

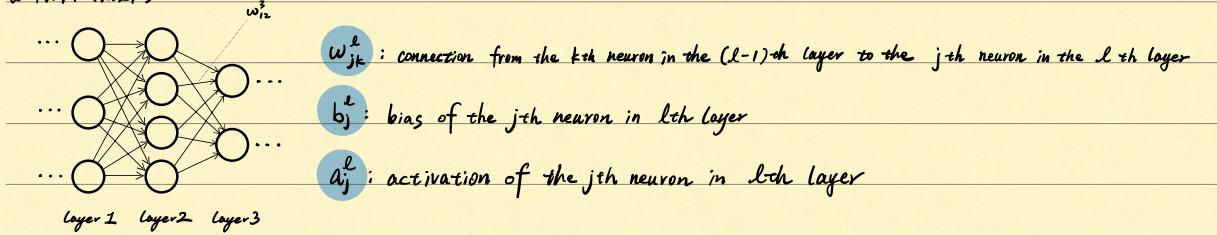
On Backpropagation

★ Goal: $\frac{\partial C}{\partial w}$ (or $\frac{\partial C}{\partial b}$) w.r.t. any weight w (or bias b) ; where C stands for Cost function

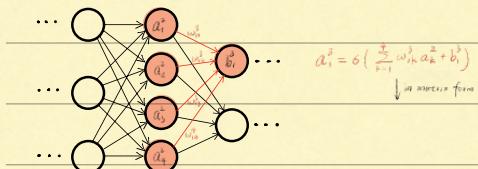
- How quickly the cost changes when we change the w & b
- Overall: BP gives us detailed insights into how changing the w & b changes overall behaviour of the network

NOTATIONS

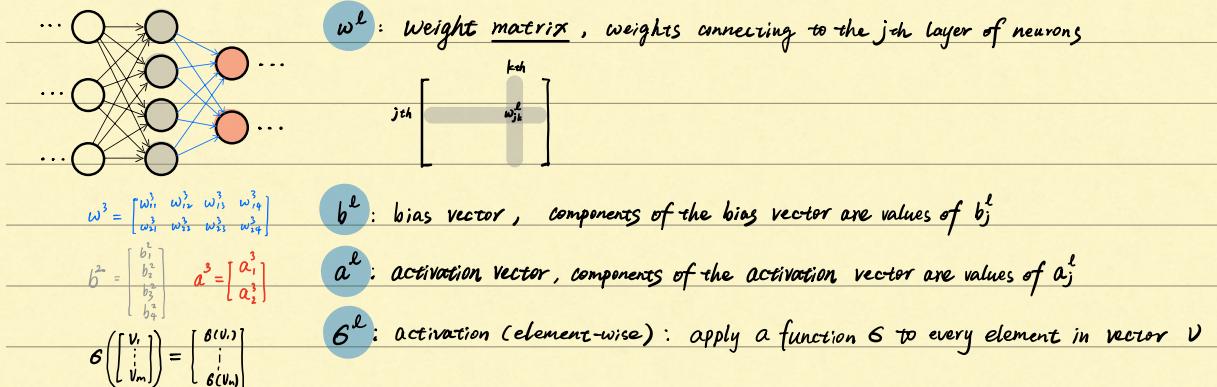
a NN: (MLP)



we have: $a_j^l = \sigma(\sum_k w_{jk}^l a_k^{l-1} + b_j^l)$ σ : activation function (e.g. Sigmoid)



vectorization:



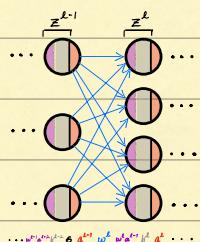
we have: $a^l = \sigma(w^l a^{l-1} + b^l) = \sigma(z^l)$

more math bias term

$$z^l := w^l a^{l-1} + b^l \quad (\text{intermediate quantity}) \quad z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l,$$

z^l is "Weighted input" to neurons in layer l

so far so good:



two assumptions we need about the cost function

example cost function:

$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$$

cost average over all training samples

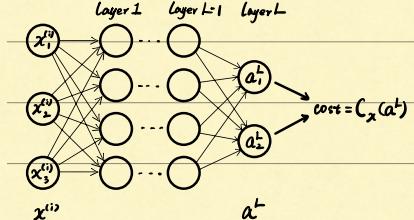
n: total number of training samples $\left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}, \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} \right)$

x: individual training sample: i-th sample $x^{(i)} \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_m^{(i)} \end{bmatrix}$

y(x): sample x's corresponding desired output (ground-truth / label ...)

