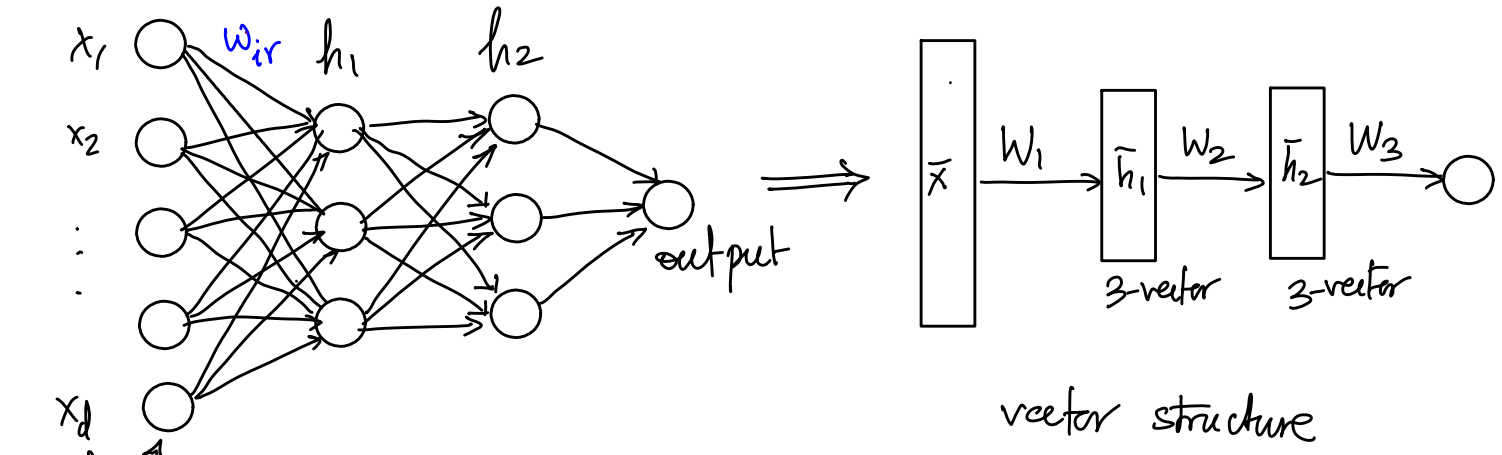


# MATH 565: Lecture 28 (04/23/2026)

Today: \* vector case of NNs

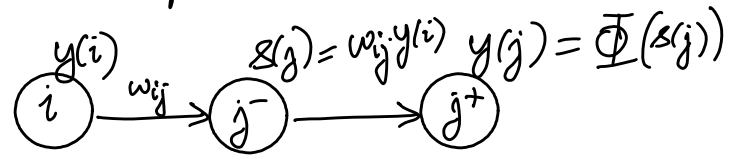
## Feedforward NNs as Vector Structures

Recall the previous illustration:



All nodes  $j$  in a layer  $l$  share the same activation function  $\Phi$ , receive inputs from all nodes in layer  $(l-1)$ , and pass output to each node in layer  $(l+1)$ .

Recall the decoupled formulation:



Generalizing to vector settings, we get

$$i \rightarrow j^- : \bar{s}^{(l)} = W \bar{y}^{(l-1)}$$

↘ vector of weighted sums

$j^- \rightarrow j^+ : \text{elementwise activation}$

$$\bar{y}^{(l)} = \Phi(\bar{s}^{(l)})$$

↗ diagonal matrix

Note: In some settings,  $\bar{y}^{(l)}$  is taken as  $\text{diag}(\Phi(\bar{s}))$  so that it has same dimensions as associated Jacobians. But we'll stick with the vector form.

