

A grayscale background image of a large, ornate classical building with a central dome and multiple columns, likely a university building.

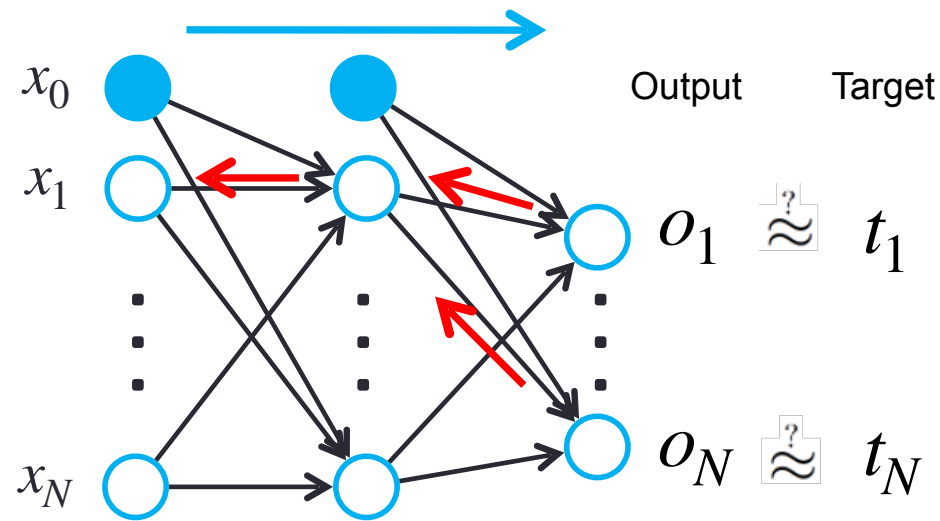
Neural Networks - part 2

Practicals Machine Learning 1, SS23

Ceca Kraišniković
Institute of Theoretical Computer Science
Graz University of Technology, Austria

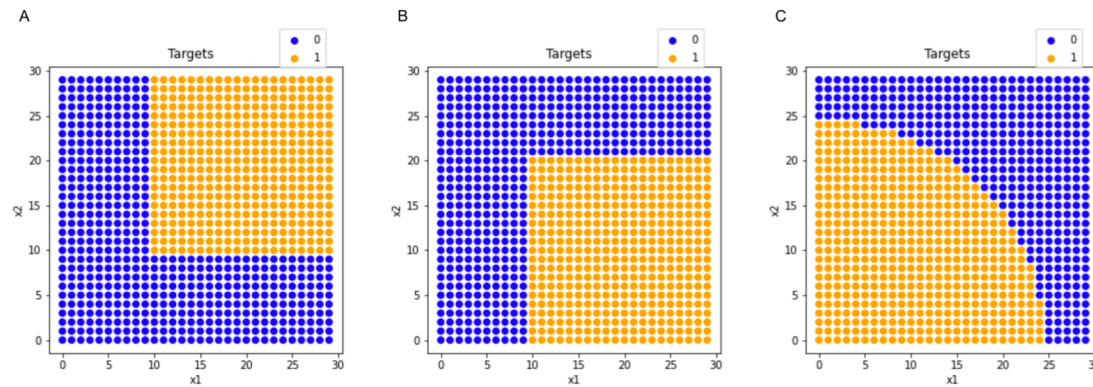
Neural Networks training - Backpropagation algorithm

- **Forward:** Calculate activations and outputs of all neurons
- **Backward:** Calculate errors and propagate them back



Neural Networks training

- Training:
 - Feed-forward pass
 - Backward pass
- Task: (binary) classification task, with 2 input features



- Example network architecture: 2 input neurons ($n_{in} = 2$), n_{hid} neurons in the hidden layer, 1 output neuron
- Appropriate loss function: binary cross entropy, for each training example i

$$-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

(Notation changes here: outputs \hat{o} are denoted now as predictions \hat{y} , and targets t as y .)