L: number of layers in the network

$a^L = a^L(x)$: output vector from the network when x is input



1st assumption: cost function C can be written as an average $C = \frac{1}{n} \sum_x C_x$

over functions C_x for individual training sample x

$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2 = \frac{1}{n} \sum_x \frac{1}{2} \|y(x) - a^L(x)\|^2 C_x$$

BP computes $\frac{\partial C_x}{\partial w}$ ($\frac{\partial C_x}{\partial b}$) for a single sample

then we recover $\frac{\partial C}{\partial w}$ ($\frac{\partial C}{\partial b}$) by averaging over training sample

2nd assumption: the cost can be written as a function of the outputs from the NN.

so far so good: $C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$ is a function of our NN, as x, y(x) are fixed parameters.

when we change weight & bias (the parameters in NN), we'll change the output of NN, given x as input

The four fundamental equations behind backpropagation

recall: our goal is $\frac{\partial C}{\partial w}$ and $\frac{\partial C}{\partial b}$,

ultimately, we need to compute partial derivatives $\frac{\partial C}{\partial w_j^l}$ and $\frac{\partial C}{\partial b_j^l}$

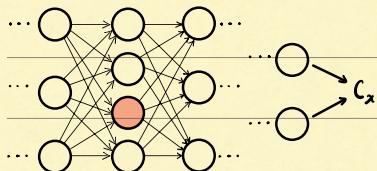
to do this, we introduce intermediate quantity δ_j^l : error in the j -th neuron in the l -th layer

BP 1° gives us a procedure to compute the error δ_j^l ,

2° and then will relate δ_j^l to $\frac{\partial C}{\partial w_{jk}^l}$ and $\frac{\partial C}{\partial b_j^l}$

layer l

Consider j -th neuron in layer l :



given a sample point x , we have weighted input z_j^l , and output $\sigma(z_j^l)$

now add a small change Δz_j^l to z_j^l , the output will be $\sigma(z_j^l + \Delta z_j^l)$

(we can use a_j^l as well, δ_j^l for convenience)

this change propagates through later layers in the network,

and finally causing the overall cost to change by an amount $\frac{\partial C}{\partial z_j^l} \Delta z_j^l$

(approx. in this case, Δz_j^l is a small value)

from $\frac{\partial C}{\partial z_j^l} \Delta z_j^l$, we can see that $\frac{\partial C}{\partial z_j^l}$ affects the magnitude of Δz_j^l

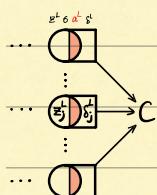
required to change the overall cost by the same amount

by this heuristic, we define $\delta_j^l := \frac{\partial C}{\partial z_j^l}$ as "a measure of the error in the neuron"

different meaning from "error" in a test (3.5% error) - for me its more of a "rate" thing.

δ^l : the vector of errors associated with layer l .

(BP1) An equation for error in the output layer, δ^L



the components of δ^L are given by

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \delta'(z_j^L)$$

measures how fast the loss is changing as a function of the j-th output activation
remember $L = \text{Clay}(n)$

Pf: Applying the chain rule:

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} = \sum_k \frac{\partial C}{\partial a_k^L} \frac{\partial a_k^L}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} \delta'(z_j^L)$$

$$\delta^L = \nabla_a C \odot \delta'(z^L)$$

• z_j^L can be obtained by one forward-pass on NN

• $\delta'(z_j^L)$ is easy to compute

• $\frac{\partial C}{\partial a_j^L}$ is easy to compute as the form of cost function C is provided (e.g. $C = \frac{1}{2} \|y(w-a^L(x))\|^2$)

write BP1 in matrix-based form:

$$\delta^L = \nabla_a C \odot \delta'(z^L) \quad , \text{ where } \odot \text{ stands for Hadamard product (elementwise product)}$$

