# Error backpropagation

goal: compute 1st & 2nd derivatives of $E(\vec{w})$ [wrt $\vec{w}$] efficiently.

Consider a general NN with arbitrary feed-forward topology, arbitrary activation f's and a general error f'n:

$$E(\vec{w}) = \sum_{n=1}^{N} E_n(\vec{w})$$

Consider first a linear model:

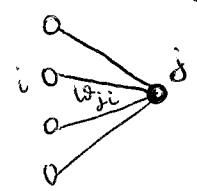$$\begin{cases} y_k = \sum_i w_{ki} x_i \, , \\ E_n = \frac{1}{2} \sum_k (y_{nk} - t_{nk})^2 \\ \quad\quad\quad\quad \underset{y_k(\vec{x}_n, \vec{w})}{} \end{cases}$$

Then $\quad \dfrac{\partial E_n}{\partial w_{ji}} = \sum_k (y_{nk} - t_{nk}) \underbrace{\dfrac{\partial y_{nk}}{\partial w_{ji}}}_{\delta_{kj} x_{ni}} =$

$$= \underbrace{(y_{nj} - t_{nj})}_{\text{"error signal"}} \underbrace{x_{ni}}_{\text{"input"}}$$

More generally,

$$a_j = \sum_i w_{ji} \underbrace{z_i}_{}$$

 activation of unit $i$ connected to unit $j$

Then $z_j = h(a_j)$

  ↑ activation of unit $j$

Now consider $\dfrac{\partial E_n}{\partial w_{ji}} = \underbrace{\dfrac{\partial E_n}{\partial a_j}}_{\overset{\text{\tiny III}}{\delta_j}} \underbrace{\dfrac{\partial a_j}{\partial w_{ji}}}_{\text{``}z_i}$ ⊜ (no sum!)

⊜ $\delta_j z_i$.

  ($z_i = 1$ if we are dealing with bias)

If $j$ is an output unit,

$$\delta_j = \dfrac{\partial E_n}{\partial a_j} = \underbrace{y_j - t_j}_{\substack{n \text{ index} \\ \text{omitted for simplicity}}} \quad \text{for regression OR} \\ \text{classification}$$

Indeed, for regression

$y_k = a_k$  (unit activation function)

  ↳ $\dfrac{\partial E}{\partial a_k} = y_k - t_k$ .

For $k=2$ classification,

$$\dfrac{\partial E_n}{\partial a_j} = - \left[ \dfrac{t_j}{y_j} \underbrace{\dfrac{\partial y_j}{\partial a_j}}_{} + (1-t_j)(-1)\dfrac{1}{1-y_j} \dfrac{\partial y_j}{\partial a_j} \right] ⊜$$

  ↗ $y_j(1-y_j)$

$y_j = \sigma(a_j)$ , sigmoid

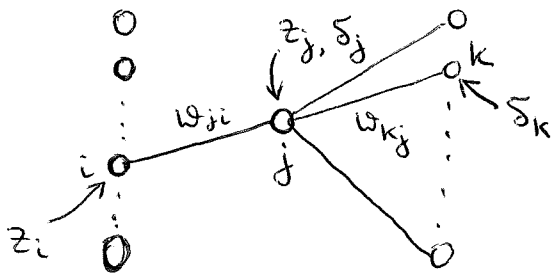⊜ $- \left[ t_j(1-y_j) - (1-t_j)y_j \right] = y_j - t_j$ .

Same for $N$ $K=2$ classif's and
$K>2$ classif'n.

———o———

If $j$ is a hidden unit,

$$\delta_j = \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k}\frac{\partial a_k}{\partial a_j} \quad \diamondsuit$$

↑
over all "future" units, hidden or
output, to which $j$ is connected



$$\diamondsuit = \sum_k \delta_k \frac{\partial}{\partial a_j}\left(\sum_i w_{ki} h(a_i)\right) = \sum_{k,i} \delta_k w_{ki} \underbrace{\frac{\partial h(a_i)}{\partial a_j}}_{h'(a_j)\delta_{ij}} =$$

$$= h'(a_j) \sum_k w_{kj} \delta_k. \qquad (*)$$

We can use $(*)$ by starting from
output units (for which $\delta_j$ is easily
computed) and computing $\delta_j$'s for hidden
units in a <u>backpropagation</u> pass.

