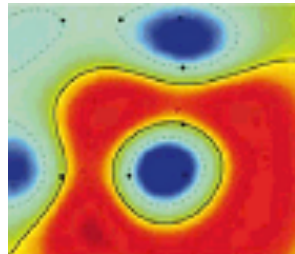
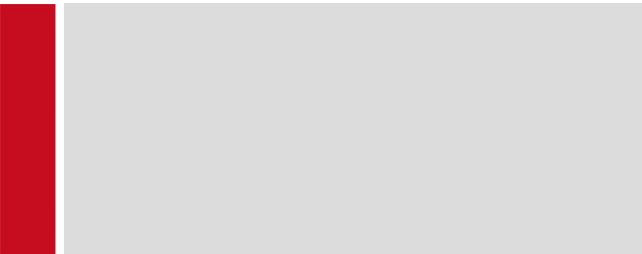




WiSe 2022/23

Deep Learning 1



Lecture 2

Error Backpropagation

- ▶ The Perceptron algorithm
- ▶ Perceptron as gradient descent
- ▶ How to compute the gradient in a neural network
- ▶ Numerical gradient computation
- ▶ The error backpropagation algorithm
 - ▶ Chain rule and multivariate chain rule
 - ▶ Worked through example
 - ▶ General equation for backpropagation
 - ▶ Vectorization of backpropagation
 - ▶ Automatic differentiation

The Perceptron

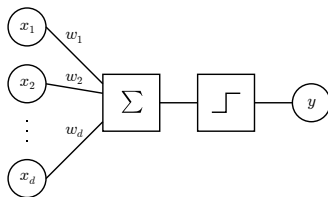


F. Rosenblatt (1928–1971)

- ▶ Proposed by F. Rosenblatt in 1958.
- ▶ Classifier that perfectly separates training data (if the data is linearly separable).
- ▶ Trained using a simple and cheap iterative procedure.
- ▶ The perceptron is an early version of artificial neural networks (with only one layer).

The Perceptron

Structure of the perceptron:



- ▶ A weighted sum of the input features:

$$\begin{aligned} z &= \sum_{i=1}^d w_i x_i + b \\ &= \mathbf{w}^\top \mathbf{x} + b \end{aligned}$$

- ▶ Followed by the sign function

$$y = \text{sign}(z)$$

Problem formulation: Let

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)} \in \mathbb{R}^d$$

be our input data and $t^{(1)}, t^{(2)}, \dots, t^{(N)} \in \{-1, 1\}$ our corresponding labels (aka. targets). The goal of the perceptron is to learn a collection of parameters (\mathbf{w}, b) such that all points are correctly classified, i.e.

$$\forall_{k=1}^N : y^{(k)} = t^{(k)}$$

The Perceptron Algorithm

Recall that for each data point, predictions of our perceptron are computed as:

$$z^{(k)} = \mathbf{w}^\top \mathbf{x}^{(k)} + b$$

$$y^{(k)} = \text{sign}(z^{(k)})$$

Perceptron algorithm

- ▶ Iterate (multiple times from $k = 1 \dots N$).
 - ▶ If $\mathbf{x}^{(k)}$ is correctly classified ($y^{(k)} = t^{(k)}$), continue.
 - ▶ If $\mathbf{x}^{(k)}$ is wrongly classified ($y^{(k)} \neq t^{(k)}$), update the perceptron:

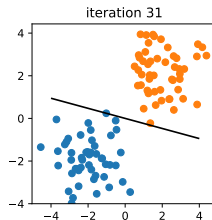
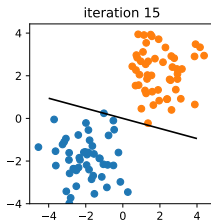
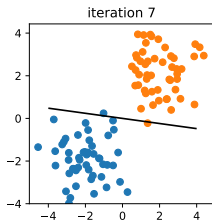
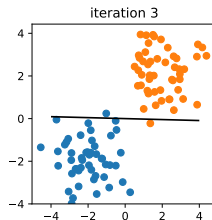
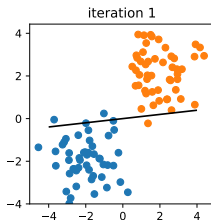
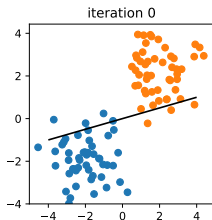
$$\mathbf{w} \leftarrow \mathbf{w} + \eta \cdot \mathbf{x}^{(k)} t^{(k)}$$

$$b \leftarrow b + \eta \cdot t^{(k)}$$

where η is a learning rate.

