

ECE/CS 559 Lecture 9 &amp; 10 T,Th 9/24, 9/26

Last time: Gradient Descent, Widrow-Hoff LMS Algorithm

$$\text{minimize } R(w) = \sum_{(x,y) \in \text{Data}} l(y, f(x; w))$$

$$w' \leftarrow w - \eta \nabla R(w)$$

$$\nabla R_w = \left[ \frac{\partial R}{\partial w_1}, \dots, \frac{\partial R}{\partial w_k} \right]$$

$$\cdot l(y, f(x; w)) = \|y - f(x; w)\|^2$$

$$f(x; w) = w^T x$$

$$W' \leftarrow W - \eta \nabla R(w) = W + 2\eta \sum_{(x,y)} (y - w^T x) x^T$$

- We will apply the chain rule multiple times:

$$\frac{\partial R}{\partial w_k} = \sum_{(x,y)} \left[ \frac{\partial l}{\partial w_k} \right]^{(*)}$$

need this

$$(?) \frac{\partial l}{\partial w_k} = \frac{\partial}{\partial w_k} l(t_k, \dots) \Big|_{t_k = g(w_k z + \dots)} \quad \begin{matrix} v_{t_k} \\ \vdots \\ \end{matrix}$$

have this from forward pass

$$(\text{two chain rules}) \quad = \left[ \frac{\partial e}{\partial t_k} \Big|_{t_k = g(v_{t_k})} \right]^{(*)} \cdot g'(v_{t_k}) \cdot \delta_{t_k}$$

need this at every neuron output

$$(*) \frac{\partial l}{\partial z} = \frac{\partial}{\partial z} l(t_i, \dots, t_k, \dots, t_m) \Big|_{t_k = g(w_k z + \dots)} \quad \begin{matrix} v_{t_k} \\ \vdots \\ \end{matrix}$$

have this from forward pass

$$(\text{two chain rules}) \quad = \sum_{k=1}^m \frac{\partial l}{\partial t_k} \Big|_{t_k = g(v_{t_k})} \cdot g'(v_{t_k}) \cdot w_k \quad \delta_{t_k}$$

- Of course, the layer of  $z$  has other (say  $n$ ) neurons. Let's index the equations for these (blue  $j$ )

$$\frac{\partial l}{\partial w_{kj}} = \delta_{t_k} \cdot \delta_j$$

$$\frac{\partial l}{\partial \delta_j} = \sum_{k=1}^m \delta_{t_k} w_{kj} \quad j=1, \dots, n$$

- If we think of  $\delta_z$  as a  $m$ -dim vector or  $w$  as a  $m \times n$ -dim matrix
- We get the linear algebraic backward equation:

$$\nabla_l = \delta_{t_k} \delta_z^T$$

outer product

$$\nabla_l = W^T \delta_t$$

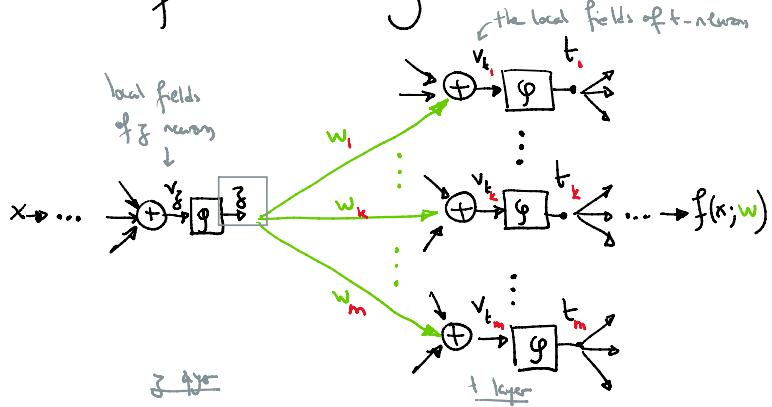
matrix multiplication

$$\delta_z = \nabla_l \circ g'(v_z)$$

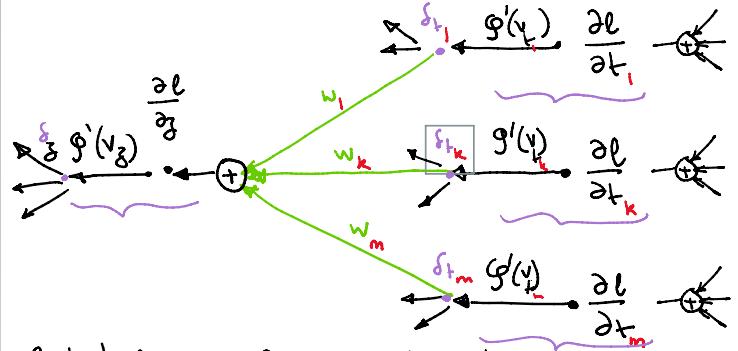
elementwise product

## Back propagation

- So far we've been calculating gradients with 1 neuron.
- What if we connect many neurons?



- To calculate (?) and (\*), we perform a forward pass to get all the outputs ( $v_{t_k}, v_t, t, \dots$ ), then we perform a backward pass:

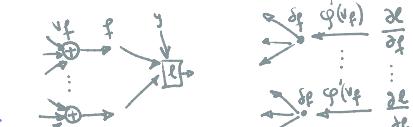

 Gradient of  $w = \text{forward signal} \times \text{backward signal} \delta_t$ 

- To complete this, we need 2 things: how to start, how to handle bias

- Start at the last layer  $\nabla_l$  depends on the choice of  $l$

 e.g. if  $l = \|f-y\|^2$  then

$$\nabla_l = 2(f-y), \delta_f = \nabla_l \circ g'(v_f)$$



- Biases are like dead-ends on the backward path:

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial b} \Big|_{t_k = g(v_k + b)} \quad \begin{matrix} v_k \\ \vdots \\ \end{matrix}$$

$\delta_t$

$$= \frac{\partial l}{\partial t_k} \Big|_{t_k = g(v_k)} \cdot g'(v_k) \cdot 1 \quad \delta_{t_k}$$

$$\nabla_l = \delta_t$$

$$\nabla_l = \delta_t \circ g'(v_t)$$

 Note: Can we find eqs to get these:  $l(t) = l(\underbrace{g(w_j+b)}_t)$   $\Rightarrow \nabla_l = [\nabla_l \circ g'(v_t)] \delta_t$ ,  $\nabla_b l = \frac{[\nabla_l \circ g'(v_t)] - 1}{\delta_t}$ ,  $\nabla_w l = W^T [\nabla_l \circ g'(v_t)] \delta_t$