

AOS 2 – Deep learning

Lecture 01: Neural networks: an introduction

Sylvain Rousseau

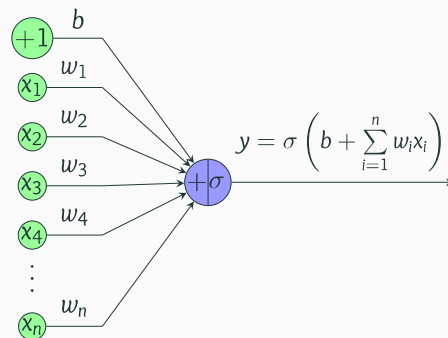
Feed forward neural network

1

What is a neuron?

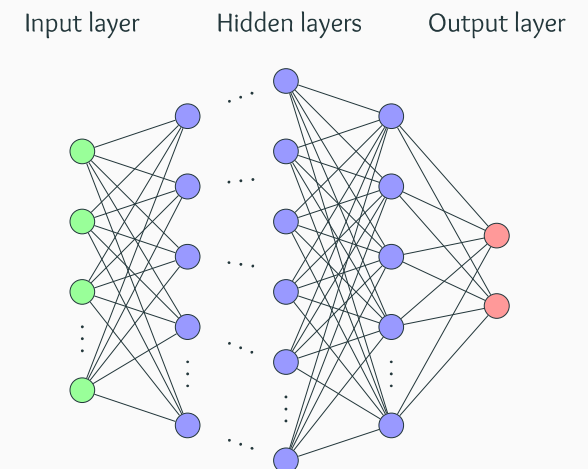
Neuron = linear transform of the input followed by a non-linearity

- Input: $\mathbf{x} = (x_i)_{i=1\dots n} \in \mathbb{R}^n$
- Scalar output: $y \in \mathbb{R}$
- +1 denotes the intercept (or bias) b
- The w_i 's are called **weights**
- σ is nonlinear function called an **activation function**



Feed forward neural network

- Stacked collection of neurons arranged in layers
- Input layer nodes are not neurons
- Output layer neurons do not have necessarily an activation function



2

3

Input layer and first hidden layer

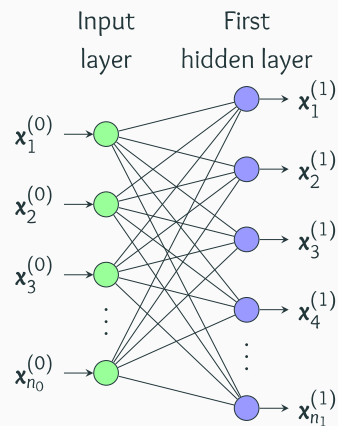
- Input: $\mathbf{x}^{(0)} = (\mathbf{x}_i^{(0)})_{i=1\dots n_0} \in \mathbb{R}^{n_0}$
- Output: $\mathbf{x}^{(1)} = (\mathbf{x}_i^{(1)})_{i=1\dots n_1} \in \mathbb{R}^{n_1}$

We have

$$\mathbf{x}_i^{(1)} = \sigma \left(\sum_{j=1}^{n_0} \mathbf{w}_{ij}^{(1)} \mathbf{x}_j^{(0)} + \mathbf{b}_i^{(1)} \right)$$

where

- $\mathbf{b}_i^{(1)}$ is the bias of neuron i (not displayed)
- $\mathbf{w}_{ij}^{(1)}$ is the j -th coefficient of i -th neuron
- Vector version $\mathbf{x}_i^{(1)} = \sigma \left(\langle \mathbf{w}_i^{(1)}, \mathbf{x}^{(0)} \rangle + \mathbf{b}_i^{(1)} \right)$

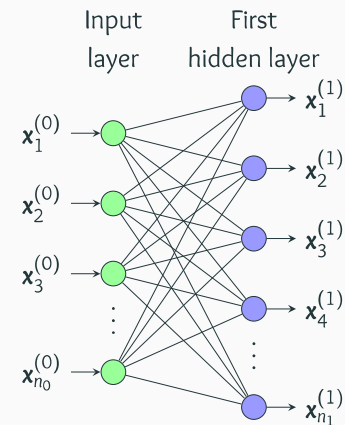


Input layer and first hidden layer

- Matrix version:

$$\mathbf{x}^{(1)} = \sigma \left(\left(\mathbf{W}^{(1)} \right)^T \mathbf{x}^{(0)} + \mathbf{b}^{(1)} \right)$$

