# CME 216, ME 343 - Spring 2020 Eric Darve, ICME



In the previous lecture, we explain the first step in the backpropagation algorithm.

We will derive the general formula for all layers.

# Let us find the derivative with respect to $W^{\left(n-1 ight)}$ to get started.

We have:

$$y=W^{(n)}\phi\odot W^{(n-1)}a^{(n-2)}$$

## Let's differentiate with respect to the (i, j) component of $W^{(n-1)}$

Let's denote  $E_{ij}$  a matrix full of zeros with a 1 at index (i, j).

## Then

$$rac{\partial y}{\partial [W^{(n-1)}]_{ij}} = W^{(n)}\,\psi^{(n-1)}\,E_{ij}\,a^{(n-2)}$$

 $\psi^{(n-1)}$  is a diagonal matrix with entries

$$[\psi^{(n-1)}]_{ii} = [\phi' \odot W^{(n-1)} a^{(n-2)}]_i$$

The scalar  $W^{(n)} \psi^{(n-1)} E_{ij} a^{(n-2)}$  can be re-written using vectors:

$$rac{\partial y}{\partial [W^{(n-1)}]_{ij}} = [W^{(n)}\psi^{(n-1)}]_i \; [a^{(n-2)}]_j$$

 $W^{(n)}\psi^{(n-1)}$  is a row vector

 $a^{(n-2)}$  is a column vector

## Let's denote

$$\delta^{(n)}=\psi^{(n-1)}[W^{(n)}]^T$$

7/22

## Using matrix notations we get

$$rac{\partial y}{\partial W^{(n-1)}} = \delta^{(n)} \ [a^{(n-2)}]^T$$

 $\delta^{(n)}$  is a column vector $[a^{(n-2)}]^T$  is a row vector. $rac{\partial y}{\partial W^{(n-1)}}$  is a matrix

Let us briefly repeat this for  $W^{(n-2)}$  to see how this works. We now have:

$$y=W^{(n)}\phi\odot W^{(n-1)}\phi\odot W^{(n-2)}a^{(n-3)}$$

 $rac{\partial y}{\partial [W^{(n-2)}]_{ij}} =$  $= [W^{(n)}\psi^{(n-1)}W^{(n-1)}\psi^{(n-2)}]_i \, [a^{(n-3)}]_j$ 

## In matrix form:

$$egin{aligned} &rac{\partial y}{\partial [W^{(n-2)}]} = \delta^{(n-1)} \, [a^{(n-3)}]^T \ &\delta^{(n-1)} = \psi^{(n-2)} [W^{(n-1)}]^T \psi^{(n-1)} [W^{(n)}]^T \end{aligned}$$

T

We can now define the general formula.

There is a recurrence relation for  $\delta^{(k)}$ .

$$egin{aligned} \delta^{(k)} &= \psi^{(k-1)} \, [W^{(k)}]^T \, \delta^{(k+1)}, & \delta^{(n+1)} \ & \ & [\psi^{(k-1)}]_{ii} = [\phi' \odot W^{(k-1)} a^{(k-2)}]_i \end{aligned}$$

12/22

= 1

# Equation for gradient with respect to matrix $W^{\left(k ight)}$

$$rac{\partial y}{\partial W^{(k)}} = \delta^{(k+1)} \, [a^{(k-1)}]^T$$

# Explicit expression for the bias

We previously said that the bias can be represented in this framework by adding a 1 at the end of the activation  $a^{(k)}$ .

With this trick, we can immediately find an explicit expression for the gradient with respect to the bias.

## Gradient with respect to the bias

$$rac{\partial y}{\partial b^{(k)}} = \delta^{(k+1)}$$

The implementation can be done in two passes called the forward and backward passes.

# In the following, for completeness we will reintroduce the biases.

# Forward pass

We compute all the activations.

k goes from 1 to n-1.

$$a^{(k)} = \phi \odot (W^{(k)} a^{(k-1)} + b^{(k)}) \ a^{(0)} = x$$

 $\phi \odot$  means that  $\phi$  is applied element-wise to  $W^{(k)}a^{(k-1)} + b^{(k)}.$ 

# We also need to save the values of the derivative: $[\psi^{(k)}]_{ii} = [\phi' \odot (W^{(k)}a^{(k-1)} + b^{(k)})]_i$ $\psi^{(k)}$ is a diagonal matrix.

20/22

# Backward pass

For k going from n to 2:

$$egin{aligned} \delta^{(k)} &= \psi^{(k-1)} \, [W^{(k)}]^T \, \delta^{(k+1)} \ \delta^{(n+1)} &= 1 \end{aligned}$$

21/22

Derivatives: from k = n to k = 1

$$egin{aligned} rac{\partial y}{\partial W^{(k)}} &= \delta^{(k+1)} \, [a^{(k-1)}]^T \ & \ rac{\partial y}{\partial b^{(k)}} &= \delta^{(k+1)} \end{aligned}$$

22/22