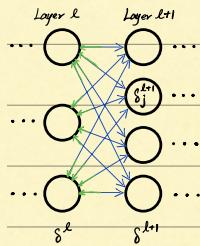
or in matrix form:

• $\nabla_a C$: gradient of C w.r.t. the output activations a^L $[\dots, \frac{\partial C}{\partial a_1^L}, \dots]$

$$\delta^L = \begin{bmatrix} \delta_1^L \\ \vdots \\ \delta_n^L \end{bmatrix} \cdot \nabla_a C \\ = \sum_i (z^L_i) \cdot \nabla_a C$$

e.g. for $C = \frac{1}{2} \|y(w-a^L(x))\|^2$, we have $\delta^L = (a^L - y) \odot \delta'(z^L)$

(BP2) An equation for the error δ^l in terms of the errors in the next layer, δ^{l+1}



In particular:

$$\delta^l = [(\omega^{l+1})^\top \delta^{l+1}] \odot \sigma'(\mathbf{z}^l)$$

"seeing the error backpropagated through the network"

gives us some kind of measure of error at the nodes of the l^{th} layer

or in matrix form: $\delta^l = \sum'(\mathbf{z}^l)(\omega^{l+1})^\top \delta^{l+1}$

Pf:

$$\begin{aligned}\delta_j^l &= \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1} \\ &= \sum_k \frac{\partial}{\partial z_j^l} \left[\sum_i w_{ki}^{l+1} \sigma(z_i^l) + b_k^{l+1} \right] \delta_k^{l+1} \\ &= \sum_k w_{kj}^{l+1} \sigma'(z_j^l) \delta_k^{l+1} \\ \delta^l &= [(\omega^{l+1})^\top \delta^{l+1}] \odot \sigma'(\mathbf{z}^l)\end{aligned}$$

By combining BP1 and BP2 we can compute the error δ^l for any layer in the network

1° compute δ^l using BP1

$$\delta^l = \sum'(\mathbf{z}^l) \nabla_a C$$

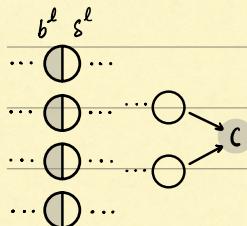
2° compute $\delta^{l-1}, \delta^{l-2}, \dots$ using BP2

$$\delta^{l-1} = \sum'(\mathbf{z}^l) (\omega^l)^\top \delta^l$$

$$= \sum'(\mathbf{z}^l) (\omega^l)^\top \sum'(\mathbf{z}^l) \nabla_a C$$

$$\delta^l = \sum'(\mathbf{z}^l) (\omega^{l+1})^\top \dots \sum'(\mathbf{z}^l) (\omega^l)^\top \sum'(\mathbf{z}^l) \nabla_a C$$

(BP3) An equation for the rate of change of the cost w.r.t. any bias in the network



In particular:

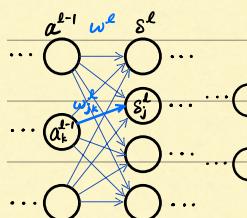
$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

δ_j^l is exactly equal to the rate of change $\frac{\partial C}{\partial b_j^l}$!

Pf:

$$\begin{aligned}\frac{\partial C}{\partial b_j^l} &= \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial b_j^l} \\ &= \delta_j^l \cdot \frac{\partial}{\partial b_j^l} \left[\sum_k w_{jk}^l \sigma(z_k^l) + b_k^l \right] \\ &= \delta_j^l\end{aligned}$$

(BP4) An equation for the rate of change of the cost w.r.t. any weights in the network



In particular:

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

Pf:

$$\begin{aligned}\frac{\partial C}{\partial w_{jk}^l} &= \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial w_{jk}^l} \\ &= \delta_j^l \cdot \frac{\partial}{\partial w_{jk}^l} \left[\sum_i w_{ik}^l a_i^{l-1} + b_k^l \right] \\ &= \delta_j^l a_k^{l-1}\end{aligned}$$