- σ is applied element-wise
- $\mathbf{W}^{(1)}$ is of size $n_0 \times n_1$
- $\mathbf{w}_i^{(1)} \in \mathbb{R}^{n_0}$ columns of $\mathbf{W}^{(1)}$
- $\mathbf{b}^{(1)} \in \mathbb{R}^{n_1}$ groups all the biases
- The parameters are $\Theta^{(1)} = \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}\}$



4

5

Two hidden layers

- n_k is the number of neurons at layer k

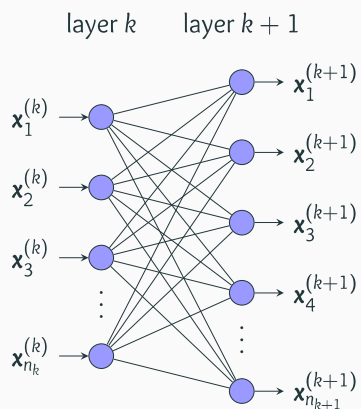
$$\mathbf{x}_i^{(k+1)} = \sigma \left(\sum_{j=1}^{n_k} \mathbf{w}_{ij}^{(k+1)} \mathbf{x}_j^{(k)} + \mathbf{b}_i^{(k+1)} \right)$$

- Condensed version

$$\mathbf{x}^{(k+1)} = \sigma \left(\left(\mathbf{W}^{(k+1)} \right)^T \mathbf{x}^{(k)} + \mathbf{b}^{(k+1)} \right)$$

- The parameters at layer $k+1$ are

$$\Theta^{(k+1)} = \{\mathbf{W}^{(k+1)}, \mathbf{b}^{(k+1)}\}$$



All layers

- Transformation at layer k :

$$F_{\Theta_k}(\mathbf{x}) = \sigma \left(\left(\mathbf{W}^{(k)} \right)^T \mathbf{x} + \mathbf{b}^{(k)} \right)$$

- Suppose the neural network has K layers (input excluded, output included)

$$\mathbf{x}^{(K)} = F_{\Theta_K}(\mathbf{x}^{(K-1)})$$

\vdots

$$\mathbf{x}^{(K)} = F_{\Theta_K} \left(F_{\Theta_{K-1}} \left(\dots F_{\Theta_1}(\mathbf{x}^{(0)}) \right) \right) = F_{\Theta}(\mathbf{x}^{(0)})$$

- Parameters are

$$\Theta = \{\mathbf{b}^{(1)}, \mathbf{W}^{(1)}, \mathbf{b}^{(2)}, \mathbf{W}^{(2)}, \dots, \mathbf{b}^{(K)}, \mathbf{W}^{(K)}\}$$

6

7

- Measure the discrepancy between
 - the neural network prediction $\hat{\mathbf{y}} = F_{\Theta}(\mathbf{x})$ with $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_{n_k})$
 - the expected output $\mathbf{y} = (y_1, \dots, y_{n_k})$
- Denoted $\ell(\hat{\mathbf{y}}, \mathbf{y})$
- Mean squared error (MSE) loss:

$$\ell(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{n_k} \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \frac{1}{n_k} \sum_{i=1}^{n_k} (\hat{y}_i - y_i)^2$$

8

- Adapted to probabilistic output and probabilistic observation (non-negative, sum to 1)
- Cross-entropy error

$$\text{CE}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{i=1}^{n_k} y_i \log \hat{y}_i$$

- For classification task, \mathbf{y} is one-hot encoded (all the y_i 's are zero except one), in that case

$$\text{CE}(\hat{\mathbf{y}}, \mathbf{y}) = - \log \hat{y}_k \quad \text{with } y_k = 1$$

9

- Parameter-free normalizing transform

$$\text{softmax} : \mathbb{R}^{n_k} \rightarrow \{(p_1, \dots, p_{n_k}) \in \mathbb{R}^{n_k}, p_i \geq 0, p_1 + \dots + p_{n_k} = 1\}$$

