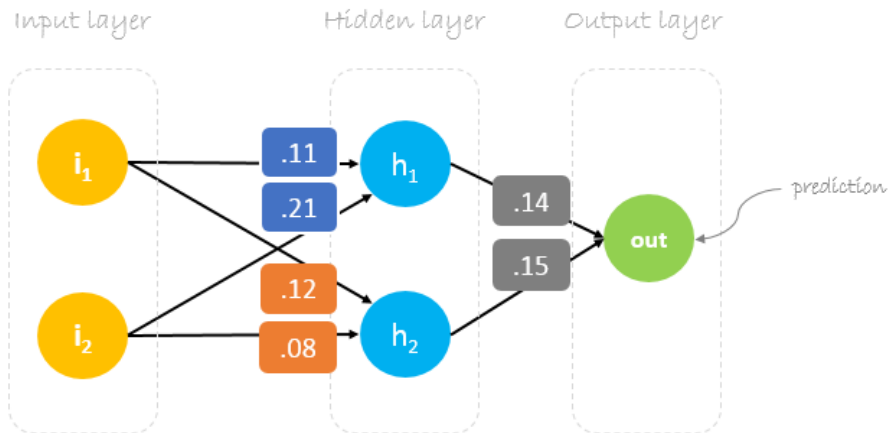


Backpropagation is the technique used to train an ANN.

Example overview



Weights, Weights, Weights



What exactly is NN Training?

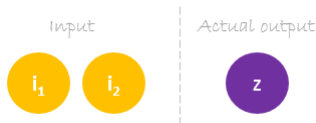
Neural network training is about finding weights that minimize prediction error.

Start training with a set of randomly generated weights.

Then, use backpropagation to update the weights in an attempt to correctly map arbitrary inputs to outputs.

Dataset

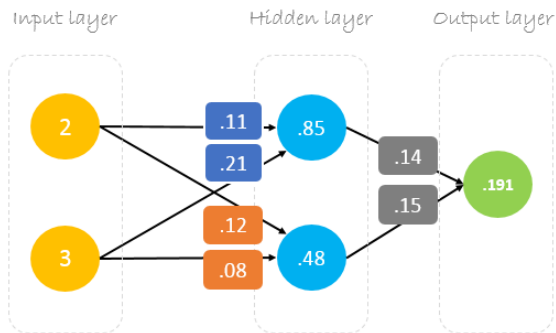
Our dataset has one sample with two inputs and one output.



Our single sample is as following inputs= $[2, 3]$ and output= $[1]$.



Forward pass



Forward Pass

$$\begin{bmatrix} 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 0.11 & 0.12 \\ 0.21 & 0.08 \end{bmatrix} = \begin{bmatrix} 0.85 & 0.48 \end{bmatrix} \cdot \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} = \begin{bmatrix} 0.191 \end{bmatrix}$$

Matrix multiplication

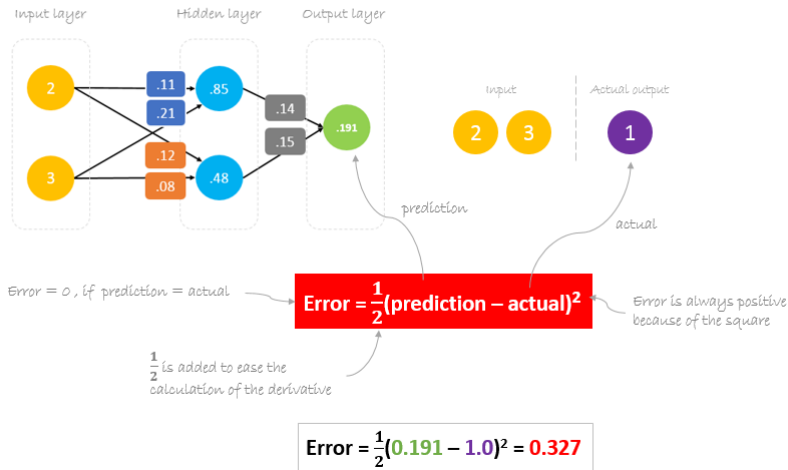
Details

$$2 \times .11 + 3 \times .21 = .85$$

$$.85 \times .14 + .48 \times .15 = .191$$

$$2 \times .12 + 3 \times .08 = .48$$

Calculate the Error



Reducing the error

prediction = out



prediction = $(h_1) w_5 + (h_2) w_6$



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

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**

The question now is how to change (update) the weights value so that the error is reduced? The answer is Backpropagation!

Reducing the error

The goal of the training is to reduce the error or the difference between prediction and actual output.

Actual output is constant, “not changing”, the only way to reduce the error is to change prediction value.

Decomposing prediction into its basic elements we can find that weights are the variable elements affecting prediction value.

In other words, in order to change prediction value, we need to change weights values.