- ▶ Stop once all examples are correctly classified.

Perceptron at Work



The Perceptron: Optimization View

Proposition

The perceptron can be seen as a gradient descent of the error function

$$\mathcal{E}(\mathbf{w}, b) = \frac{1}{N} \sum_{k=1}^N \underbrace{\max(0, -z^{(k)} t^{(k)})}_{\mathcal{E}_k(\mathbf{w}, b)}$$

Also known as the Hinge Loss.

Proof.

$$\begin{aligned} \mathbf{w} - \eta \frac{\partial \mathcal{E}_k}{\partial \mathbf{w}} &= \mathbf{w} - \eta \cdot 1_{-z^{(k)} t^{(k)} > 0} \cdot \left(- \frac{\partial z^{(k)}}{\partial \mathbf{w}} t^{(k)} \right) \\ &= \mathbf{w} - \eta \cdot 1_{y^{(k)} \neq t^{(k)}} \cdot \left(- \frac{\partial z^{(k)}}{\partial \mathbf{w}} t^{(k)} \right) \\ &= \mathbf{w} + \eta \cdot 1_{y^{(k)} \neq t^{(k)}} \cdot \mathbf{x}^{(k)} t^{(k)} \end{aligned}$$

which is the parameter update equation of the perceptron algorithm. We proceed similarly for the parameter b . □

Perceptron vs. Nearest Mean Classifier

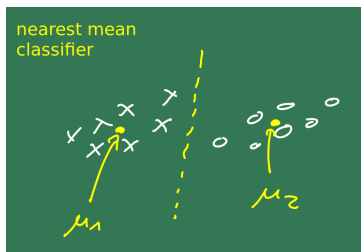
The nearest mean classifier first builds the means of two classes:

$$\mu_1 = \frac{1}{|\mathcal{C}_1|} \sum_{k \in \mathcal{C}_1} \mathbf{x}^{(k)} \quad \text{and} \quad \mu_2 = \frac{1}{|\mathcal{C}_2|} \sum_{k \in \mathcal{C}_2} \mathbf{x}^{(k)}$$

Then it predicts:

class 1 if $\|\mathbf{x} - \mu_1\| < \|\mathbf{x} - \mu_2\|$

class 2 if $\|\mathbf{x} - \mu_1\| > \|\mathbf{x} - \mu_2\|$



Perceptron vs. Nearest Mean Classifier

Equivalent formulation of the nearest mean classifier:

$$\text{class 1} \quad \text{if } \|x - \mu_1\|^2 < \|x - \mu_2\|^2$$

$$\text{class 2} \quad \text{if } \|x - \mu_1\|^2 > \|x - \mu_2\|^2$$

(raising distances to the square does not change the decision).

The nearest mean classifier can then be further developed to become expressible as a linear classifier:

$$\begin{aligned} y &= \text{sign}(\|x - \mu_2\|^2 - \|x - \mu_1\|^2) \\ &= \text{sign}(\|x\|^2 - 2\mu_2^\top x + \|\mu_2\|^2 - \|x\|^2 + 2\mu_1^\top x - \|\mu_1\|^2) \\ &= \text{sign}(\underbrace{2(\mu_1 - \mu_2)^\top}_{w} x + \underbrace{\|\mu_2\|^2 - \|\mu_1\|^2}_{b}) \end{aligned}$$

where $y = 1$ corresponds to class 1 and $y = -1$ corresponds to class 2.

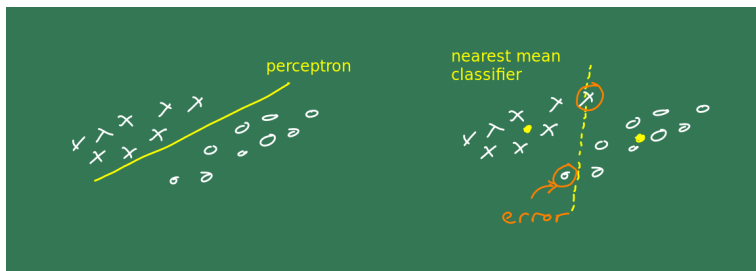
Question: How does the perceptron and nearest mean classifier differ?

Perceptron vs. Nearest Mean Classifier

Question: Both the perceptron and the nearest mean classifier are linear classifiers. Then, how do the perceptron and nearest mean classifier differ?