- $\mathbf{z} = \text{softmax}(\mathbf{x})$ defined by

$$z_i = \frac{\exp x_i}{\sum_{j=1}^{n_k} \exp x_j}$$

- $\text{softmax}(\mathbf{x} + c\mathbf{1}) = \text{softmax}(\mathbf{x})$, with $c \in \mathbb{R}$
- If $x_i \geq x_j$ then $\text{softmax}(\mathbf{x})_i \geq \text{softmax}(\mathbf{x})_j$
- Really an “arg softmax” rather than a “softmax”

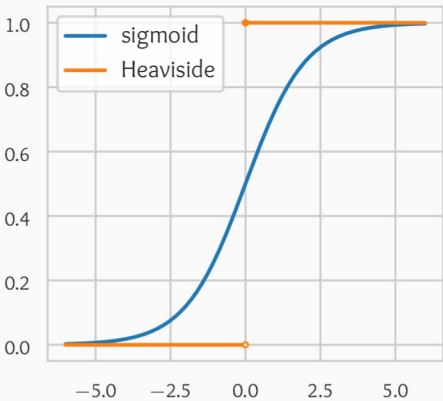
10

Activation functions

Logistic function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- Saturates at 0 and 1 when $x \rightarrow \pm\infty$
- Smooth version of Heaviside function
- $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
- Killing gradients

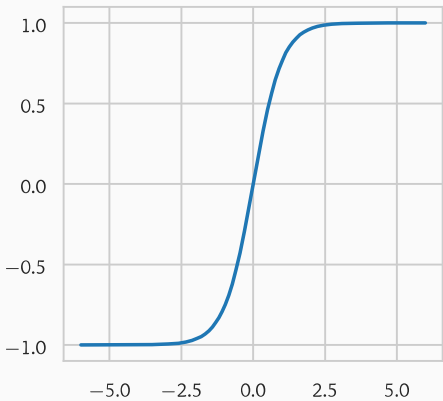


11

Hyperbolic tangent

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- Symmetric sigmoid
- Saturates at ± 1 when $x \rightarrow \pm\infty$
- $\tanh'(x) = 1 - \tanh^2(x)$
- Linearly related to the sigmoid function by
$$\tanh(x) = 2\sigma(2x) - 1$$
- Killing gradients

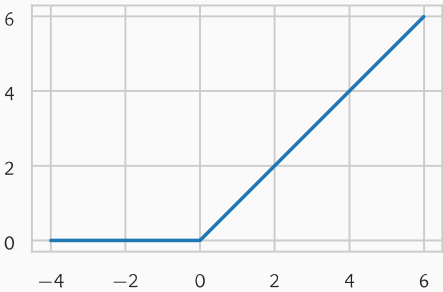


12

Rectified linear unit (ReLU)

$$\text{ReLU}(x) = \max(x, 0)$$

- Learns faster than sigmoid-like activation function
- Simpler to compute
- Provide sparsity of activations
- Dead ReLU: never activated across whole training set.
- Non negative activation function: zig-zag learning

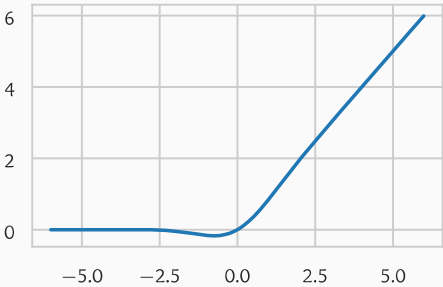


13

Gaussian Error Linear Unit (GELU), see Hendrycks and Gimpel 2023

$$\text{GELU}(x) = x\Phi(x)$$

- where Φ is the standard Gaussian cumulative distribution function
- Smoother version of ReLU
 - Mostly used in transformers



14

- Use ReLU
- Sigmoid not used anymore
- Prefer hyperbolic tangent to sigmoid

15

Gradient descent algorithms

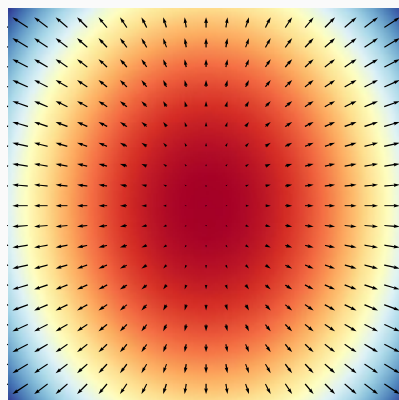
Gradient of a (scalar) function

- Differentiable function

$$f(x, y) = 1.5x^2 + y^2 + 2$$

- The gradient

$$\nabla_{(x,y)} f = \begin{pmatrix} 3x \\ 2y \end{pmatrix}$$



16

Gradient descent: motivation I

Gradient descent step improves current solution:

- Suppose \mathcal{L} is a differentiable function we want to minimize

$$\arg \min_{\theta \in \Theta} \mathcal{L}(\theta)$$

- Starting from the first order Taylor expansion. For a small $\|\mathbf{h}\|$ we have

$$\mathcal{L}(\theta + \mathbf{h}) \approx \mathcal{L}(\theta) + \langle \nabla_{\theta} \mathcal{L}, \mathbf{h} \rangle$$

- Choose $\mathbf{h} = -\eta \nabla_{\theta} \mathcal{L}$ (the gradient descent step), we have

$$\mathcal{L}(\theta - \eta \nabla_{\theta} \mathcal{L}) \approx \mathcal{L}(\theta) - \eta \|\nabla_{\theta} \mathcal{L}\|^2 < \mathcal{L}(\theta)$$

- $\theta - \eta \nabla_{\theta} \mathcal{L}$ is better than θ

17

Gradient descent step yields best update of linearized and regularized objective function

- Looking for the best $\theta' = \theta + \mathbf{h}$ around fixed θ
- Instead of minimizing $\mathcal{L}(\theta')$, we minimize

$$\mathcal{L}(\theta') \approx \mathcal{L}(\theta) + \langle \mathbf{h}, \nabla_{\theta} \mathcal{L} \rangle \quad (1)$$

- θ' should stay close to θ for (1) to hold, we penalize by $\|\theta - \theta'\| = \|\mathbf{h}\|$

$$\mathcal{L}(\theta) + \langle \mathbf{h}, \nabla_{\theta} \mathcal{L} \rangle + \frac{1}{2\eta} \|\mathbf{h}\|^2$$

- Minimizing w.r.t \mathbf{h} gives

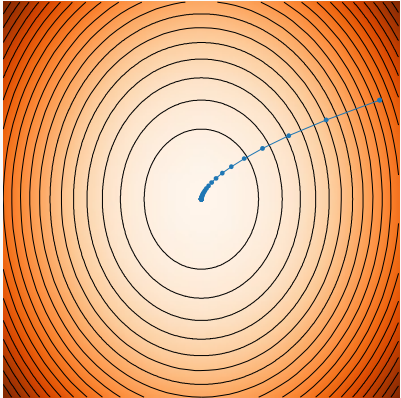
$$\theta' = \theta - \eta \nabla_{\theta} \mathcal{L}$$

18

- Starting point θ_0
- Gradient descent step

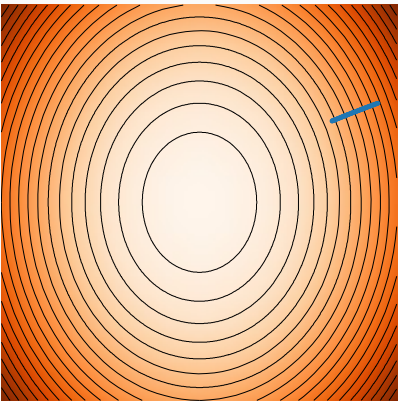
$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \mathcal{L}$$

- η is the *learning rate*
- Learning rate too small: slow convergence
- Learning rate too high: fluctuate around minimum or even diverge

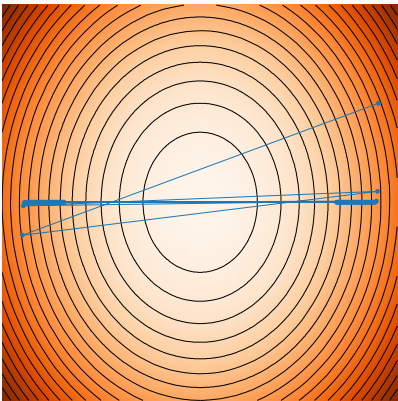


(a) $\eta = 0.1$

19



(a) $\eta = 0.001$



(b) $\eta = 0.665757$

20

Backpropagation algorithm

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f'_3 \circ f_2 \circ f_1)(\theta) \cdot (f'_2 \circ f_1)(\theta) \cdot f'_1(\theta) \\ &= f'_3((f_2 \circ f_1)(\theta)) \cdot f'_2(f_1(\theta)) \cdot f'_1(\theta) \end{aligned}$$

