

Loss function

can be sum of squared errors

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

↑ ↑
 given training value prediction

chain rule

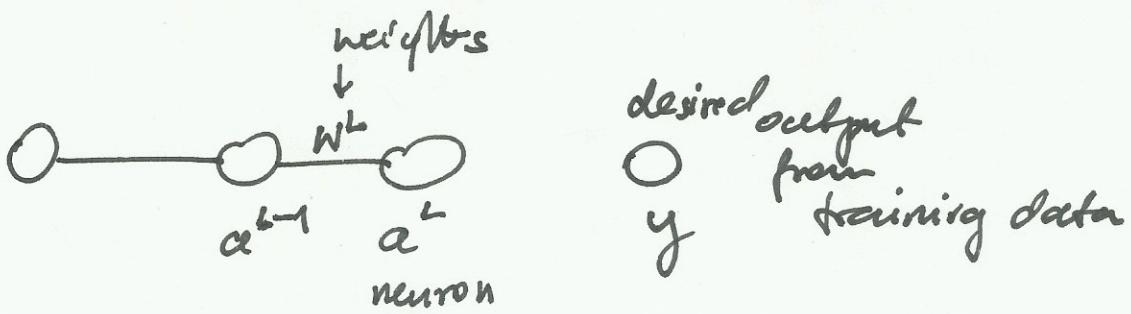
$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial W}, \text{ where } z = Wx + b$$

↑ ↑ ↑
 derivative by weights weights biases
 of neural network layers

After transformation

$$\frac{\partial L}{\partial W} = 2(y - \hat{y}) \cdot \text{derivative-sign}(z) \cdot \text{layer}^{L-1}$$

└ proof on the next page

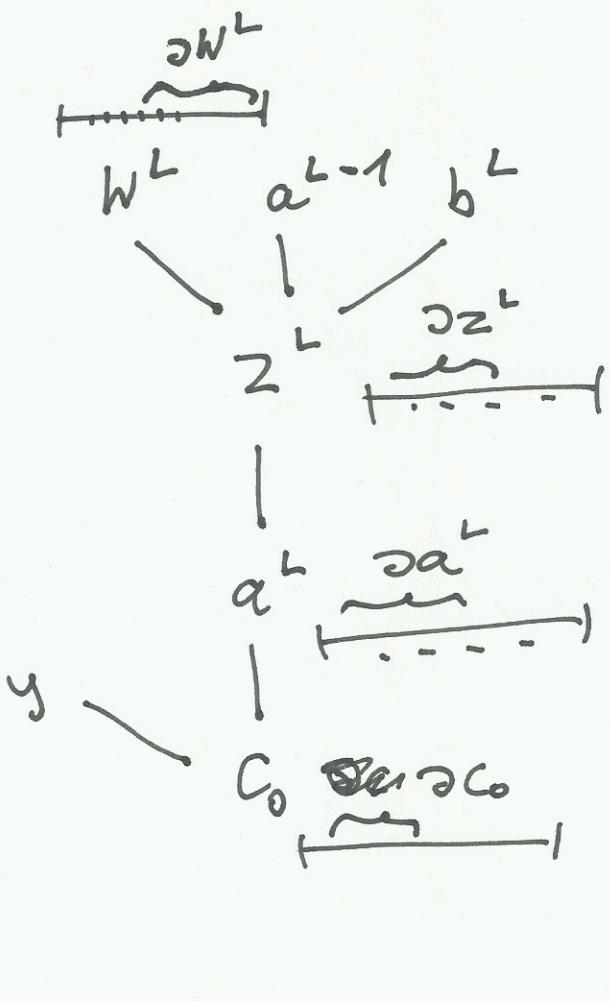


cost function

$$C_o = (\alpha^L - y)^2$$

value of neuron L $\underbrace{z^L}_{\alpha^L = G(w^L \cdot \alpha^{L-1} + b^L)}$; $w^L \cdot \alpha^{L-1} + b^L = z^L$