- ▶ Perceptron always separates data when it is linearly separable. This is not always the case for the nearest mean classifier, especially when the data of each class is elongated.

Example:



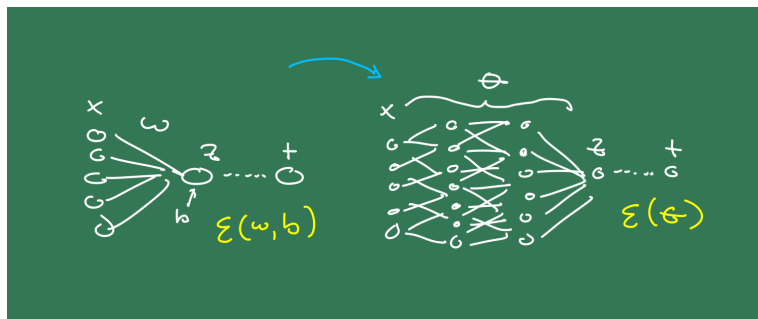
Nearest mean classifier describes the data, perceptron minimizes the error.

From the Perceptron Algorithm to Training Deep Networks

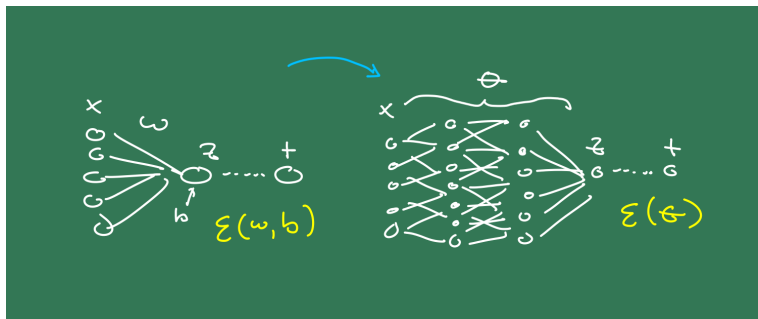
Recap: The perceptron can be seen as a gradient descent of the error function

$$\mathcal{E}(\mathbf{w}, b) = \frac{1}{N} \sum_{k=1}^N \underbrace{\max(0, -z^{(k)} t^{(k)})}_{\mathcal{E}_k(\mathbf{w}, b)}$$

Idea: Generalize the formulation where z is not the output of the perceptron, but of any multilayer neural network.



From the Perceptron Algorithm to Training Deep Networks



Question: How hard is it to compute the gradient of the newly defined error function w.r.t. the model parameters?

$$\theta = \theta - \eta \frac{\partial \mathcal{E}}{\partial \theta}$$

Numerical differentiation

Formula for numerical differentiation:

$$\forall_t : \frac{\partial \mathcal{E}}{\partial \theta_t} = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{E}(\theta + \epsilon \cdot \delta_t) - \mathcal{E}(\theta)}{\epsilon}$$

where δ_t is an indicator vector for the parameter t .

Properties:

- ▶ Can be applied to any error function \mathcal{E} (not necessary the error of a neural network).
- ▶ Need to evaluate the function as many times as there are parameters (\rightarrow slow when the number of parameters is large).
- ▶ A neural network typically has between 10^3 and 10^9 parameters (\rightarrow numerical differentiation unfeasible).
- ▶ Still useful as a unit-test for verifying gradient computation.
- ▶ Because ϵ and the numerator are very small, for numerical differentiation to work, one must use high-precision (e.g. float64 rather than float32).

The Chain Rule

Suppose that some parameter of interest θ_q (one element of the parameter vector θ) is linked to the output of the network through some sequence of functions.

$$\theta_q \rightarrow a \rightarrow b \rightarrow z$$

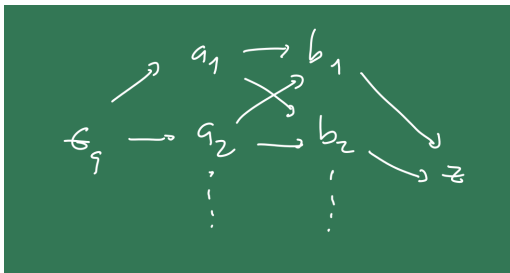
The chain rule for derivatives states that

$$\frac{\partial z}{\partial \theta_q} = \frac{\partial a}{\partial \theta_q} \frac{\partial b}{\partial a} \frac{\partial z}{\partial b}$$

i.e. the derivative w.r.t. the parameter of interest is the product of local derivatives along the path connecting θ_q to z .