So..Backpropagation..

Backpropagation, short for “backward propagation of errors” is a mechanism used to update the weights using gradient descent.

BP computes the gradient of the error function with respect to the neural network's weights.

The calculation proceeds backwards through the network.

Gradient Descent

an iterative optimization algorithm for finding the minimum of a function; in our case we want to minimize the error function. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point.

Weight Update Formula

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

Diagram illustrating the Weight Update Formula:

- $*W_x$: New weight
- W_x : Old weight
- a : Learning rate
- $\left(\frac{\partial \text{Error}}{\partial W_x} \right)$: Derivative of Error with respect to weight

Update Example

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

Optionally, we multiply the derivative of the error function by a selected number to make sure that the new updated weight is minimizing the error function; this is called **learning rate**.

Chain Rule

$$\begin{aligned}
 \frac{\partial \text{Error}}{\partial W_6} &= \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial W_6} \quad \leftarrow \text{chain rule} \\
 \text{Error} &= \frac{1}{2}(\text{prediction} - \text{actual})^2 \\
 \frac{\partial \text{Error}}{\partial W_6} &= \frac{1}{2}(\text{prediction} - \text{actual})^2 * \frac{\partial (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6}{\partial W_6} \\
 \text{prediction} &= (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6 \\
 \frac{\partial \text{Error}}{\partial W_6} &= 2 * \frac{1}{2}(\text{prediction} - \text{actual}) \frac{\partial (\text{prediction} - \text{actual})}{\partial \text{prediction}} * (i_1 w_3 + i_2 w_4) \quad \leftarrow h_2 = i_1 w_3 + i_2 w_4 \\
 \frac{\partial \text{Error}}{\partial W_6} &= (\text{prediction} - \text{actual}) * (h_2) \\
 \Delta &= \text{prediction} - \text{actual} \quad \leftarrow \text{delta} \\
 \frac{\partial \text{Error}}{\partial W_6} &= \Delta h_2
 \end{aligned}$$

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

Similarly, we can derive the update formula for w_5 and any other weights existing between the output and the hidden layer.

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

Weights associated with Hidden Layer

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

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

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

$$h_1 = i_1 w_1 + i_2 w_2$$

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

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

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

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

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

Moving backward to update w_1 , w_2 , w_3 and w_4 existing between input and hidden layer, the partial derivative for the error function with respect to w_1 , for example,

Weights associated with Hidden Layer

updated weights

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

We can rewrite the update formulas in matrices as following

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

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

Backward Pass

Using derived formulas we can find the new weights.

L

Learning rate: is a hyperparameter which means that we need to manually guess its value.

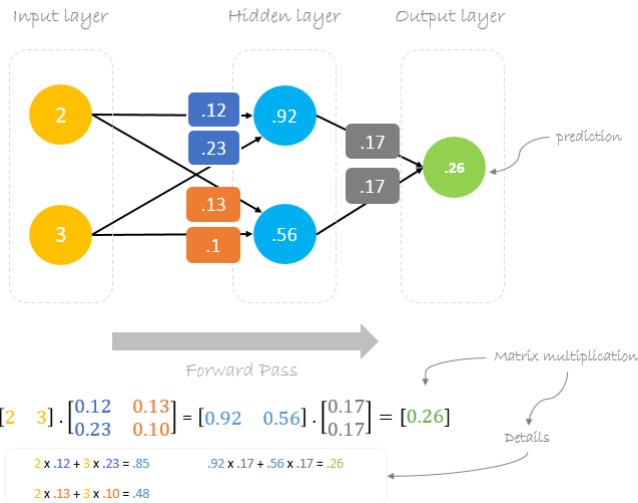
$$\Delta = 0.191 - 1 = -0.809 \quad \leftarrow \text{Delta} = \text{prediction} - \text{actual}$$

$$a = 0.05 \quad \leftarrow \text{Learning rate, we smartly guess this number}$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 0.85 \\ 0.48 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - \begin{bmatrix} -0.034 \\ -0.019 \end{bmatrix} = \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.14 & 0.15 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - \begin{bmatrix} -0.011 & -0.012 \\ -0.017 & -0.018 \end{bmatrix} = \begin{bmatrix} .12 & .13 \\ .23 & .10 \end{bmatrix}$$

Backward Pass



Consider a system with two states and two actions. You perform actions and observe the rewards and transitions listed below. Each step lists the current state, reward, action and resulting transition as S_i ; $R = r$; a_k : $S_i \rightarrow S_j$. Perform Q-learning using a learning rate of $\alpha = 0.5$ and a discount factor of $\gamma = 0.5$ for each step. The Q-table entries are initialized to zero.

Q	S_1	S_2
a_1	0	0
a_2	0	0

Table: S_1 ; $R = -10$; a_1 : $S_1 \rightarrow S_2$