## Sensitivity Recursion

$$\bar{\delta}(j^-, t) = \bar{\delta}(j^+, t) \odot \bar{\Phi}'(\bar{z})$$

↘ Hadamard or elementwise product

It is standard notation to use  $\bar{\delta}^-$  to denote  $\bar{\delta}(j^-, t)$  or  $\bar{\delta}^{-(l)}$  for the sensitivities at level  $l$ . Hence, the recursion is given as

$$\bar{\delta}^{-(l)} = \bar{\delta}^{+(l)} \odot \bar{\Phi}'(\bar{z}^{(l)}).$$

Plugging in the recursive expression for  $\bar{\delta}^{+(l)}$  gives us the backward pass equation as

$$\bar{\delta}_{\text{prev}} = W^T \bar{\delta}^- \quad \text{or, with levels,}$$

$$\bar{\delta}^{+(l-1)} = W^{(l)T} \bar{\delta}^{-(l)}.$$

We next formally prove these recursive relations for vector backpropagation.

**Theorem 15** (Vector backpropagation) Let  $G_c = (N, \mathbb{R})$  be an  $N$ -layer feedforward NN with  $\bar{y}^{(0)} = \bar{x}$  and  $\bar{z}^{(l)} = W^{(l)} \bar{y}^{(l-1)}$ ,  $\bar{y}^{(l)} = \Phi(\bar{z}^{(l)})$ , with  $\Phi$  applied elementwise, for  $l=1, \dots, N$ , and loss function  $J(\bar{y}^{(N)})$ . Let  $\bar{\delta}^{-(l)} = \frac{\partial J}{\partial \bar{z}^{(l)}}$ ,  $\bar{\delta}^{+(l)} = \frac{\partial J}{\partial \bar{y}^{(l)}}$  be the

pre- and post-activation sensitivity vectors at layer  $l$ . Then the following relations hold for  $l=1, \dots, N$ :

$$\bar{\delta}^{-(l)} = \bar{\delta}^{+(l)} \odot \Phi'(\bar{z}^{(l)}) \quad \text{--- (1)}$$

$$\bar{\delta}^{+(l)} = W^{(l+1)T} \bar{\delta}^{-(l+1)} \quad \text{--- (2)}$$

$$\frac{\partial J}{\partial W^{(l)}} = \bar{\delta}^{-(l)} \bar{y}^{(l-1)T} \quad \text{--- (3)}$$

with initialization  $\bar{\delta}^{+(N)} = \nabla_{\bar{y}^{(N)}} J$ .

Proof We apply chain rule backwards from layer  $N$ .

At the output layer,

$$\bar{\delta}^{+(N)} = \nabla_{\bar{y}^{(N)}} J \quad \text{by definition.}$$

$$\Rightarrow \bar{\delta}^{-(N)} = \frac{\partial J}{\partial \bar{z}^{(N)}} = \frac{\partial J}{\partial \bar{y}^{(N)}} \odot \frac{\partial \bar{y}^{(N)}}{\partial \bar{z}^{(N)}}$$

$$= \bar{\delta}^{+(N)} \odot \Phi'(\bar{z}^{(N)}),$$

giving (1) for  $l=N$ .

## Induction

Assume  $\bar{s}^{-(l+1)}$  is known for some  $l$  in  $1 \leq l < N$ . As  $J$  depends on  $\bar{y}^{(l)}$  only through  $\bar{s}^{(l+1)} = W^{(l+1)} \bar{y}^{(l)}$ , chain rule gives

$$\begin{aligned} \bar{\delta}^{+(l)} &= \frac{\partial J}{\partial \bar{y}^{(l)}} = \left( \frac{\partial \bar{s}^{(l+1)}}{\partial \bar{y}^{(l)}} \right)^T \frac{\partial J}{\partial \bar{s}^{(l+1)}} \xrightarrow{W^{(l+1)}} \\ &= W^{(l+1)T} \bar{\delta}^{-(l+1)}, \text{ giving (2).} \end{aligned}$$

Applying entrywise activation derivative then gives

$$\bar{\delta}^{-(l)} = \frac{\partial J}{\partial \bar{s}^{(l)}} = \frac{\partial J}{\partial \bar{y}^{(l)}} \odot \Phi'(\bar{s}^{(l)})$$

$$\Rightarrow \bar{\delta}^{-(l)} = \bar{\delta}^{+(l)} \odot \Phi'(\bar{s}^{(l)}),$$

giving (1) at layer  $l$ .