Summary: equations of backpropagation

$$\delta^L = \nabla_C C \odot \sigma'(z^L) \quad (\text{BP 1})$$

$$\delta^L = [(w^{L+1})^T \delta^{L+1}] \odot \sigma'(z^L) \quad (\text{BP 2})$$

$$\frac{\partial C}{\partial b_j} = \delta_j \quad (\text{BP 3})$$

$$\frac{\partial C}{\partial w_{jk}} = a_k^{L-1} \delta_j^L \quad (\text{BP 4})$$

Some insights on learning with BP

in BP1: term $\sigma'(z_j^L) \rightarrow 0$ when $z_j^L \rightarrow \infty$; a weight in the final layer will learn slowly if the neuron is either low activation (~ 0) or high activation (~ 1), it's common to say the output neuron has saturated, as a result the weight has stopped learning (or learning slowly)

in BP2: similar for earlier layers: δ_j^L is likely to get small if the neuron is near saturation

in BP4: when a_k^{L-1} is small, $\frac{\partial C}{\partial w_{jk}}$ also tends to be small

in context of gradient descent, the weights output from low-activation neurons learn slowly

affect

We can turn the observations around:

e.g. activation function σ s.t. σ' is always positive & never gets close to zero

this would prevent the slow-down of learning that occurs when ordinary sigmoid neurons saturate

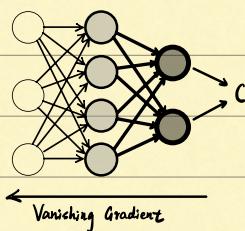
Vanishing/exploding gradient:

consider error vector δ^L , we have $\delta^L = \Sigma'(z^L)(w^{L+1})^T \dots \Sigma'(z^1)(w^1)^T \Sigma'(z^0) \nabla_C C$

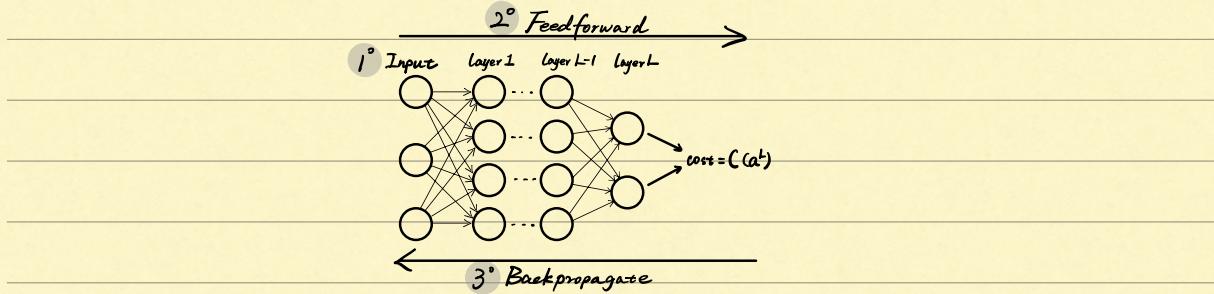
notably, there are multiple products of these derivatives of the activation function ($\Sigma'(z^L)$)

this observation and the insights above will lead to the problem of the vanishing/exploding gradient

which is a fundamental issue arising from NN's architecture and the choice of activation functions



The Backpropagation algorithm



1° Input x : set the corresponding activation a' for the input layer

2° Feedforward: For each $l = 2, 3, \dots, L$, compute $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$

3° Output error δ^L : Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$

4° Backpropagate the error: For each $l = L-1, L-2, \dots, 2$, compute $\delta^l = [(w^{l+1})^T \delta^{l+1}] \odot \sigma'(z^l)$

5° Output: The gradient of the cost function is given by $\frac{\partial C}{\partial w_{jk}^l} = a_j^{l-1} \delta_j^l$ and $\frac{\partial C}{\partial b_j^l} = \delta_j^l$

Combining BP with mini-batch SGD:

given a mini-batch of m training examples :

1° Input a set of training examples

2° for each training example x :

- set the corresponding activation $\tilde{a}^{x,l}$ for the input layer

- Feedforward: For each $l = 2, 3, \dots, L$, compute $\tilde{z}^{x,l} = w^l \tilde{a}^{x,l-1} + b^l$ and $\tilde{a}^{x,l} = \sigma(\tilde{z}^{x,l})$