Finally, use $\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i$ to
compute all gradients.

If needed, finish with

$$\frac{\partial E}{\partial w_{ji}} = \sum_n \frac{\partial E_n}{\partial w_{ji}}.$$

The above derivation can be easily generalized to several types of activation functions $h(\cdot)$.

Ex. Consider a two-layer network as before, and focus on regression:

$$\text{output units} \implies y_k = a_k$$

$$\text{hidden units} \implies h(a_k) = \tanh(a_k)$$

Note that $\dfrac{d\tanh(a)}{da} = 1 - \tanh^2(a)$

For pattern $n$, $E_n = \dfrac{1}{2} \sum_{k=1}^{K} (y_k - t_k)^2$   ← # output units

Forward pass:

$$\begin{cases} a_j = \sum_{i=0}^{D} w_{ji}^{(1)} x_i & \leftarrow \text{\# input units} \\ \\ \qquad\qquad \uparrow \text{inputs} \\ \\ z_j = \tanh(a_j) \\ \\ \qquad \swarrow \text{\# hidden units} \\ \\ y_k = \sum_{j=0}^{M} w_{kj}^{(2)} z_j \end{cases}$$

Next, compute $\delta_k = y_k - t_k$

Use (*) to find

$$\delta_j = (1 - z_j^2) \sum_{k=1}^{K} w_{kj}^{(2)} \delta_k \qquad \text{for hidden units}$$

Finally, compute

$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \delta_j x_i \quad , \quad \frac{\partial E_n}{\partial w_{kj}^{(2)}} = \delta_k z_j \qquad (**)$$

and sum over $n$.

Note that backpropagation is $\theta(W)$, where $W$ is the total # weights & biases. Most effort goes into evaluating

$$a_j = \sum_i w_{ji} z_i$$

Forward propagation (computing $E_n$) is $\theta(W)$ as well. Note that finite difference:

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \varepsilon) - E_n(w_{ji} - \varepsilon)}{2\varepsilon} + \theta(\varepsilon^2)$$

is $\theta(W^2)$ since each $E_n$ computation is $\theta(W)$ & one needs $W$ of those. However, it can be used for numerical checks.

We can check these results by evaluating derivatives of $E_n$ without using backpropagation:

$$E_n = \frac{1}{2} \sum_{k=1}^{K} \underbrace{(y_k - t_k)^2}_{\text{index } n \text{ omitted for clarity}}, \quad \text{where}$$

$$y_k = \sum_{j=0}^{M} w_{kj}^{(2)} \tanh\left(\sum_{i=0}^{D} w_{ji}^{(1)} x_i\right).$$

Then $\dfrac{\partial E_n}{\partial w_{k'j'}^{(2)}} = \sum_{k=1}^{K} (y_k - t_k) \dfrac{\partial y_k}{\partial w_{k'j'}^{(2)}} =$

$$= \sum_{k=1}^{K} (y_k - t_k) \sum_{j=0}^{M} \underbrace{\frac{\partial w_{kj}^{(2)}}{\partial w_{k'j'}^{(2)}}}_{\delta_{kk'}\,\delta_{jj'}} \underbrace{\tanh\left(\underbrace{\sum_{i=0}^{D} w_{ji}^{(1)} x_i}_{a_j}\right)}_{z_j} \ \textcircled{=}$$

$$\textcircled{=}\ \underbrace{(y_{k'} - t_{k'})}_{\delta_{k'}} z_{j'} = \underline{\underline{\delta_{k'} z_{j'}}}, \quad \text{same as } (**)$$

