

Neural networks

Motivation:

- we want a powerful model which can describe arbitrary functional dependencies
- ↳ usually translates into large number of tunable / trainable parameters
- ↳ but has to be computationally efficient (evaluation of prediction $f_w(x)$ and gradients $\nabla_w f_w(x)$)

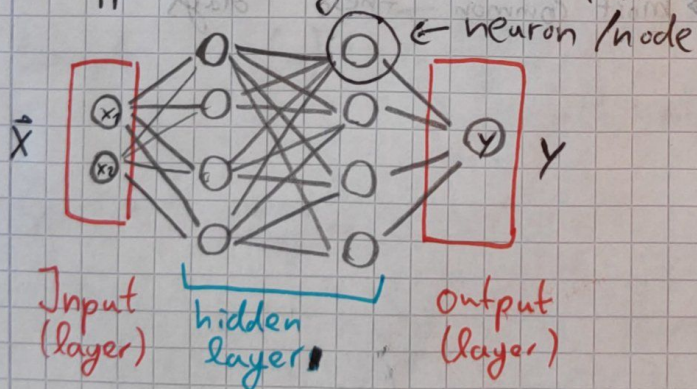
↳ ~~for training~~
for training (next section)

Idea:

- stack many simple operations on top of each other!

Fully connected dense neural networks - dense layer

- ↳ proto-type of neural networks
- ↳ appears in many architectures (almost all)



hidden layer: linear operation + non-linear operation

$$\vec{z} = W \cdot \vec{x} + \vec{b}$$

weight matrix of tunable parameters

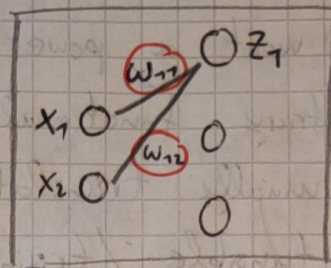
$$y = \sigma(z)$$

↑ non-linear function

weight matrix $W = \begin{pmatrix} w_{11} & w_{12} & \dots \\ w_{21} & w_{22} & \dots \\ \vdots & \vdots & \ddots \\ \vdots & \vdots & \dots & w_{MN} \end{pmatrix} \left. \vphantom{\begin{pmatrix} w_{11} & w_{12} & \dots \\ w_{21} & w_{22} & \dots \\ \vdots & \vdots & \ddots \\ \vdots & \vdots & \dots & w_{MN} \end{pmatrix}} \right\} N$

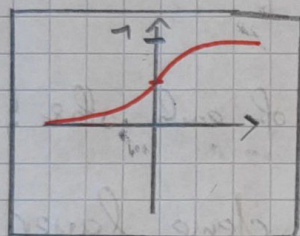
M

N : dimension of input \vec{x}
 M : number of units/neurons

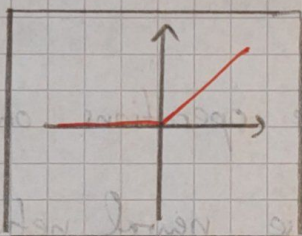


Layered structure: $\dots \sigma(w_2 \sigma(w_1 \vec{x} + \vec{b}_1) + \vec{b}_2) \dots$

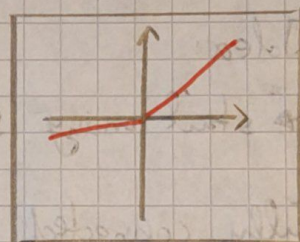
Non linearity? ← activation function



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



ReLU
 (rectified linear unit)



leaky ReLU

→ historically important, only used for last layer in classification (or softmax: higher dim. version)

→ most common these days

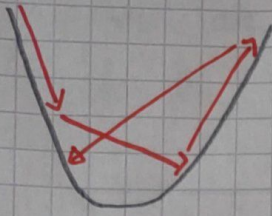
to wrap up:

→ trainable parameters: $w_{ij}^l, b_j^l \leftarrow \text{layer}$

→ but what about:

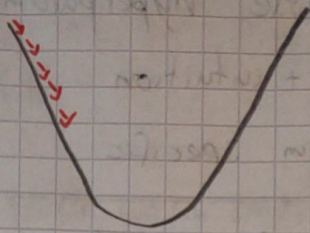
- ↳ number of inner layers
- ↳ number of units per layer
- ↳ what activation function

} hyperparameters!
 → have to be chosen by hand



large learning rate

→ can't reach minimum



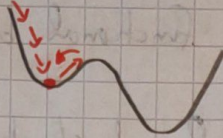
small learning rate

→ too many steps, expensive

problems:

→ stuck in local minimum

→ gradients are expensive for large datasets



⇒ Stochastic gradient descent

$$L(w) = \sum_{i=1}^N l_i(w) \quad | \quad \text{example } l_i(w) = |f(x_i) - y_i|^2$$

Idea: Approximate gradients with batches of data

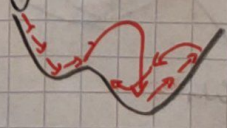
$$\begin{aligned} \nabla_w L(w) &= \nabla_w \sum_{i=1}^N l_i(w) & \longrightarrow & \sum_{i=1}^M \nabla_w l_i(w) \\ &= \sum_{i=1}^N \nabla_w l_i(w) & & \text{with } M < N \end{aligned}$$

M: minibatch size

Epoch: One loop over all minibatches

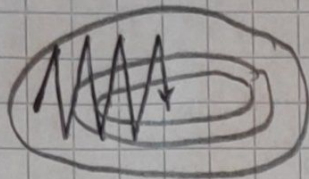
→ adds stochasticity which can help the algorithm to get out of a local minimum

→ less expensive!

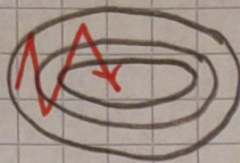


• there are more algorithms but they build up on the idea of SGD

↳ SGD + momentum: speeds up in direction with small but consistent gradients!



SGD



SGD + momentum

↳ Adam: very commonly used; uses ^{in addition} $(\nabla_w L(w))^2$
(gradient squared)

~~Summary so far:~~

~~↳ we know how basic neural networks work~~

~~→ provide us with powerful models~~

~~↳ we know how to train them~~

~~We have~~

to wrap up:

→ Training is just finding a minimum of the loss function

→ we have efficient algorithms to train neural networks:

SGD, Adam, ...