The Multivariate Chain Rule

In practice, some parameter of interest may be linked to the output of the network through multiple path (formed by all neurons in layers between θ_q and z):

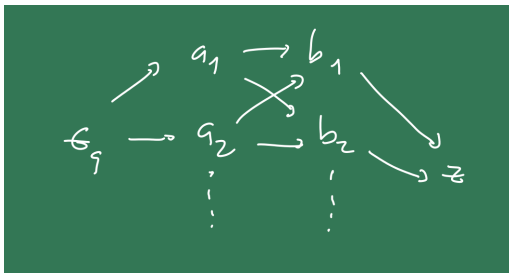


The chain rule can be extended to this multivariate scenario by enumerating all the paths between θ_q and z

$$\frac{\partial z}{\partial \theta_q} = \sum_i \sum_j \frac{\partial a_i}{\partial \theta_q} \frac{\partial b_j}{\partial a_i} \frac{\partial z}{\partial b_j}$$

where \sum_i and \sum_j run over the indices of the nodes in the corresponding layers. This is a nested sum (its complexity grows exponentially with the number of layers).

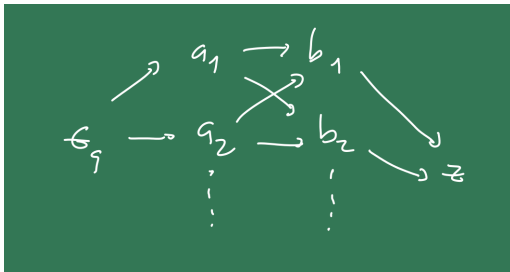
Factor Structure in the Multivariate Chain Rule



$$\frac{\partial z}{\partial \theta_q} = \sum_i \frac{\partial a_i}{\partial \theta_q} \underbrace{\sum_j \frac{\partial b_j}{\partial a_i} \frac{\partial z}{\partial b_j}}_{\delta_i}$$

- ▶ Computation can be rewritten in a way that summing operation can be performed incrementally.
- ▶ Intermediate computation can be reused for different paths, and for different parameters for which we would like to compute the gradient.

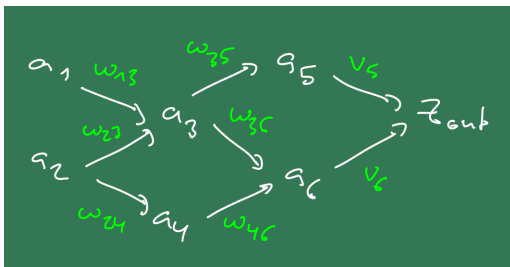
Factor Structure in the Multivariate Chain Rule



$$\frac{\partial z}{\partial \theta_q} = \sum_i \frac{\partial a_i}{\partial \theta_q} \underbrace{\sum_j \frac{\partial b_j}{\partial a_i} \frac{\partial z}{\partial b_j}}_{\delta_i}$$

- ▶ Overall, the resulting gradient computation (w.r.t. of all parameters in the network) becomes linear with the size of the network (\Rightarrow fast!)
- ▶ The algorithm is known as the Error Backpropagation algorithm (Rumelhart 1986).

Worked through Example



Forward pass:

$$a_1 = x_1$$

$$a_2 = x_2$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$a_3 = \tanh(z_3)$$

$$z_4 = a_2 w_{24}$$

$$a_4 = \tanh(z_4)$$

$$z_5 = a_3 w_{35}$$

$$a_5 = \tanh(z_5)$$

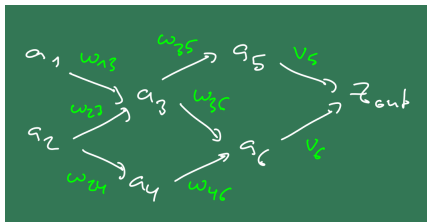
$$z_6 = a_3 w_{36} + a_4 w_{46}$$

$$a_6 = \tanh(z_6)$$

$$z_{out} = a_5 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

Worked through Example



$$a_1 = x_1$$

$$a_2 = x_2$$

$$a_3 = \tanh(z_3)$$

$$a_4 = \tanh(z_4)$$

$$a_5 = \tanh(z_5)$$

$$a_6 = \tanh(z_6)$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$z_4 = a_2 w_{24}$$

$$z_5 = a_3 w_{35}$$

$$z_6 = a_3 w_{36} + a_4 w_{46}$$

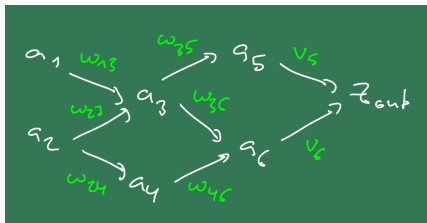
$$z_{out} = a_4 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