Neural Networks training

- FF pass

$$z_1(x) = W_1^T x + b_1$$

$$h_1(z_1) = \text{ReLu}(z_1)$$

$$z_2(h_1) = W_2^T h_1 + b_2$$

$$h_2(z_2) = \sigma(z_2) = \hat{y}$$

- Calculate loss for each example, and accumulate it (at init, $\mathcal{L} = 0$)

$$\mathcal{L} += -y \log(h_2) - (1 - y) \log(1 - h_2)$$

Neural Networks training

- FF pass

$$z_1(x) = W_1^T x + b_1$$

$$h_1(z_1) = \text{ReLu}(z_1)$$

$$z_2(h_1) = W_2^T h_1 + b_2$$

$$h_2(z_2) = \sigma(z_2) = \hat{y}$$

Variables and their dimensions (*shapes*):

$$x : (n_in, 1)$$

$$W_1 : (n_in, n_hid)$$

$$b_1 : (n_hid, 1)$$

$$z_1 : (n_hid, 1)$$

$$h_1 : (n_hid, 1)$$

$$W_2 : (n_hid, n_out)$$

$$z_2 : (n_out, 1)$$

$$h_2 : (n_out, 1)$$

$$b_2 : (n_out, 1)$$

- Calculate loss for each example, and accumulate it (at init, $\mathcal{L} = 0$)

$$\mathcal{L} += -y \log(h_2) - (1 - y) \log(1 - h_2)$$

Neural Networks training

- FF pass

$$z_1(x) = W_1^T x + b_1$$

$$h_1(z_1) = \text{ReLu}(z_1)$$

$$z_2(h_1) = W_2^T h_1 + b_2$$

$$h_2(z_2) = \sigma(z_2) = \hat{y}$$

Variables and their dimensions (*shapes*):

$$x : (n_in, 1)$$

$$W_1 : (n_in, n_hid)$$

$$b_1 : (n_hid, 1)$$

$$z_1 : (n_hid, 1)$$

$$h_1 : (n_hid, 1)$$

$$W_2 : (n_hid, n_out)$$

$$z_2 : (n_out, 1)$$

$$h_2 : (n_out, 1)$$

$$b_2 : (n_out, 1)$$

- Calculate loss for each example, and accumulate it (at init, $\mathcal{L} = 0$)

$$\mathcal{L} += -y \log(h_2) - (1 - y) \log(1 - h_2)$$

Variables to optimize:

$$W_1 : (n_in, n_hid)$$

$$b_1 : (n_hid, 1)$$

$$W_2 : (n_hid, n_out)$$

$$b_2 : (n_out, 1)$$

Neural Networks training

- Backward pass

$$\frac{\partial \mathcal{L}}{\partial W_2} = ?$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = ?$$

$$\frac{\partial \mathcal{L}}{\partial W_1} = ?$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = ?$$

8 Matrix Calculus

Scalar by Vector ($\mathbb{R}^n \rightarrow \mathbb{R}$ mapping): $\frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1} \quad \dots \quad \frac{\partial y}{\partial x_n} \right]$ (one row, n columns)

Vector by Vector ($\mathbb{R}^n \rightarrow \mathbb{R}^m$ mapping), Jacobian matrix:
$$\begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix},$$

(m rows, n columns)

Scalar by Matrix: $\frac{\partial y}{\partial A} = \begin{bmatrix} \frac{\partial y}{\partial A_{11}} & \frac{\partial y}{\partial A_{12}} & \dots & \frac{\partial y}{\partial A_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial A_{m1}} & \frac{\partial y}{\partial A_{m2}} & \dots & \frac{\partial y}{\partial A_{mn}} \end{bmatrix},$

Vector by Matrix: ... (not needed in this tutorial).

Matrix Calculus

- **Shape rule:**
 - What dimensions do we need to get as a final result?
 - Scalar by vector
 - Vector by vector
 - Scalar by matrix
- **Dimension balancing:**
 - Enforcing shapes (when deriving over a linear transformation)
 - Should work in most practical settings

Matrix Calculus

Scalar by Vector ($\mathbb{R}^n \rightarrow \mathbb{R}$ mapping): $\frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1} \quad \dots \quad \frac{\partial y}{\partial x_n} \right]$, (one row, n columns)

Vector by Vector ($\mathbb{R}^n \rightarrow \mathbb{R}^m$ mapping), Jacobian matrix:

$$\begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix},$$

(m rows, n columns)

Scalar by Matrix: $\frac{\partial y}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial y}{\partial A_{11}} & \frac{\partial y}{\partial A_{12}} & \dots & \frac{\partial y}{\partial A_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial A_{m1}} & \frac{\partial y}{\partial A_{m2}} & \dots & \frac{\partial y}{\partial A_{mn}} \end{bmatrix}$, (dimensions of A)

Vector by Matrix: ... (not needed in this tutorial).

