

Machine Learning

Backpropagation 1

Mostafa S. Ibrahim

Teaching, Training and Coaching for more than a decade!

Artificial Intelligence & Computer Vision Researcher

PhD from Simon Fraser University - Canada

Bachelor / MSc from Cairo University - Egypt

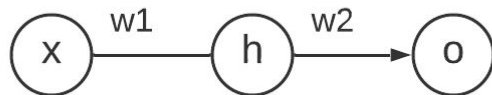
Ex-(Software Engineer / ICPC World Finalist)



© 2023 All rights reserved.

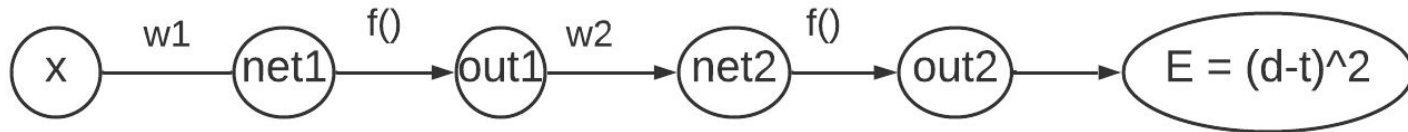
Please do not reproduce or redistribute this work without permission from the author

Recall a Trivial NN network



Assume h and o are followed by activation $f(a) = a^3$

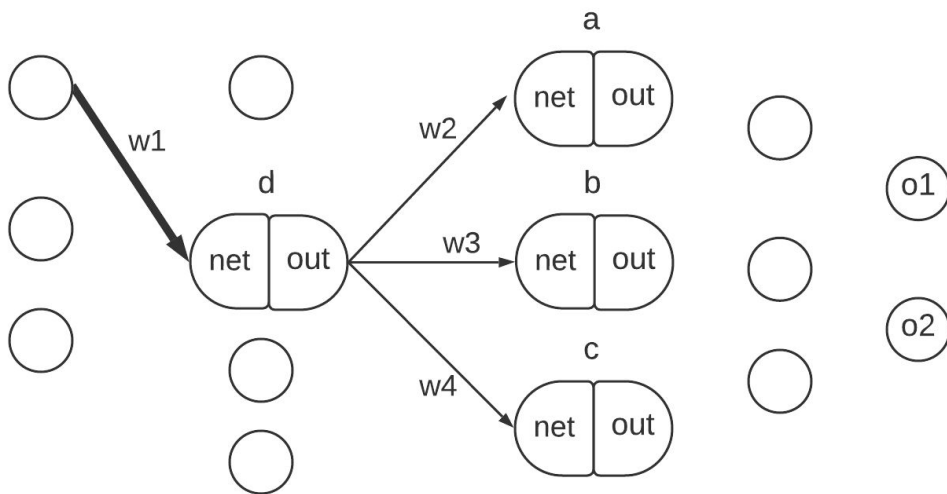
$E = (o-t)^2$
Compute $\partial E / \partial w1$



- $\partial E / \partial w1 = \partial E / \partial out2 * \partial out2 / \partial net2 * \partial net2 / \partial out1 * \partial out1 / \partial net1 * \partial net1 / \partial w1$
- If we **cached** the results of $\partial E / \partial net1$
then $\partial E / \partial w1 = \partial E / \partial net1 * \partial net1 / \partial w1$

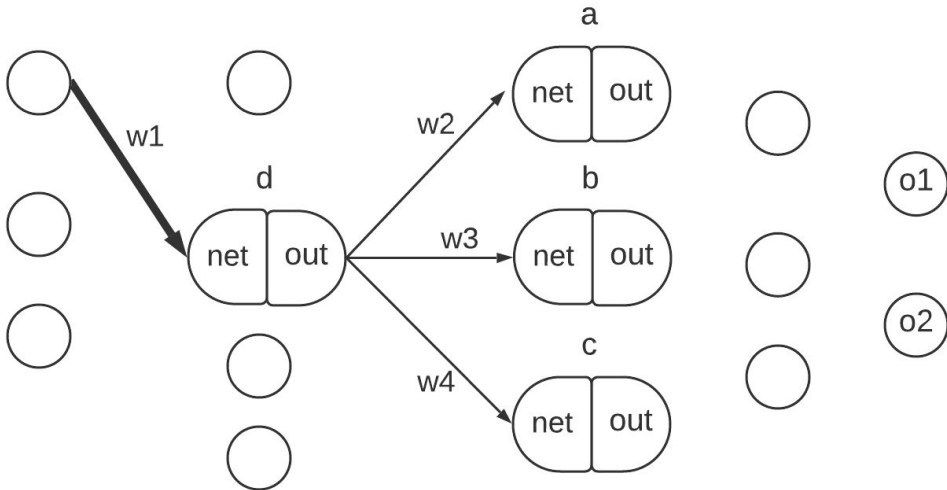
Compute $\partial E / \partial w_1$

- Assume we have 5-layers neural network
- It ends with 2 output nodes where the total error is computed using MSE
- We would like to compute $\partial E / \partial w_1$
- To do this, we consider all paths from w_1 to an output nodes
 - Empty nodes in layers 1 and 2 don't matter
- w_1 is linked to node **d**
- Node d is linked to nodes: **a, b, c**
- a, b, c are linked to more nodes that we should consider!