21

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f'_3 \circ f_2 \circ f_1)(\theta) \cdot (f'_2 \circ f_1)(\theta) \cdot f'_1(\theta) \\ &= f'_3((f_2 \circ f_1)(\theta)) \cdot f'_2(f_1(\theta)) \cdot f'_1(\theta) \end{aligned}$$

- Forward and backward pass

$$\theta \xrightarrow{f_1} f_1(\theta)$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f'_3 \circ f_2 \circ f_1)(\theta) \cdot (f'_2 \circ f_1)(\theta) \cdot f'_1(\theta) \\ &= f'_3((f_2 \circ f_1)(\theta)) \cdot f'_2(f_1(\theta)) \cdot f'_1(\theta) \end{aligned}$$

- Forward and backward pass

$$\theta \xrightarrow{f_1} f_1(\theta) \xrightarrow{f_2} f_2(f_1(\theta))$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f'_3 \circ f_2 \circ f_1)(\theta) \cdot (f'_2 \circ f_1)(\theta) \cdot f'_1(\theta) \\ &= f'_3((f_2 \circ f_1)(\theta)) \cdot f'_2(f_1(\theta)) \cdot f'_1(\theta) \end{aligned}$$

- Forward and backward pass

$$\theta \xrightarrow{f_1} f_1(\theta) \xrightarrow{f_2} f_2(f_1(\theta)) \xrightarrow{f_3} f_3(f_2(f_1(\theta))) = F(\theta)$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\begin{array}{ccccccc} \theta & \xrightarrow{f_1} & f_1(\theta) & \xrightarrow{f_2} & f_2(f_1(\theta)) & \xrightarrow{f_3} & f_3(f_2(f_1(\theta))) = F(\theta) \\ & & & & \downarrow & & \\ & & & & f_3'(f_2(f_1(\theta))) & & \end{array}$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\begin{array}{ccccccc} \theta & \xrightarrow{f_1} & f_1(\theta) & \xrightarrow{f_2} & f_2(f_1(\theta)) & \xrightarrow{f_3} & f_3(f_2(f_1(\theta))) = F(\theta) \\ & & \downarrow & & \downarrow & & \\ & & f_2'(f_1(\theta)) & \xleftarrow{\times} & f_3'(f_2(f_1(\theta))) & & \end{array}$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\begin{array}{ccccccc} \theta & \xrightarrow{f_1} & f_1(\theta) & \xrightarrow{f_2} & f_2(f_1(\theta)) & \xrightarrow{f_3} & f_3(f_2(f_1(\theta))) = F(\theta) \\ \downarrow & & \downarrow & & \downarrow & & \\ F'(\theta) = f_1'(\theta) & \xleftarrow{\times} & f_2'(f_1(\theta)) & \xleftarrow{\times} & f_3'(f_2(f_1(\theta))) & & \end{array}$$

Backpropagation

- Generalizable to \mathbb{R}^n to \mathbb{R}^p functions using jacobians
- Generalizable to functions having their own set of parameters
- Generalizable to a computational graph (DAG)

Total loss on the minibatch \mathcal{B} (whole training set for now)

$$\mathcal{L}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(F_{\Theta}(\mathbf{x}), \mathbf{y})$$

Differentiating w.r.t any scalar parameter θ (any $\mathbf{w}_{ij}^{(k)}$ or $\mathbf{b}_j^{(k)}$)

$$\frac{\partial \mathcal{L}_{\mathcal{B}}}{\partial \theta} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \frac{\partial \ell(\mathbf{x}^{(K)}, \mathbf{y})}{\partial \theta} \quad \text{with } \mathbf{x}^{(K)} = \hat{\mathbf{y}} = F_{\Theta}(\mathbf{x})$$

It suffices to compute

$$\frac{\partial \ell(\mathbf{x}^{(K)}, \mathbf{y})}{\partial \theta}$$

23

• MSE loss is defined by

$$\ell(\mathbf{x}^{(K)}, \mathbf{y}) = \frac{1}{n_K} \sum_{i=1}^{n_K} (\mathbf{x}_i^{(K)} - \mathbf{y}_i)^2$$