Since  $J$  depends on  $W^{(l)}$  only through  $\bar{s}^{(l)} = W^{(l)} \bar{y}^{(l-1)}$ , differentiating entrywise gives

$$\frac{\partial J}{\partial W_{jk}^{(l)}} = \sum_i \frac{\partial J}{\partial s_i^{(l)}} \cdot \frac{\partial s_i^{(l)}}{\partial W_{jk}^{(l)}} = \bar{\delta}_j^{-(l)} \cdot y_k^{(l-1)}, \text{ since}$$

$$\frac{\partial s_i^{(l)}}{\partial W_{jk}^{(l)}} = y_k^{(l-1)} \cdot \delta_{ij} \quad \text{Assembling over all } (j,k) \text{ entries}$$

gives  $\bar{\delta}_j^{-(l)} \cdot \bar{y}^{(l-1)T}$ , proving (3). □

We next present a result that is directly implied by the above recursions and characterizes some limitations on trainability of NNs just based on their architectures.

**Theorem 16** (Exponential decay of sensitivities) Suppose  $\|\Phi'\|_\infty \leq r$  and  $\|W^{(l)}\|_2 \leq \sigma$  for all layers  $l=1, \dots, N$ .

Then  $\|\bar{\delta}^{-(l)}\|_2 \leq (\sigma r)^{N-l} \|\bar{\delta}^{-(N)}\|_2$

Spectral norm:  $\|W\|_2 = \sup_{\|\bar{x}\|_2=1} \|W\bar{x}\|_2$  is the largest singular value of  $W$ . We note that

$\|AB\|_2 \leq \|A\|_2 \|B\|_2$ , called the submultiplicativity of spectral norm.

Proof (1) and (2)  $\Rightarrow$

$$\bar{\delta}^{-(l)} = (W^{(l+1)T} \bar{\delta}^{-(l+1)}) \odot \Phi'(\bar{s}^{(l)})$$

$$\Rightarrow \|\bar{\delta}^{-(l)}\|_2 \leq \|W^{(l+1)T} \bar{\delta}^{-(l+1)}\|_2 \cdot \|\Phi'(\bar{s}^{(l)})\|_\infty \leq r$$

using the standard Hadamard bound inequality:

$$\|\bar{u} \odot \bar{v}\|_2 \leq \|\bar{u}\|_2 \|\bar{v}\|_\infty$$

$$\leq \|W^{(l+1)}\|_2 \|\bar{\delta}^{-(l+1)}\|_2 \cdot r$$

$\hookrightarrow$  by submultiplicativity of spectral norm

$$\Rightarrow \|\bar{\delta}^{-(l)}\|_2 \leq \sigma r \|\bar{\delta}^{-(l+1)}\|_2$$

Applying this relation (or contraction) inductively gives

$$\|\bar{\delta}^{-(l)}\|_2 \leq (\sigma r)^{N-l} \|\bar{\delta}^{-(N)}\|_2$$



So, if  $\sigma r < 1$ , then  $\|\bar{s}^{-(l)}\|_2$  decays exponentially as we go from layer  $N$  toward the input layer ( $l=1$ ).

Hence, (3)  $\Rightarrow \frac{\partial J}{\partial W^{(l)}}$  decays exponentially for early layers, rendering it hard or impossible to train the NN with GD.

Theorem 16 suggests the following design choices for NNs.

\* ReLU activation:  $\Phi'(s) = 1 (s > 0) \Rightarrow r = 1$ . So  $\Phi$  does not contribute to the contraction of sensitivities. In contrast, sigmoid activation has  $\Phi' \leq \frac{1}{4} \Rightarrow r \leq \frac{1}{4} \Rightarrow$  decay factor is  $(\frac{1}{4})^{N-l}$

\* Choosing  $W^{(l)}$  such that  $\|W^{(l)}\|_2 = 1$  keeps  $\sigma = 1$  (so, no contraction or amplification). Combined with ReLU activation, we get  $\sigma r = 1$ , and the sensitivities are not forced to decay just by the NN architecture alone.

\* Residual NNs with skip connections could bypass this issue.

