# MLP: arch, backpropagation & symmetries
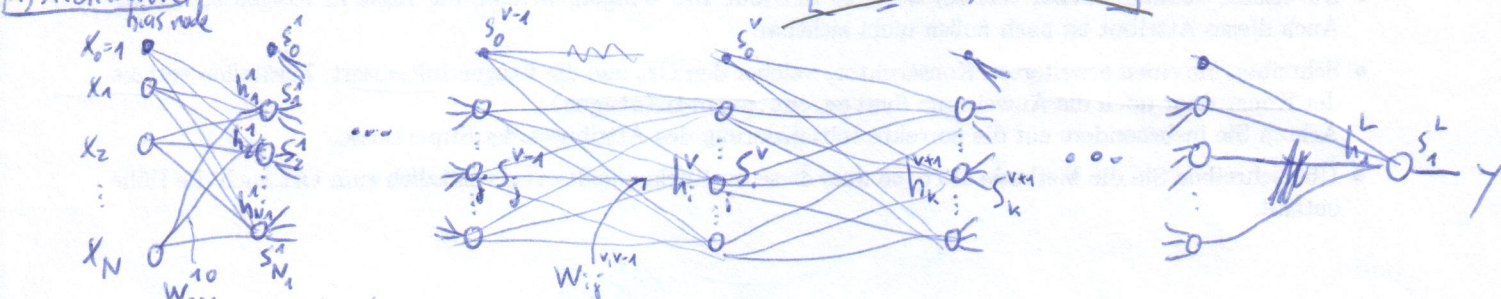
## a) Derivation of backpropagation (incl. architecture of MLP and gradient descent)

~~"will be discussed at the end"~~

feed forward MLP: $Y(\underline{x}) \in \mathbb{R}$ [or ~~$Y(\underline{x}) \in \mathbb{R}^M$~~] for $\underline{x} \in \mathbb{R}^N$

**1) Architecture**



| layer: | input | first hidden | ... | v-1th hidden | v-th hidden | v+1-th hidden | ... | L-1-th hidden | output |
|---|---|---|---|---|---|---|---|---|---|
| index: | 0 | 1 | | v-1 | v | v+1 | | L-1 | L |
| size (no. neurons) | N | $N_1$ | | $N_{v-1}$ | $N_v$ | $N_{v+1}$ | | $N_{L-1}$ | |
| transfer function | | $f_1$ | | $f_{v-1}$ | $f_v$ | $f_{v+1}$ | | $f_{L-1}$ | $f_L$ |

often $f_v = f_{hidden}$ for $v = 1, ..., L-1$ (tanh or rectLU) typical examples

$f_{output}$ (often: linear, logistic or sigmoid, softmax)

**2) Evaluation by forward propagation**

input → 1st layer:

$s_2^1 = f_1(h_2^1)$ with $h_2^1 = \sum_{k=0}^{N} x_k W_{2k}^{1,0}$ ← from layer 0 to 1

↑ output of neuron 2 in layer 1, input to that neuron

from neuron k (of layer 0) to neuron 2 (of layer 1)

In general: for each layer $v = 1, ..., L$:

$$\boxed{s_i^v = f_v(h_i^v) \text{ with } h_i^v = \sum_{k=0}^{N_{v-1}} s_k^{v-1} W_{ik}^{v,v-1} \quad , i = 1, ..., N_v}$$

↳ output of neuron k in layer v-1

including the input → 1st layer by $s_n^0 = x_k$, $k = 0, ..., N = N_0$; bias nodes: $s_0^v = 1 \; \forall v$

Output: $\underline{Y = s_1^L = f_L(h_1^L)}$ with $h_1^L = \sum_{k=0}^{N_{L-1}} s_k^{L-1} W_{1k}^{L,L-1}$

$\to$ **MLP**: function $Y(\underline{x}; \underline{w})$ parametrized by $\underline{w} = \left\{ W_{ik}^{v,v-1} \right\}_{\substack{v=1,...,L \\ i=1,...,N_v \\ k=0,...,N_{v-1}}}$

next: how to learn parameter values $\underline{w}$ from data $\left\{ \left( \underline{x}^{(\alpha)}, Y_T^{(\alpha)} \right) \right\}_{\alpha=1}^{P}$