- Differentiating with respect to some scalar parameter $\theta = \mathbf{b}_q^{(k)}$ or $\theta = \mathbf{w}_{pq}^{(k)}$
- ($\mathbf{x}_i^{(K)}$ and \mathbf{y}_i are scalars)

$$\begin{aligned} \frac{\partial \ell(\mathbf{x}^{(K)}, \mathbf{y})}{\partial \theta} &= \frac{1}{n_K} \sum_{i=1}^{n_K} \frac{\partial (\mathbf{x}_i^{(K)} - \mathbf{y}_i)^2}{\partial \mathbf{x}_i^{(K)}} \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta} \\ &= \frac{2}{n_K} \sum_{i=1}^{n_K} (\mathbf{x}_i^{(K)} - \mathbf{y}_i) \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta} \end{aligned}$$

- We need to know $\frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta}$ now!

24

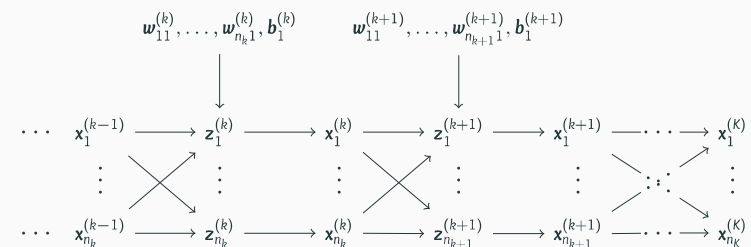
- Cross entropy loss: $\ell(\mathbf{x}^{(K)}, \mathbf{y}) = -\log \sigma_y(\mathbf{x}^{(K)})$
- Differentiating with respect to some scalar parameter $\theta = \mathbf{b}_q^{(k)}$ or $\theta = \mathbf{w}_{pq}^{(k)}$

$$\begin{aligned} \frac{\partial \ell(\mathbf{x}^{(K)}, \mathbf{y})}{\partial \theta} &= -\sum_{i=1}^{n_K} \frac{\partial \log \sigma_y(\mathbf{x}^{(K)})}{\partial \mathbf{x}_i^{(K)}} \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta} \\ &= -\sum_{i=1}^{n_K} (\mathbf{y}_i - \sigma_i(\mathbf{x})) \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta} \end{aligned}$$

- We need to know $\frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta}$ now!

25

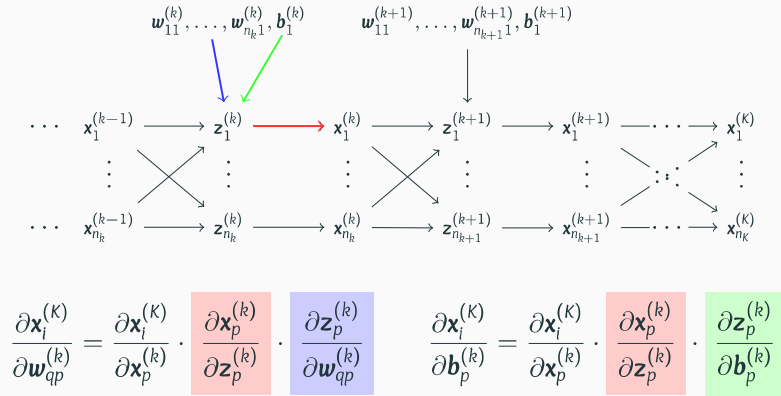
- Define $\mathbf{z}_i^{(k)} = \langle \mathbf{w}_i^{(k)}, \mathbf{x}^{(k-1)} \rangle + \mathbf{b}_i^{(k)}$ so that we have $\mathbf{x}_i^{(k)} = \sigma(\mathbf{z}_i^{(k)})$
- Computational graph for $\mathbf{x}_1^{(k-1)}, \mathbf{x}_1^{(k)}, \mathbf{x}_1^{(k+1)}$ only!