Neural Networks training

- Backward pass

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial W_2}$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial b_2}$$

$$\frac{\partial \mathcal{L}}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial z_1} \frac{\partial z_1}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial W_1}$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial \mathcal{L}}{\partial z_1} \frac{\partial z_1}{\partial b_1} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

Dimensions:

$$W_2 : (n_{hid}, n_{out})$$

$$b_2 : (n_{out}, 1)$$

$$W_1 : (n_{in}, n_{hid})$$

$$b_1 : (n_{hid}, 1)$$

- FF pass

$$z_1(x) = W_1^T x + b_1$$

$$h_1(z_1) = \text{ReLU}(z_1)$$

$$z_2(h_1) = W_2^T h_1 + b_2$$

$$h_2(z_2) = \sigma(z_2) = \hat{y}$$

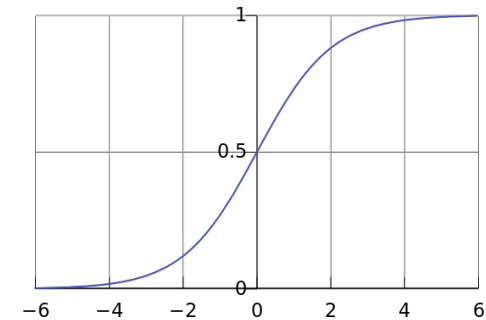
$$\mathcal{L}(\hat{y}, y)$$

Recap: Sigmoid function

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

$$\begin{aligned} \boxed{\frac{d\sigma(x)}{dx}} &= \frac{e^x(e^x + 1) - e^x e^x}{(e^x + 1)^2} \\ &= \frac{e^x}{(e^x + 1)^2} \\ &= \frac{e^x}{e^x + 1} \cdot \frac{1}{e^x + 1} \\ &= \frac{e^x}{e^x + 1} \cdot \frac{1}{1 + e^x - e^x} \\ &= \frac{e^x + 1}{e^x + 1} \cdot \frac{e^x + 1}{e^x + 1} \\ &= \boxed{\sigma(x)(1 - \sigma(x))} \end{aligned}$$

Sigmoid (logistic) function:



Neural Networks training

- Backward pass

$$\frac{\partial \mathcal{L}}{\partial z_2} = ?$$

First, rewrite \mathcal{L} :

$$\begin{aligned}\mathcal{L}(z_2) &= -y \log(\sigma(z_2)) - (1 - y) \log(1 - \sigma(z_2)) = \\ &= -y \log\left(\frac{e^{z_2}}{1 + e^{z_2}}\right) - (1 - y) \log\left(1 - \frac{e^{z_2}}{1 + e^{z_2}}\right) = \\ &= -y \log(e^{z_2}) + y \log(1 + e^{z_2}) - (1 - y) \log(1) + (1 - y) \log(1 + e^{z_2}) = \\ &= -yz_2 + \log(1 + e^{z_2})\end{aligned}$$

Then, find derivative w.r.t. z_2 :

$$\begin{aligned}&\left(\text{using } \frac{d(\log x)}{dx} = \frac{1}{x}\right) \\ \frac{\partial \mathcal{L}}{\partial z_2} &= -y + \frac{1}{1 + e^{z_2}} e^{z_2} = -y + \sigma(z_2) = h_2 - y\end{aligned}$$

Activation function gradients

- Activation functions are applied element-wise.
- The input and the output are both vectors, and of the same dimension.
- Case: Vector by Vector.

$$h = f(z), \mathbb{R}^n \rightarrow \mathbb{R}^n,$$

$$\frac{\partial h}{\partial z} = \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{\partial h_n}{\partial z_n} \end{bmatrix}$$

($n \times n$ matrix, with elements on the main diagonal, 0s elsewhere)