# MLP: hit arch., backprop., symmetries (cont'd)

## 3) Parameter optimization through gradient descent

Aim: find optimal parameter values $\underline{w}$ by minimizing the training cost:

$$E^T(\underline{w}) = \frac{1}{P} \sum_{\alpha=1}^{P} e^{(\alpha)}(\underline{w}) \overset{!}{=} \min_{\underline{w}}$$



Typical (neural networks) approach: descent the gradient, stopping at stationary point $\frac{\partial E^T}{\partial \underline{w}}(\underline{w}) = \underline{0}$:

initialize weights with $\underline{w}^{(0)}$

until convergence: $\underline{w}^{(n+1)} = \underline{w}^{(n)} - \eta \frac{\partial E^T}{\partial \underline{w}}(\underline{w}^{(n)}) \overset{\text{linearity of gradient}}{=} \underline{w}^{(n)} - \eta \frac{1}{P} \sum_{\alpha=1}^{P} \frac{\partial e^{(\alpha)}}{\partial \underline{w}}(\underline{w}^{(n)})$

— learning step

descent

gradient of error of datapt. $\alpha$

$\rightarrow$ per iteration $n$, per data point $\alpha$ required:

$$\frac{\partial e^{(\alpha)}}{\partial \underline{w}} = \frac{\partial e^{(\alpha)}}{\partial Y} \frac{\partial Y}{\partial \underline{w}}$$ (crossed out)

dep. on cost fct., e.g.

ii) $\frac{\partial e^{(\alpha)}}{\partial Y} = \frac{\partial}{\partial Y}\left(\frac{1}{2}\left[Y(\underline{x}^{(\alpha)};\underline{w}) - Y_T^{(\alpha)}\right]^2\right)$

$= Y(\underline{x}^{(\alpha)};\underline{w}) - Y_T^{(\alpha)}$

cross entropy: homework

$$\frac{\partial e^{(\alpha)}}{\partial \underline{w}}\left(Y(\underline{x}^{(\alpha)};\underline{w}), Y_T^{(\alpha)}\right) = \frac{\partial e^{(\alpha)}}{\partial Y} \cdot \frac{\partial Y}{\partial \underline{w}}(\underline{x}^{(\alpha)};\underline{w})$$

dep. on neural network

Backpropagation: efficient evaluation of using the chain rule: $O(\#weights)$
algorithm

runtime complexity instead of $O(\#weights^2)$!

$$\#weights = N_1(N+1) + \sum_{v=2}^{L} N_v(N_{v-1}+1) = 10N^2 + 10N$$

e.g. $L=10, N_v = N$

## 4) Backpropagation algorithm

component-wise (simplest):

$$\frac{\partial y}{\partial w_{ij}^{v,v-1}} = \frac{\partial y}{\partial h_i^v} \frac{\partial h_i^v}{\partial w_{ij}^{v,v-1}} = \delta_i^v s_j^{v-1} = \begin{cases} \delta_i^v f_{v-1}(h_j^{v-1}) & ,v=2,...,L \\ \delta_i^v x_j & ,v=1 \\ \delta_i^v & ,j=0 \end{cases} = \begin{cases} \delta_i^v & ,j=0 \\ \delta_i^v f_{v-1}(h_j^{v-1}) & else \end{cases}$$