Next, $\dfrac{\partial E_n}{\partial w_{j'i'}^{(1)}} = \sum_{k=1}^{K} (y_k - t_k) \dfrac{\partial y_k}{\partial w_{j'i'}^{(1)}} =$

$$= \sum_{k=1}^{K} \underbrace{(y_k - t_k)}_{\delta_k} \sum_{j=0}^{M} w_{kj}^{(2)} \underbrace{\tanh'(a_j)}_{1 - z_j^2} \sum_{i=0}^{D} x_i \underbrace{\frac{\partial w_{ji}^{(1)}}{\partial w_{j'i'}^{(1)}}}_{\delta_{jj'}\,\delta_{ii'}} \boxed{=}$$

$$\boxdot \sum_{k=1}^{K} \delta_k \, w_{kj'}^{(2)} \, (1 - z_{j'}^2) \, x_{i'} =$$

$$= x_{i'} \, (1 - z_{j'}^2) \underbrace{\sum_{k=1}^{k} \delta_k \, w_{kj'}^{(2)}}_{\delta_{j'}} = \delta_{j'} \, x_{i'} \, , \quad \text{same as } (**)$$

# Jacobian matrix

$$J_{ki} = \frac{\partial y_k}{\partial x_i}$$ 
local sensitivity of outputs to changes in inputs

$$\Delta y_k = \sum_i \underbrace{\frac{\partial y_k}{\partial x_i}}_{J_{ki}} \underbrace{\Delta x_i}_{\substack{\text{small} \\ \text{input perturbations}}} \qquad \left| \frac{\Delta x_i}{x_i} \right| \ll 1$$

$J_{ki}$ can be evaluated using backpropagation:

$$J_{ki} = \frac{\partial y_k}{\partial x_i} = \sum_j \frac{\partial y_k}{\partial a_j} \underbrace{\frac{\partial a_j}{\partial x_i}}_{w_{ji}} \qquad (*)$$

over all "future" units $j$ connected to unit $i$

Further, $\dfrac{\partial y_k}{\partial a_j} = \sum_\ell \dfrac{\partial y_k}{\partial a_\ell} \dfrac{\partial a_\ell}{\partial a_j}$ ⊜

over all "future" units $\ell$ connected to unit $j$

$$⊜ \sum_\ell \frac{\partial y_k}{\partial a_\ell} \frac{\partial}{\partial a_j} \sum_i w_{\ell i} \underbrace{h(a_i)}_{z_i} =$$

$$= \sum_{\ell, i} \frac{\partial y_k}{\partial a_\ell} w_{\ell i} h'(a_i) \delta_{ij} = h'(a_j) \sum_\ell w_{\ell j} \frac{\partial y_k}{\partial a_\ell}.$$

Start from output units, e.g. with sigmoid activation functions:

$$\frac{\partial y_k}{\partial a_j} = \sigma'(a_k) \underbrace{\frac{\partial d_k}{\partial a_j}}_{\delta_{kj}} = \sigma(a_j)(1-\sigma(a_j))\,\delta_{kj}.$$

## Algorithm:

1. Choose $\bar{x}_n$ and forward-propagate

2. For each $k$, start with output units and backpropagate $\Rightarrow$ find $\frac{\partial y_k}{\partial a_j}$, $\forall j$

3. Find $J_{ki}$ using (*)

Can be checked against finite difference:

$$\frac{\partial y_k}{\partial x_i} = \frac{y_k(x_i + \varepsilon) - y_k(x_i - \varepsilon)}{2\varepsilon} + \mathcal{O}(\varepsilon^2)$$

needs 2D forward propagations $(D = \#\text{ inputs})$

### Hessian matrix

$$H_{ij} = \frac{\partial^2 E}{\partial w_i \, \partial w_j}$$

$i, j = 1, \dots, W \underset{\uparrow}{\text{total \# weights \& biases}}$

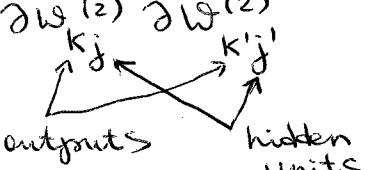all weights/biases relabeled using one consecutive index.

