

MATH 565: Lecture 27 (04/21/2026)

Today: * back-propagation in neural networks

Recall In CGs,

$$z(i, j) = \frac{\partial y(j)}{\partial y(i)} \quad \text{edge derivative}$$

$$S(i, t) = \sum_{j \in A(i)} S(j, t) \cdot z(i, j) \quad \text{sensitivity recursion}$$

We now consider the case of feed forward neural networks or NNs in more detail. These are DACGs. But, unlike in generic CGs, each node in an NN typically applies an activation function Φ to a weighted sum $s(j)$ of its inputs.

For a single edge (i, j) , we have

$$y(j) = \Phi(w_{ij} y(i)).$$

More generally,

$$y(j) = \Phi \left(\sum_{i \in A(i)} w_{ij} y(i) \right) = \Phi(s(j))$$

\rightarrow tail nodes at end of in-arcs of node j

\downarrow weighted sum

$s(j)$: pre-activation at node j

LOAML uses $a(j)$, but that could be confusing, since one may confuse that with the post-activation function. Some other books use $z(j)$.

Edge derivatives

We use the chain rule through pre-activation:

$$z(i,j) = \frac{\partial y(j)}{\partial y(i)} = \frac{\partial \Phi(s_j)}{\partial s_j} \cdot \frac{\partial s_j}{\partial y(i)}$$

$$\Rightarrow z(i,j) = \Phi'(s_j) \cdot w_{ij} \quad \text{--- (1)}$$

In words, edge derivative of $(i,j) \in E$ in the NN is the product of the derivative of activation at pre-activation of node j and the synaptic weight w_{ij} .

Sensitivity Recursion

$$S(i,t) = \frac{\partial y(t)}{\partial y(i)} \equiv \text{sensitivity}$$

→ how sensitive is the output function to a unit increase in $y(i)$ value.

The recursion works just as before:

$$S(i,t) = \sum_{j \in A(i)} S(j,t) \cdot z(i,j)$$

$$\Rightarrow S(i,t) = \sum_{j \in A(i)} S(j,t) \cdot \Phi'(s_j) \cdot w_{ij} \quad \text{--- (2)}$$

with initialization (or boundary condition) $S(t,t) = 1$.

w_{ij} here "carries the signal" backward through the weight.

$\Phi'(s_j)$: "gates" how much gradient passes through node j .

⇒ When $\Phi' \approx 0$, the gate is "closed", and gives rise to the problem of vanishing gradients.

Derivatives of common activation functions

Common activation functions used in NNs: (as seen in Lecture 24)

$\Phi(v) = \frac{1}{1+e^{-v}}$ sigmoid function

$\Phi(v) = \frac{e^{2v}-1}{e^{2v}+1}$ tanh function

hyperbolic tangent

$\Phi(v) = \max\{v, 0\}$ ReLU (Rectified linear unit)

$\Phi(v) = \max\{\min\{v, 1\}, -1\}$ hard tanh function

We compute the derivative in each case.

1. Sigmoid: $\Phi(s) = \frac{1}{1+e^{-s}}$

$\Rightarrow \frac{\partial \Phi(s)}{\partial s} = \frac{-1 \cdot -1 \cdot e^{-s}}{(1+e^{-s})^2} = \frac{1}{1+e^{-s}} \left(1 - \frac{1}{1+e^{-s}}\right) = \Phi(s)(1-\Phi(s))$

2. tanh: $\Phi(s) = \frac{e^{2s}-1}{e^{2s}+1}$

$\Rightarrow \frac{\partial \Phi}{\partial s} = \frac{(e^{2s}+1)2 \cdot e^{2s} - (e^{2s}-1)2e^{2s}}{(e^{2s}+1)^2} = \frac{4e^{2s}}{(e^{2s}+1)^2}$
 $= 1 - \left(\frac{e^{2s}-1}{e^{2s}+1}\right)^2 = 1 - (\Phi(s))^2$

3. ReLU $\Phi(s) = \max\{s, 0\}$.

$\Rightarrow \Phi'(s) = \begin{cases} 1 & \text{if } s > 0, \\ 0 & \text{o.w.} \end{cases}$

4. hard tanh $\Phi(s) = \max\{\min\{s, 1\}, -1\}$

$\Rightarrow \Phi'(s) = \begin{cases} 1 & \text{if } s \in [-1, 1] \\ 0 & \text{o.w.} \end{cases}$

or just as a constant!

We can express Φ' conveniently in terms of Φ itself in each case!

Proof of backpropagation recursion

Lemma 14 Let $G_c = (N, E_c)$ be a DAG with a single output node t , and let $S(i, t) = \frac{\partial y(t)}{\partial y(i)}$ $\forall i \in N$.

Then S satisfies the recursion

$$S(i, t) = \sum_{j \in A(i)} S(j, t) \cdot z(i, j), \text{ with } S(t, t) = 1$$

where $z(i, j) = \frac{\partial y(j)}{\partial y(i)}$.