activation
function (here sigmoid)



we want to know
how changes in weights
influence cost function
(in other words predictor
accuracy)

$$\frac{\partial C_o}{\partial w^L} = \frac{\partial z^L}{\partial w^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_o}{\partial a^L}$$

as per
chain
rule

$$\begin{aligned} \frac{\partial C_o}{\partial a^L} &= \frac{\partial (\alpha^L - y)^2}{\partial \alpha^L} = \\ &= 2(\alpha^L - y) \end{aligned}$$

$$\frac{\partial \alpha^L}{\partial z^L} = \frac{\partial \sigma(z^L)}{\partial z^L} = \text{sigmoid-deriv}(z^L)$$

$$\frac{\partial z^L}{\partial w^L} = \frac{\partial (w^L \cdot \alpha^{L-1} + b^L)}{\partial w^L} = \alpha^{L-1}$$

so

$$\frac{\partial C_0}{\partial w^L} = 2(\alpha^L - y) \cdot \text{sigm.-deriv}(z^L) \cdot \alpha^{L-1}$$

for network on page 4
 A) output to hidden layers
 α^{L-1} = layer 1

$$z^L = w^L \alpha^{L-1} + b^L, \text{ assuming that } b^L = 0 \text{ we have}$$

$$z^L = w^L \cdot \alpha^{L-1} = \text{output}$$

that's why

layer 1.1 $\circ 2(\text{output} - y) \cdot \text{sigm.-deriv.}(\text{output})$
 makes sense layer 1 position went to front
 because of numpy dot function

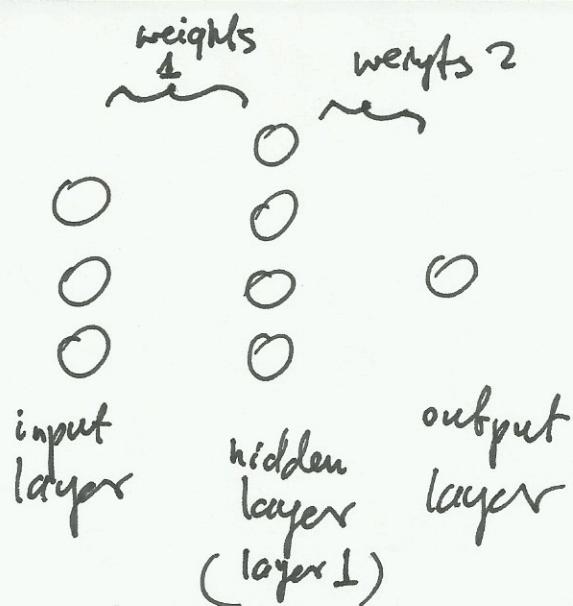
B) hidden layer to input

$$\alpha^{L-1} = \text{input}$$

$$\text{by analogy } z^L = \text{layer 1}$$

→ see p. 4

input.1 $\circ \{ \text{difference between prediction and true value}$
 moved by one layer $\cdot \text{sigmoid-deriv}(\text{layer 1}) \}$



change of weights 2:

$$d_2 = \text{layer 1.T} \circ 2(y - \hat{y}) \cdot \text{sigmoid-derivative}(\hat{y})$$

↑
 values kept
 in neurons
 at layer 1
 (must be transposed
 to correctly do operation
 on vectors)

$$\text{weights 2} = \text{weights 2} + d_2$$

change of weights 1:

$$d_1 = \text{input.T} \circ \left\{ \overbrace{\left[2(y - \hat{y}) \cdot \text{sigmoid-derivative}(\hat{y}) \circ \text{weights 2.T} \right]}^{\substack{\text{difference between} \\ \text{network prediction and the true} \\ \text{value moved to layer 1 by} \\ \text{weights 2 multiplication}}} \cdot \text{sigmoid-derivative (layer 1)} \right\}$$