(1) $H_{ij}$ used in some non-linear optimization algorithms

(2) $H_{ij}$ is used in Bayesian neural networks

# Exact evaluation:

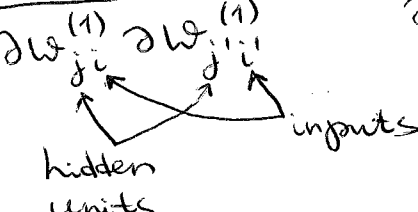Consider a two-layer network for simplicity.

Define $\quad \delta_k = \dfrac{\partial E_n}{\partial a_k}$, $\quad M_{kk'} = \dfrac{\partial^2 E_n}{\partial a_k \, \partial a_{k'}}$

$\underbrace{\qquad\qquad}_{\text{as before}}$ here, $k$ & $k'$ can be ~~...~~ hidden or output nodes

(a) Both weights in the second layer: (hidden)

$$\frac{\partial^2 E_n}{\partial w_{kj}^{(2)} \, \partial w_{k'j'}^{(2)}} = \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} \frac{\partial a_k}{\partial w_{kj}^{(2)}} \frac{\partial a_{k'}}{\partial w_{k'j'}^{(2)}} =$$
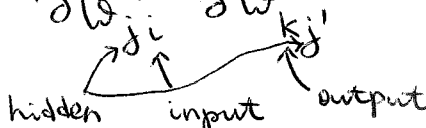
$w_{kj}^{(2)}$ ↑ outputs ↗ ↖ hidden units $w_{k'j'}^{(2)}$

$$= M_{kk'} \, z_j \, z_{j'}.$$

(b) Both weights in the first (input) layer:

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \, \partial w_{j'i'}^{(1)}} = \frac{\partial^2 E_n}{\partial a_j \, \partial a_{j'}} \underbrace{\frac{\partial a_j}{\partial w_{ji}^{(1)}}}_{x_i} \underbrace{\frac{\partial a_{j'}}{\partial w_{j'i'}^{(1)}}}_{x_{i'}} \textcircled{=} \text{~~...~~}$$

$w_{ji}^{(1)}$ ↑ hidden units ↗ ↖ inputs $w_{j'i'}^{(1)}$

$$\textcircled{=} \; M_{jj'} \, x_i \, x_{i'}.$$

(c) One weight in each layer:

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \, \partial w_{kj'}^{(2)}} = \frac{\partial}{\partial w_{ji}^{(1)}} \left[ \frac{\partial E_n}{\partial a_k} \underbrace{\frac{\partial a_k}{\partial w_{kj'}^{(2)}}}_{z_{j'}} \right] =$$

$w_{ji}^{(1)}$ ↗ ↑ hidden   input   ↖ output $w_{kj'}^{(2)}$

$$= \underbrace{\frac{\partial^2 E_n}{\partial a_j \partial a_k}}_{M_{jk} = M_{kj}} x_i \, z_{j'} + \underbrace{\frac{\partial E_n}{\partial a_k}}_{\delta_k} \frac{\partial z_{j'}}{\partial w_{ji}^{(1)}} \textcircled{=}$$

$$\textcircled{=} \; M_{jk} \, x_i \, z_{j'} + \delta_k \, h'(a_j) \, \delta_{jj'} \, x_i.$$

Thus, we need to backpropagate $M_{kk'}$.