Backward pass:

$$\delta_{out} = \frac{\partial \mathcal{E}}{\partial z_{out}} = 1_{\{-z_{out} t > 0\}} \cdot (-t)$$

Worked through Example



$$a_1 = x_1$$

$$a_2 = x_2$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$a_3 = \tanh(z_3)$$

$$z_4 = a_2 w_{24}$$

$$a_4 = \tanh(z_4)$$

$$z_5 = a_3 w_{35}$$

$$a_5 = \tanh(z_5)$$

$$z_6 = a_3 w_{36} + a_4 w_{46}$$

$$a_6 = \tanh(z_6)$$

$$z_{out} = a_4 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

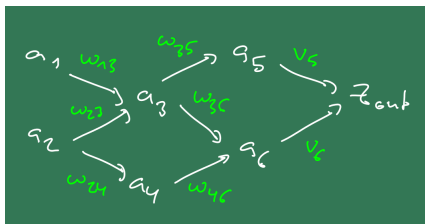
Backward pass:

$$\delta_{out} = \frac{\partial \mathcal{E}}{\partial z_{out}} = 1_{\{-z_{out} t > 0\}} \cdot (-t)$$

$$\frac{\partial \mathcal{E}}{\partial v_6} = \frac{\partial z_{out}}{\partial v_6} \frac{\partial \mathcal{E}}{\partial z_{out}} = a_6 \cdot \delta_{out}$$

$$\frac{\partial \mathcal{E}}{\partial v_5} = \frac{\partial z_{out}}{\partial v_5} \frac{\partial \mathcal{E}}{\partial z_{out}} = a_5 \cdot \delta_{out}$$

Worked through Example



$$a_1 = x_1$$

$$a_2 = x_2$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$a_3 = \tanh(z_3)$$

$$z_4 = a_2 w_{24}$$

$$a_4 = \tanh(z_4)$$

$$z_5 = a_3 w_{35}$$

$$a_5 = \tanh(z_5)$$

$$z_6 = a_3 w_{36} + a_4 w_{46}$$

$$a_6 = \tanh(z_6)$$

$$z_{out} = a_4 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

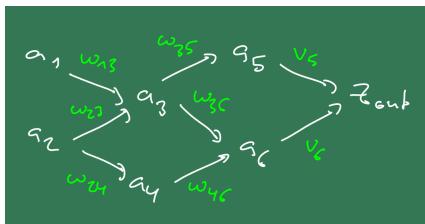
Backward pass:

$$\delta_{out} = \frac{\partial \mathcal{E}}{\partial z_{out}} = 1_{\{-z_{out} t > 0\}} \cdot (-t)$$

$$\delta_6 = \frac{\partial \mathcal{E}}{\partial a_6} = \frac{\partial z_{out}}{\partial a_6} \frac{\partial \mathcal{E}}{\partial z_{out}} = v_6 \cdot \delta_{out}$$

$$\delta_5 = \frac{\partial \mathcal{E}}{\partial a_5} = \frac{\partial z_{out}}{\partial a_5} \frac{\partial \mathcal{E}}{\partial z_{out}} = v_5 \cdot \delta_{out}$$

Worked through Example



$$a_1 = x_1$$

$$a_2 = x_2$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$z_4 = a_2 w_{24}$$

$$z_5 = a_3 w_{35}$$

$$z_6 = a_3 w_{36} + a_4 w_{46}$$

$$z_{out} = a_4 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

$$a_3 = \tanh(z_3)$$

$$a_4 = \tanh(z_4)$$

$$a_5 = \tanh(z_5)$$

$$a_6 = \tanh(z_6)$$

Backward pass:

$$\delta_6 = \frac{\partial \mathcal{E}}{\partial a_6} = \frac{\partial z_{out}}{\partial a_6} \frac{\partial \mathcal{E}}{\partial z_{out}} = v_6 \cdot \delta_{out}$$