- Output error $\delta^{x,l}$: Compute the vector $\delta^{x,l} = \nabla_a C \odot \sigma'(\tilde{z}^{x,l})$

- Backpropagate the error: For each $l = L-1, L-2, \dots, 2$, compute $\delta^{x,l} = [(w^{l+1})^T \delta^{x,l+1}] \odot \sigma'(\tilde{z}^{x,l})$

3° Gradient descent: For each $l = L, L-1, \dots, 2$,

update the weights : $w^l \rightarrow w^l - \eta \cdot \frac{1}{m} \sum_x \delta^{x,l} (\tilde{a}^{x,l-1})^T$

biases : $b^l \rightarrow b^l - \eta \cdot \frac{1}{m} \sum_x \delta^{x,l}$

matrix-based approach

represent a mini-batch input by a matrix $X = [x^{(1)}, \dots, x^{(m)}] = \begin{bmatrix} x_1^{(1)} & \cdots & x_1^{(m)} \\ \vdots & & \vdots \\ x_n^{(1)} & \cdots & x_n^{(m)} \end{bmatrix}$

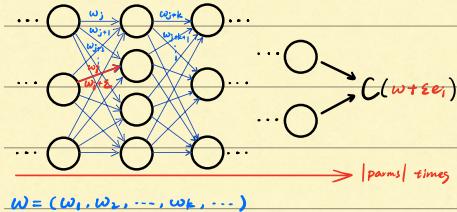
$$a' = [w' \cdot \sigma(X) + b']$$

introduce a new axis "↔" $\left[\dots \overset{\leftrightarrow}{a_j^{x,1}} \dots \right]_{out \times m} = \left[\overset{\leftrightarrow}{w_{ij}}, \dots, \overset{\leftrightarrow}{w_{ij}} \right]_{out \times n} \begin{bmatrix} \sigma(x_1^{(1)}) & \cdots & \sigma(x_1^{(m)}) \\ \vdots & & \vdots \\ \sigma(x_n^{(1)}) & \cdots & \sigma(x_n^{(m)}) \end{bmatrix}_{n \times m} + \left[\overset{\leftrightarrow}{b_j} \dots \overset{\leftrightarrow}{b_j} \right]_{out \times m}$

$$A^l = \sigma(w^l A^{l-1} + b^l)$$

In what sense is backpropagation a fast algorithm

another approach to computing the gradient:



still, our goal is computing $\frac{\partial C}{\partial w_j}$ for every j .

by definition, we have $\frac{\partial C}{\partial w_j} = \lim_{\varepsilon \rightarrow 0^+} \frac{C(w + \varepsilon e_i) - C(w)}{\varepsilon}$, where e_i is the unit vector in j -th direction.

now we can estimate $\frac{\partial C}{\partial w_j}$ by $\frac{C(w_j) - C(w'_j)}{w_j - w'_j}$, where w_j, w'_j are 2 slightly differ values of w_j

and it's same for $\frac{\partial C}{\partial b}$.

a promising method, but the implementation turns out to be extremely slow.

e.g. a NN with million weights. then for each w_j , we need to compute $C(w + \varepsilon e_i)$

thus we needs to compute the cost function a million different times, requiring a million forward passes through the network (per training example)!

In contrast, backpropagation enables us to compute all partial derivatives $\frac{\partial C}{\partial w_j}$

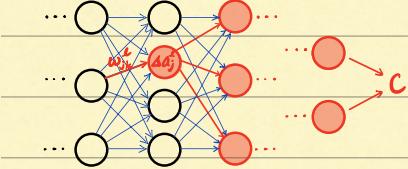
using just one forward pass through the network,

followed by one backward pass through the network.

Backpropagation : the big picture

more intuition how changing the w & b changes overall behaviour of the network?

when we add a small change to w_{jk}^l , $w_{jk}^l \rightarrow w_{jk}^l + \Delta w_{jk}^l$, this will cause changes in following layers



now focus on the change ΔC in the cost, $\Delta C \approx \frac{\partial C}{\partial w_{jk}^l} \Delta w_{jk}^l$

this suggest that a possible approach to computing $\frac{\partial C}{\partial w_{jk}^l}$:

track how a small change in w_{jk}^l propagates to cause a small change in C .