Indeed,
$$\frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial \overbrace{a_k}^{=\sum_\ell w_{k\ell} h(a_\ell)}}{\partial a_j} = h'(a_j) \sum_k \frac{\partial E_n}{\partial a_k} w_{kj}$$

$\uparrow$ hidden $\Big($ downstream of & connected to $j$

$$\underbrace{\sum_\ell w_{k\ell} h'(a_\ell) \frac{\partial a_\ell}{\partial a_j}}_{\displaystyle = w_{kj} h'(a_j)} = $$

$$= w_{kj} h'(a_j). \qquad \delta_{\ell j}$$

Now,
$$\frac{\partial^2 E_n}{\partial a_i \partial a_j} = h''(a_i)\, \delta_{ij} \sum_k \frac{\partial E_n}{\partial a_k} w_{kj} +$$

$\uparrow$ hidden

$$+ h'(a_j) \sum_{k,k'} \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} w_{kj} \underbrace{\frac{\partial \overbrace{a_{k'}}^{\sum_\ell w_{k'\ell} h(a_\ell)}}{\partial a_i}}_{} \text{\textcircled{=}}$$

$k,k'$ $\uparrow$ downstream of $j$ $\quad\uparrow$ downstream of $i$

$$\sum_\ell h'(a_\ell)\, \delta_{\ell i}\, w_{k'\ell} =$$
$$= h'(a_i)\, w_{k'i}$$

$$\text{\textcircled{=}}\ h''(a_i)\, \delta_{ij} \sum_k \underbrace{\frac{\partial E_n}{\partial a_k}}_{\delta_k} w_{kj} + h'(a_i) h'(a_j) \sum_{k,k'} \underbrace{\frac{\partial^2 E_n}{\partial a_k \partial a_{k'}}}_{M_{kk'}} w_{kj} w_{k'i}$$

$$(**)$$

Compute $M_{ij}$ recursively starting from output nodes. For ex., for regression we get:

output nodes $\begin{cases} \dfrac{\partial E_n}{\partial a_j} = \dfrac{\partial E_n}{\partial y_j} = y_j - t_j, \\[2mm] \dfrac{\partial^2 E_n}{\partial a_i \partial a_j} = \delta_{ij}. \end{cases}$

$\uparrow$ $a_i = y_i$ here

Then we can compute $M_{ij}$ for hidden nodes using $(**)$.

$-9-$

Finally, in the "mixed" case:

$$\frac{\partial}{\partial a_i} \frac{\partial E_n}{\partial y_j} = \frac{\partial}{\partial a_i}(y_j - t_j) = \sum_k \underbrace{\frac{\partial(y_j - t_j)}{\partial a_k}}_{y_k} \underbrace{\frac{\partial a_k}{\partial a_i}}_{w_{ki} h'(a_i)} \overset{(=)}{}$$

↑ hidden    ↑ output                ↑ output

$$\overset{(=)}{} \sum_k \delta_{jk} \, w_{ki} \, h'(a_i) = h'(a_i) \, w_{ji}.$$

Thus we can compute all $M_{ij}$ & therefore all $H_{ij}$. This is an $O(W^2)$ computation.

———— o ————

This can be checked against finite differences:

$$\frac{\partial^2 E}{\partial w_{ji} \, \partial w_{\ell k}} = \frac{1}{4\xi^2}\Big[ E(w_{ji}+\xi, \, w_{\ell k}+\xi) - $$

$$- E(w_{ji}+\xi, \, w_{\ell k}-\xi) - E(w_{ji}-\xi, \, w_{\ell k}+\xi) + $$

$$+ E(w_{ji}-\xi, \, w_{\ell k}-\xi) \Big] + O(\xi^2)$$

↑ 4 forward propagations with

$O(W)$ operations each $\times$ $W^2$ Hessian elements ⟹

⟹ $O(W^3)$ computation.

⌐ However, we can use a mixed approach:

$$\frac{\partial^2 E}{\partial w_{ji} \, \partial w_{\ell k}} = \frac{1}{2\xi}\Big[ \frac{\partial E}{\partial w_{ji}}\Big|_{w_{\ell k}+\xi} - \frac{\partial E}{\partial w_{ji}}\Big|_{w_{\ell k}-\xi} \Big] + O(\xi^2).$$

compute using backpropagation

in $O(W)$ steps $\times$ $W$ weights to be perturbed $(w_{\ell k} \pm \xi)$:

$O(W^2)$ computation just as with explicit backpropagation above.