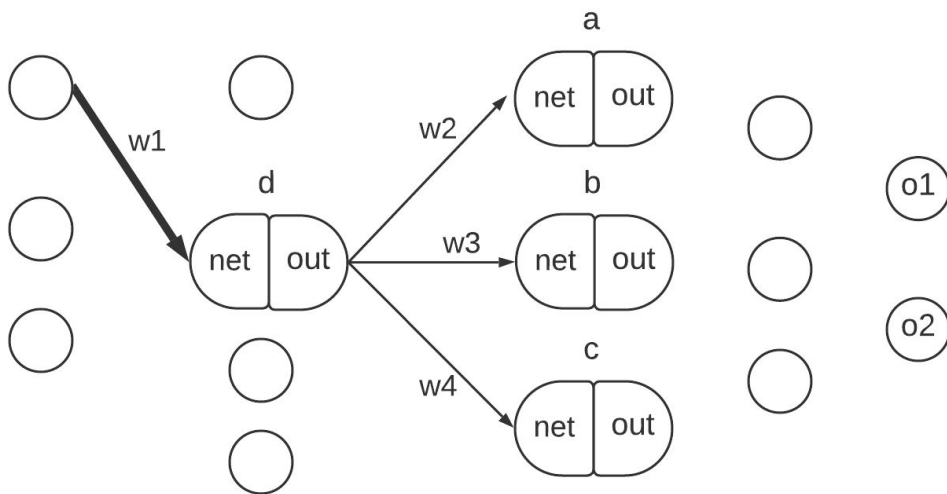
Compute $\partial E / \partial w_1$ using the Chain Rule

- $\partial E / \partial w_1 = \partial E / \partial d_net * \partial d_net / \partial w_1$
 - ?
- $\partial E / \partial d_net = \partial E / \partial d_out * \partial d_out / \partial d_net$
 - ?



Compute $\partial E / \partial w_1$ using the Chain Rule

- $\partial E / \partial w_1 = \partial E / \partial d_net * \partial d_net / \partial w_1$
 - $\partial d_net / \partial w_1 = x_1$ [net = $w_1 \mathbf{x_1} + w_2 x_2 + w_3 x_3$]
- $\partial E / \partial d_net = \partial E / \partial d_out * \partial d_out / \partial d_net$
 - $\partial d_out / \partial d_net$ is just the direct derivative of the activation function
 - If $f(net) = out = net^2$,
then $\partial d_out / \partial d_net = 2net$
 - If $f(net) = out = \text{sigmoid}(net)$,
then $\partial d_out / \partial d_net = out * (1-out)$
 - **Tip: Once we compute this value, we should cache it**

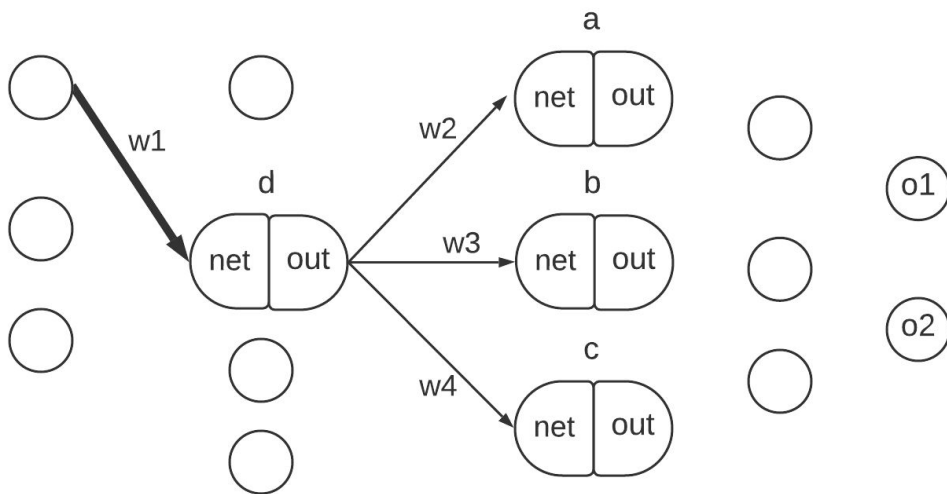


Compute $\partial E / \partial w_1$: Multivariate Chain Rule

- What about $\partial E / \partial d_{\text{out}}$?
- Node d_{out} is connected to a_{net} , b_{net} and c_{net}
- We need to process all derivative paths from d_{out} to them
- $\partial E / \partial d_{\text{out}} =$