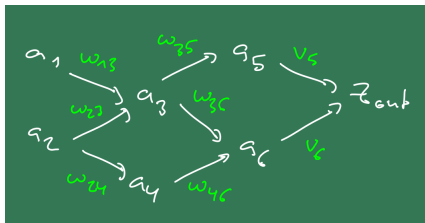
$$\delta_5 = \frac{\partial \mathcal{E}}{\partial a_5} = \frac{\partial z_{out}}{\partial a_5} \frac{\partial \mathcal{E}}{\partial z_{out}} = v_5 \cdot \delta_{out}$$

$$\frac{\partial \mathcal{E}}{\partial w_{46}} = \frac{\partial z_6}{\partial w_{46}} \frac{\partial a_6}{\partial z_6} \frac{\partial \mathcal{E}}{\partial a_6} = a_4 \cdot \tanh'(z_6) \cdot \delta_6$$

$$\frac{\partial \mathcal{E}}{\partial w_{36}} = \frac{\partial z_6}{\partial w_{36}} \frac{\partial a_6}{\partial z_6} \frac{\partial \mathcal{E}}{\partial a_6} = a_3 \cdot \tanh'(z_6) \cdot \delta_6$$

$$\frac{\partial \mathcal{E}}{\partial w_{35}} = \frac{\partial z_5}{\partial w_{35}} \frac{\partial a_5}{\partial z_5} \frac{\partial \mathcal{E}}{\partial a_5} = a_3 \cdot \tanh'(z_5) \cdot \delta_5$$

Worked through Example



$$a_1 = x_1$$

$$a_2 = x_2$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$z_4 = a_2 w_{24}$$

$$z_5 = a_3 w_{35}$$

$$z_6 = a_3 w_{36} + a_4 w_{46}$$

$$z_{out} = a_4 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

$$a_3 = \tanh(z_3)$$

$$a_4 = \tanh(z_4)$$

$$a_5 = \tanh(z_5)$$

$$a_6 = \tanh(z_6)$$

Backward pass:

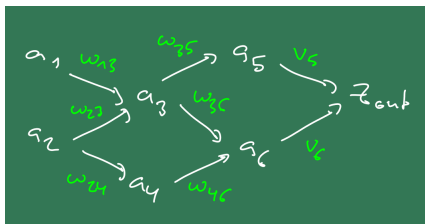
$$\delta_6 = \frac{\partial \mathcal{E}}{\partial a_6} = \frac{\partial z_{out}}{\partial a_6} \frac{\partial \mathcal{E}}{\partial z_{out}} = v_6 \cdot \delta_{out}$$

$$\delta_5 = \frac{\partial \mathcal{E}}{\partial a_5} = \frac{\partial z_{out}}{\partial a_5} \frac{\partial \mathcal{E}}{\partial z_{out}} = v_5 \cdot \delta_{out}$$

$$\delta_4 = \frac{\partial \mathcal{E}}{\partial a_4} = \frac{\partial z_6}{\partial a_4} \frac{\partial a_6}{\partial z_6} \frac{\partial \mathcal{E}}{\partial a_6} = w_{46} \cdot \tanh'(z_6) \cdot \delta_6$$

$$\delta_3 = \frac{\partial \mathcal{E}}{\partial a_3} = \frac{\partial z_6}{\partial a_3} \frac{\partial a_6}{\partial z_6} \frac{\partial \mathcal{E}}{\partial a_6} + \frac{\partial z_5}{\partial a_3} \frac{\partial a_5}{\partial z_5} \frac{\partial \mathcal{E}}{\partial a_5} = w_{36} \cdot \tanh'(z_6) \cdot \delta_6 + w_{35} \cdot \tanh'(z_5) \cdot \delta_5$$

Worked through Example



$$a_1 = x_1$$

$$a_2 = x_2$$

$$z_3 = a_1 w_{13} + a_2 w_{23}$$

$$z_4 = a_2 w_{24}$$

$$z_5 = a_3 w_{35}$$

$$z_6 = a_3 w_{36} + a_4 w_{46}$$

$$z_{out} = a_5 v_5 + a_6 v_6$$

$$\mathcal{E} = \max(0, -z_{out} t)$$

$$a_3 = \tanh(z_3)$$

$$a_4 = \tanh(z_4)$$

$$a_5 = \tanh(z_5)$$

$$a_6 = \tanh(z_6)$$

Backward pass:

$$\delta_4 = \frac{\partial \mathcal{E}}{\partial a_4} = \frac{\partial z_6}{\partial a_4} \frac{\partial a_6}{\partial z_6} \frac{\partial \mathcal{E}}{\partial a_6} = w_{46} \cdot \tanh'(z_6) \cdot \delta_6$$

$$\delta_3 = \frac{\partial \mathcal{E}}{\partial a_3} = \frac{\partial z_6}{\partial a_3} \frac{\partial a_6}{\partial z_6} \frac{\partial \mathcal{E}}{\partial a_6} + \frac{\partial z_5}{\partial a_3} \frac{\partial a_5}{\partial z_5} \frac{\partial \mathcal{E}}{\partial a_5} = w_{36} \cdot \tanh'(z_6) \cdot \delta_6 + w_{35} \cdot \tanh'(z_5) \cdot \delta_5$$

$$\frac{\partial \mathcal{E}}{\partial w_{24}} = \frac{\partial z_4}{\partial w_{24}} \frac{\partial a_4}{\partial z_4} \frac{\partial \mathcal{E}}{\partial a_4} = a_2 \cdot \tanh'(z_4) \cdot \delta_4$$

$$\frac{\partial \mathcal{E}}{\partial w_{23}} = \frac{\partial z_3}{\partial w_{23}} \frac{\partial a_3}{\partial z_3} \frac{\partial \mathcal{E}}{\partial a_3} = a_2 \cdot \tanh'(z_3) \cdot \delta_3$$

$$\frac{\partial \mathcal{E}}{\partial w_{13}} = \frac{\partial z_3}{\partial w_{13}} \frac{\partial a_3}{\partial z_3} \frac{\partial \mathcal{E}}{\partial a_3} = a_1 \cdot \tanh'(z_3) \cdot \delta_3$$

Choice of Nonlinear Activation Function

In practice, for training to proceed, the nonlinear function must be chosen in a way that:

- ▶ Its gradient is defined (almost) everywhere.
- ▶ There is a significant portion of the input domain where the gradient is non-zero.
- ▶ Gradient must be informative i.e. indicate decrease/increase of the activation function.

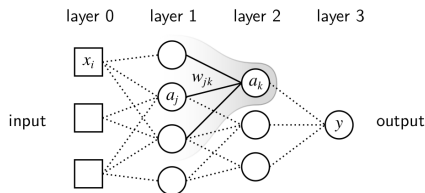
Example of commonly used activation functions:

- ▶ $g(z) = \exp(z)/(1 + \exp(z))$
- ▶ $g(z) = \tanh(z)$
- ▶ $g(z) = \max(0, z)$

Example of problematic activation functions:

- ▶ $g(z) = \max(0, z - 100)$
- ▶ $g(z) = 1_{z>0}$
- ▶ $g(z) = \sin(100 \cdot z)$

Formalization for a Standard Neural Network



The gradient of the error can be propagated from layer to layer using the chain rule:

$$\underbrace{\frac{\partial \mathcal{E}}{\partial a_j}}_{\delta_j} = \sum_k \underbrace{\frac{\partial a_k}{\partial a_j}}_{w_{jk} g'(z_k)} \cdot \underbrace{\frac{\partial \mathcal{E}}{\partial a_k}}_{\delta_k}$$

And gradients w.r.t. parameters at each layer can be extracted as:

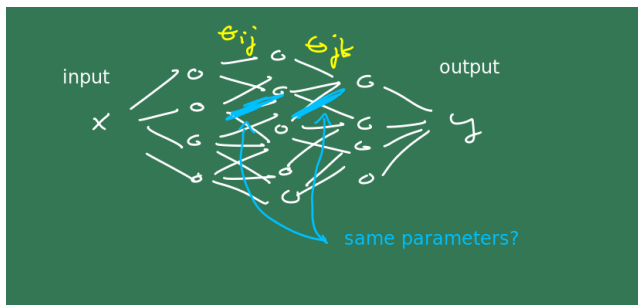
$$\frac{\partial \mathcal{E}}{\partial w_{jk}} = \sum_k \underbrace{\frac{\partial a_k}{\partial w_{jk}}}_{a_j g'(z_k)} \cdot \underbrace{\frac{\partial \mathcal{E}}{\partial a_k}}_{\delta_k}$$