Proof We proceed by induction on $d(i) =$ length of longest directed path from i to t . → in terms of # arcs

Base case $d(i) = 0 \Rightarrow i = t$. $S(t, t) = 1$. ✓

Induction Assume result holds for $\forall i' \in N$ with $d(i') < d(i)$.

As G_c is a DAG, every directed path from i to t passes through exactly one $j \in A(i)$. By multivariate chain rule,

$$\begin{aligned} \frac{\partial y(t)}{\partial y(i)} &= \sum_{j \in A(i)} \frac{\partial y(t)}{\partial y(j)} \cdot \frac{\partial y(j)}{\partial y(i)} \\ &= \sum_{j \in A(i)} S(j, t) \cdot z(i, j) \end{aligned}$$

as each $j \in A(i)$ has $d(j) < d(i)$ and hence the induction hypothesis applies. □

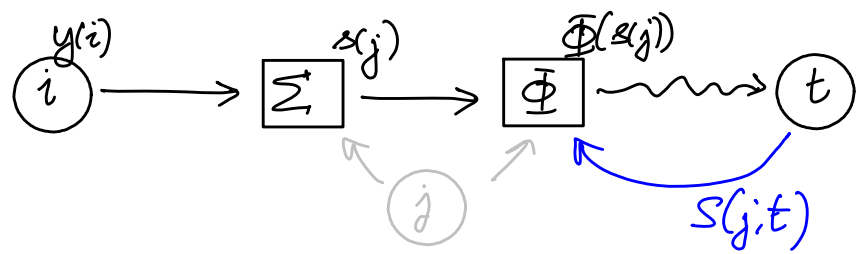
Gradient w.r.t. edge weights

Assume G_c has a single output node t , and $y(t) = J$, the loss function. Then

$$\frac{\partial y(t)}{\partial w_{ij}} = \frac{\partial y(t)}{\partial y(j)} \cdot \frac{\partial y(j)}{\partial w_{ij}} = S(j,t) \cdot \Phi'(z(j)) \cdot y(i)$$

This gradient breaks into the product of three terms

- * downstream sensitivity, $S(j,t)$
- * local gate $\Phi'(z(j))$
- * upstream (incoming) activation $y(i)$



If $\Phi'(z(j)) \approx 0$, the gate is closed. Each such edge multiplies $S(i,t)$ by a very small number. Hence, over many layers, $S(i,t)$ becomes really small.

⇒ NN cannot learn!

One good reason to use ReLU, as $\Phi'(s) = 1$ for $s > 0$!

Note: $y(i)$ is available from forward pass, while $S(j,t)$ is computed using recursion running backwards.

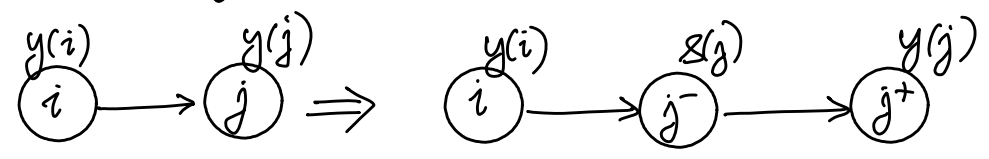
⇒ Forward and backward passes together give all needed quantities for gradient computation.

Decoupled Formulation

As in previous figure, we split each node j into two subnodes:

* linear subnode j^- : computes weighted sum
$$s(j) = \sum_{i \in AI(j)} w_{ij} y(i)$$

* activation subnode j^+ : computes output
$$y(j) = \Phi(s(j))$$



So, each $(i, j) \in E$ is replaced by a 2-edge path through sum and activation subnodes.

Edge Derivative in Decoupled NN

$$z(i, j^-) = \frac{\partial s(j)}{\partial y(i)} = w_{ij} \rightarrow \text{just weight; no activation}$$

into the linear subnode

$$z(j^-, j^+) = \frac{\partial y(j)}{\partial s(j)} = \Phi'(s(j)) \rightarrow \text{just activation, no weight}$$

through the activation subnode.

Sensitivity Recursion in Decoupled NN

Recall, $S(i, t) = \sum_{j \in A(i)} S(j, t) \cdot z(i, j)$. This recursion

applies to j^- and j^+ separately here.

$$S(j^-, t) = S(j^+, t) \cdot \Phi'(z(j))$$

$$\begin{aligned} S(i, t) &= \sum_{j \in A(i)} S(j^-, t) \cdot w_{ij} \\ &= \sum_{j \in A(i)} S(j^+, t) \cdot \Phi'(z(j)) \cdot w_{ij} \end{aligned}$$

Gradient w.r.t. w_{ij}

Note that w_{ij} appears only in (i, j^-) .

$$\begin{aligned} \Rightarrow \frac{\partial y(t)}{\partial w_{ij}} &= S(j^-, t) \cdot y(i) \\ &= S(j^+, t) \cdot \Phi'(z(j)) \cdot y(i) \end{aligned}$$

These formulae can easily generalize to the vector case in layerwise NNs.