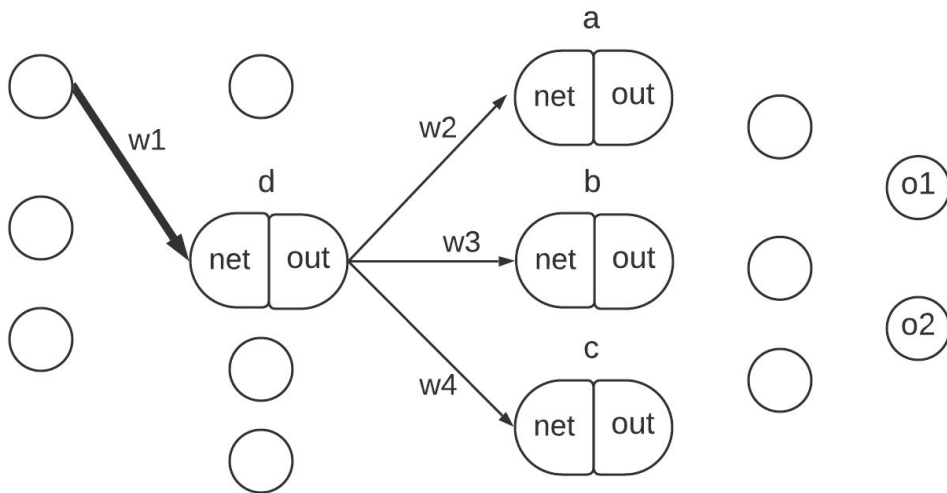
$$\begin{aligned} & \partial E / \partial a_{\text{net}} * \partial a_{\text{net}} / \partial d_{\text{out}} + \\ & \partial E / \partial b_{\text{net}} * \partial b_{\text{net}} / \partial d_{\text{out}} + \\ & \partial E / \partial c_{\text{net}} * \partial c_{\text{net}} / \partial d_{\text{out}} \end{aligned}$$

- $\partial a_{\text{net}} / \partial d_{\text{out}} = w_2$
- $\partial b_{\text{net}} / \partial d_{\text{out}} = w_3$
- $\partial c_{\text{net}} / \partial d_{\text{out}} = w_4$



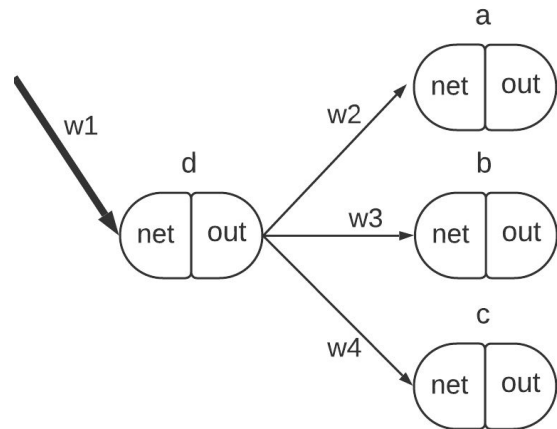
Compute $\partial E / \partial w_1$: Multivariate Chain Rule

- $\partial E / \partial d_{\text{out}} = \partial E / \partial a_{\text{net}} * w_2 + \partial E / \partial a_{\text{net}} * w_3 + \partial E / \partial a_{\text{net}} * w_4$
 - Let's generalize: $\sum \partial E / \partial \text{LayerNode_net} * \text{connection_weight}$
- Now we just moved from one layer to the next layer
- We can recursively keep doing that to enumerate all paths
- But what if we already **cached** these 3 values?
- Then we just use them!
- This is backpropagation!



Backpropagation derivative 3 formulas

- Assume we have interest in node d and outgoing edges of it are (weight = $w[i]$ and connection = $n[i]$)
- $\partial E / \partial d_out = \sum w[i] \times \partial E / \partial n[i]_net$ ($\partial n[i]_net$ cached)
- $\partial E / \partial d_net = \partial E / \partial d_out * \partial d_out / \partial d_net$
 - $\partial d_out / \partial d_net$ is just the direct derivative of the activation function
 - $out * (1-out)$ for sigmoid
 - $1 - out^2$ for tanh
 - $2net$ for quadratic polynomial
 - Cache $\partial E / \partial d_net$
- $\partial E / \partial w = \partial E / \partial d_net * \partial d_net / \partial w$
 - $\partial d_net / \partial w$ is the input value to this node
 - Either input x_i or node_out



Base Case

- We determined the general (recursive) formulas
- We need to compute the base case, which is the output layer
- Typically the activation function for regression in the output layer is identity
 - $\text{identity}(x) = x$
 - If you are regressing specific range, e.g. 0-1, you can use sigmoid
- Luckily, this is easy with regression
- Assume the error is MSE: $E = \frac{1}{2} (\text{output} - \text{target})^2$
 - This is vector notation. It is also the same for each output node O
- Then for some final $\partial E / \partial o_{\text{out}} = \frac{1}{2} * 2 * (\text{out} - \text{target}) * 1$

Relevant Materials

- Backpropagation: [StatQuest](#), 3Blue1Brown ([v1](#), [v2](#))
- A step-by-step Backpropagation Example: [Article](#)
- The Backpropagation [Algorithm](#)
- Yes you should [understand backprop](#)

“Acquire knowledge and impart it to the people.”

“Seek knowledge from the Cradle to the Grave.”