Matrix Formulation

Neuron-wise	Layer-wise
$\delta_j = \sum_k w_{jk} g'(z_k) \cdot \delta_k$	$\boldsymbol{\delta}^{(l-1)} = W^{(l-1,l)} \cdot (g'(\mathbf{z}^{(l)}) \odot \boldsymbol{\delta}^{(l)})$
$\frac{\partial \mathcal{E}}{\partial w_{jk}} = \sum_k a_j g'(z_k) \delta_k$	$\frac{\partial \mathcal{E}}{\partial W^{(l-1,l)}} = \mathbf{a} \cdot (g'(\mathbf{z}^{(l)}) \odot \boldsymbol{\delta}^{(l)})^\top$

- ▶ Backpropagation equations can be written as matrix-vector products, or outer products.
- ▶ Further vectorization can be achieved by computing the gradient for multiple data points at once, in which case, we have matrix-matrix products.

Neural Networks with Shared Parameters



Automatic Differentiation:

- ▶ Generates automatically backpropagation equations from the forward equations.
- ▶ Automatic differentiation became widely available in neural network libraries (PyTorch, Tensorflow, JAX, etc.)

Consequences:

- ▶ In practice we do not need to do backpropagation anymore. We just need to program the forward pass, and the backward pass comes for free.
- ▶ This has enabled researchers to develop neural networks that are way more complex, and with much more heterogeneous structures (e.g. ResNet, Yolo, transformers, etc.).
- ▶ Only in few cases, it is still useful to express the gradient analytically (e.g. to analyze theoretically the stability of a gradient descent procedure, such as the vanishing/exploding gradients problem in recurrent neural networks).

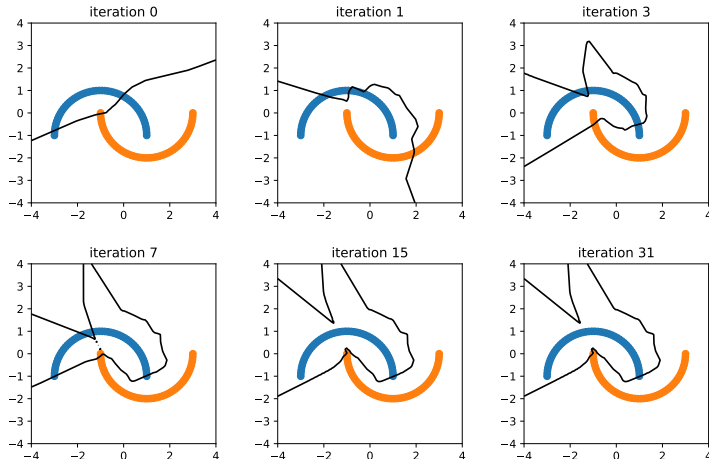
Simple algorithm for training a neural network

Basic gradient descent algorithm:

- ▶ Initialize vector of parameters θ at random.
- ▶ Iterate for T steps:
 - ▶ Compute the forward pass
 - ▶ Extract the gradient $\partial E / \partial \theta$ using backpropagation
 - ▶ Perform a gradient step $\theta = \theta - \gamma \partial E / \partial \theta$

where γ is a learning rate that needs to be set by the user.

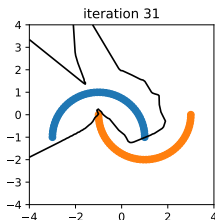
Neural Network at Work



Observation:

- ▶ After enough iterations, all points are on the correct side of the decision boundary (like for the perceptron, but even if the data is not separable).

... still not very fast, and not an optimal decision boundary.



Consideration for next lectures:

Optimization (Lecture 3)

- ▶ How to make training faster? (especially important if we consider large problems with many input variables.)

Regularization (Lecture 4)

- ▶ Decision function doesn't look nice. Unlikely to work well for new data points. Can we introduce a mechanism in the learning procedure that promotes more regular and well-generalizing decision functions?

Summary

- ▶ Error of a classifier can be minimized using gradient descent (e.g. Perceptron, neural network + backpropagation).
- ▶ Error backpropagation is computationally efficient way of computing the gradient (much faster than using the limit formulation of the derivative).
- ▶ Error backpropagation is an application of the multivariate chain rule, where the different terms can be factored due to the structure of the neural network graph.
- ▶ In practice, we most of the time do not need to program error backpropagation manually, and we can instead use automatic differentiation techniques available in most modern neural network libraries.
- ▶ Error backpropagation only extracts the gradient. This is not enough for training a neural network successfully, one still needs to make sure the gradient descent is carried out efficiently (lecture 3) and that the learned model is robust to generalize well to new data points (lecture 4).