26

Gradient of last layer w.r.t parameters I

Computational graph for neuron 1 at layer k ($p = 1$):



27

Gradient of last layer w.r.t parameters II

Given that $\mathbf{x}_p^{(k)} = \sigma(\mathbf{z}_p^{(k)})$, we have $\frac{\partial \mathbf{x}_p^{(k)}}{\partial \mathbf{z}_p^{(k)}} = \sigma'(\mathbf{z}_p^{(k)})$

And $\mathbf{z}_p^{(k)} = \langle \mathbf{w}_p^{(k)}, \mathbf{x}^{(k-1)} \rangle + \mathbf{b}_p^{(k)}$, so $\frac{\partial \mathbf{z}_p^{(k)}}{\partial \mathbf{w}_{qp}^{(k)}} = \mathbf{x}_q^{(k-1)}$ and $\frac{\partial \mathbf{z}_p^{(k)}}{\partial \mathbf{b}_p^{(k)}} = 1$

Finally

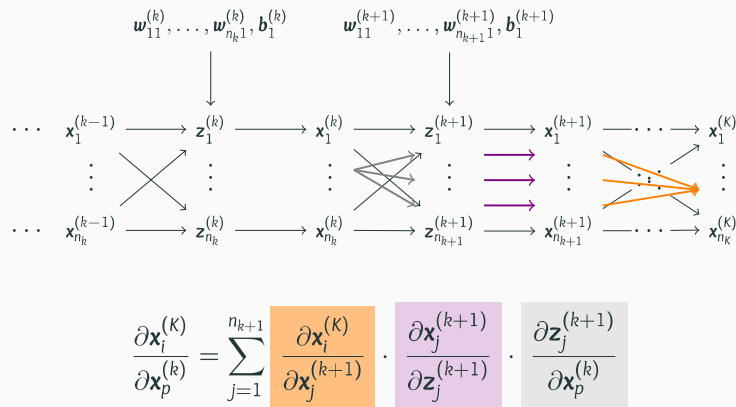
$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{w}_{qp}^{(k)}} = \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)}) \cdot \mathbf{x}_q^{(k-1)}$$

$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{b}_p^{(k)}} = \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)})$$

Need to compute $\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}}$ now!

28

Gradient of last layer w.r.t other layer



29

Backpropagating from next layer

Given that $\mathbf{x}_j^{(k+1)} = \sigma(\mathbf{z}_j^{(k+1)})$, we have

$$\frac{\partial \mathbf{x}_j^{(k+1)}}{\partial \mathbf{z}_j^{(k+1)}} = \sigma'(\mathbf{z}_j^{(k+1)})$$

And $\mathbf{z}_j^{(k+1)} = \langle \mathbf{w}_j^{(k+1)}, \mathbf{x}^{(k)} \rangle + \mathbf{b}_j^{(k+1)}$, so

$$\frac{\partial \mathbf{z}_j^{(k+1)}}{\partial \mathbf{x}_p^{(k)}} = \mathbf{w}_{pj}^{(k+1)}$$

Replacing in last equation, we have

$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} = \sum_{j=1}^{n_{k+1}} \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_j^{(k+1)}} \cdot \sigma'(\mathbf{z}_j^{(k+1)}) \cdot \mathbf{w}_{pj}^{(k+1)}$$

30

- Backpropagation equations for parameters $\mathbf{w}_{qp}^{(k)}$ and $\mathbf{b}_p^{(k)}$

$$\frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{w}_{qp}^{(k)}} = \frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)}) \cdot \mathbf{x}_q^{(k-1)} \quad (2)$$

$$\frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{b}_p^{(k)}} = \frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)}) \quad (3)$$

- Backpropagation equation

$$\frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_p^{(k)}} = \sum_{j=1}^{n_{k+1}} \frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_j^{(k+1)}} \cdot \sigma'(\mathbf{z}_j^{(k+1)}) \cdot \mathbf{w}_{pj}^{(k+1)} \quad (4)$$

- Only one sweep backward is necessary to compute all gradients

- [1] Sergey Ioffe and Christian Szegedy. “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”. Feb. 10, 2015. arXiv: 1502.03167 [cs]. URL: <http://arxiv.org/abs/1502.03167> (visited on 11/06/2017).
- [2] Ian Goodfellow et al. *Deep Learning*. Vol. 1. MIT press Cambridge, 2016.
- [3] Aston Zhang et al. *Dive into Deep Learning*. 2020.
- [4] Dan Hendrycks and Kevin Gimpel. *Gaussian Error Linear Units (GELUs)*. June 5, 2023. DOI: 10.48550/arXiv.1606.08415. arXiv: 1606.08415 [cs]. URL: <http://arxiv.org/abs/1606.08415> (visited on 11/07/2023). Pre-published.