Δw_{jk}^l causes a small change Δa_j^l in the activation of the j -th neuron in the l -th layer:

$$\Delta a_j^l \approx \frac{\partial a_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$$

Δa_j^l will cause changes in all the activations in the next layer, consider the q -th neuron:

$$\Delta a_q^{l+1} \approx \frac{\partial a_q^{l+1}}{\partial a_j^l} \Delta a_j^l$$

$$\approx \frac{\partial a_q^{l+1}}{\partial a_j^l} \frac{\partial a_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$$

now consider a single "path": $a_j^l, a_q^{l+1}, \dots, a_n^{L-1}, a_m^L$, we have

$$\Delta C_0 \approx \frac{\partial C}{\partial a_m^L} \frac{\partial a_m^L}{\partial a_{n-1}^{L-1}} \frac{\partial a_{n-1}^{L-1}}{\partial a_p^{L-2}} \dots \frac{\partial a_p^{L-1}}{\partial a_j^l} \frac{\partial a_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$$



for total change in C :

$$\Delta C \approx \sum_{m,n,\dots,q} \frac{\partial C}{\partial a_m^L} \frac{\partial a_m^L}{\partial a_{n-1}^{L-1}} \frac{\partial a_{n-1}^{L-1}}{\partial a_p^{L-2}} \dots \frac{\partial a_p^{L-1}}{\partial a_j^l} \frac{\partial a_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$$

$$\Rightarrow \frac{\partial C}{\partial w_{jk}^l} = \sum_{m,n,\dots,q} \frac{\partial C}{\partial a_m^L} \frac{\partial a_m^L}{\partial a_{n-1}^{L-1}} \frac{\partial a_{n-1}^{L-1}}{\partial a_p^{L-2}} \dots \frac{\partial a_p^{L-1}}{\partial a_j^l} \frac{\partial a_j^l}{\partial w_{jk}^l}$$

$$= \sum_{m,n,\dots,q} \underbrace{\frac{\partial C}{\partial a_m^L} \delta'(z_m^L) w_{mn}^L}_{\text{red bracket}} \cdot \underbrace{\delta'(z_{n-1}^{L-1}) w_{n-1,n}^{L-1}}_{\text{red bracket}} \cdots \underbrace{\dots}_{\text{red bracket}} \underbrace{\delta'(z_q^{L+1}) w_{qj}^{L+1}}_{\text{red bracket}} \cdot \underbrace{\delta'(z_j^l) a_{j-1}^{l-1}}_{\text{red bracket}}$$

$$= \sum_m \underbrace{\frac{\partial C}{\partial a_m^L} \delta'(z_m^L)}_{\text{red bracket}} \sum_n w_{mn}^L \cdot \delta'(z_{n-1}^{L-1}) \sum_p \dots \sum_q \underbrace{w_{qj}^{L+1} \delta'(z_q^{L+1})}_{\text{red bracket}} \cdot \underbrace{w_{qj}^{L+1} \delta'(z_j^l) a_{j-1}^{l-1}}_{\text{red bracket}}$$

$$= \sum_m \delta_m^L \sum_n w_{mn}^L \cdot \delta'(z_{n-1}^{L-1}) \sum_p \dots \sum_q \underbrace{w_{qj}^{L+2} \delta'(z_q^{L+1})}_{\text{red bracket}} \cdot \underbrace{w_{qj}^{L+1} \delta'(z_j^l) a_{j-1}^{l-1}}_{\text{red bracket}}$$

$$= \sum_n \delta_n^{L-1} \sum_p \dots \sum_q \underbrace{w_{qj}^{L+2} \delta'(z_q^{L+1})}_{\text{red bracket}} \cdot \underbrace{w_{qj}^{L+1} \delta'(z_j^l) a_{j-1}^{l-1}}_{\text{red bracket}}$$

$$= \sum_q \delta_q^{L-1} \cdot w_{qj}^{L+1} \delta'(z_j^l) a_{j-1}^{l-1}$$

$$= \delta_k^{L-1} a_k^{L-1}$$

this is exactly BP 2!

In this context, BP is a "layer-wise" order to compute the summation!