3 cases: - "normal" weight
- input weight
- bias node weight

"local error
of neuron $(i,v)$": $\delta_i^v :=$     $= s_j^{v-1}$

$$\delta_i^v = \frac{\partial y}{\partial h_i^v} = \frac{\partial y}{\partial s_i^v} \frac{\partial s_i^v}{\partial h_i^v} = f_v'(h_i^v) \cdot \sum_{k=1}^{N_{v+1}} \frac{\partial y}{\partial h_k^{v+1}} \frac{\partial h_k^{v+1}}{\partial s_i^v} = f_v'(h_i^v) \sum_{k=1}^{N_{v+1}} \delta_k^{v+1} w_{ki}^{v+1,v}$$

$$= f_v'(h_i^v) \quad y\left(h_1^{v+1}(s_i^v),...,h_{N_{v+1}}^{v+1}(s_i^v)\right) \quad \text{general chain rule}$$

by def. $h_j^0 := x_j$
and $f_0 := id$
↑ identity function $id(x)=x$

$\rightarrow$ recursive relation of local errors: $\{\delta_k^{v+1}\}_{k=1}^{N_{v+1}}$ known $\Rightarrow \{\delta_i^v\}_{i=1}^{N_v}$ can be calculated
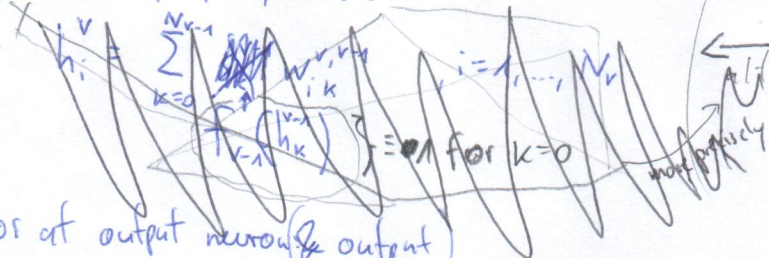(i.e. backpropagation of local errors)

## backprop algorithm:

given: $\underline{X}, \underline{Y}_T, \underline{W}$

define: $h_k^0 = x_k$ for $k=0,...,N$; $f_0=id$; $N_0 = N$

forward propagation:

for $v=1,2,...,L$ do



$$h_i^v = \sum_{k=1}^{N_{v-1}} f_{v-1}(h_k^{v-1}) w_{ik}^{v,v-1} + w_{i0}^{v,v-1}$$

local errors at output neuron & output:

$$\delta_1^L = f_L'(h_1^L)$$
$$y = f_L(h_1^L)$$

backward propagation:

for $v=L-1,...,1$ do

$$\delta_i^v = f_v'(h_i^v) \sum_{k=1}^{N_{v+1}} \delta_k^{v+1} w_{ki}^{v+1,v} \quad, i=1,...,N_v$$

$$\frac{\partial y}{\partial w_{ij}^{v,v-1}} = \begin{cases} \delta_i^v f_{v-1}(h_j^{v-1}) & , else \\ \delta_i^v & , j=0 \end{cases}$$

gradient of network output:

$$\frac{\partial y}{\partial w_{ij}^{v,v-1}} = \quad , v=1,...,L \quad, i=1,...,N_v \quad, j=0,...,N_{v-1}$$

# MLP: ass, backprop, symmetries (cont'd)

**b) Consequences of parameter space symmetries.**

i) <u>permutation of ~~weight~~ indices within a layer</u>

Let

$\boxed{\text{neuron}}$

$\overset{\text{[M]}}{}$ $\underline{w}$ be given parameter vector, then $\underline{\tilde{w}}$ ~~ass nax~~ will make the same cost $E^{T}$ ass $\underline{w}$

if $\tilde{w}_{ij}^{r,v-1} = w_{\pi(i)j}^{v,v-1}$ for some permutation $\pi$, <u>and</u> $\tilde{w}_{k,i}^{v+1,v} = w_{k,\pi(i)}^{v+1,v}$

$\rightarrow$ neuron $\pi(i)$ $(\underline{w})$ $\rightarrow$ $i$ $(\underline{\tilde{w}})$

$\Rightarrow N_v!$ permutations per hidden layer $v$

ii) <u>sign reversal across layers</u>

Let $\underline{w}$ be given parameter, then $\underline{\tilde{w}}$ will produce a model with same cost.

if $\tilde{w}_{ij}^{v,v-1} = -w_{ij}^{v,v-1}$ and $\tilde{w}_{ki}^{v+1,v} = -w_{ki}^{v+1,v}$ and $f_{hidden} = \tanh$

(follows from $\tanh(-h) = -\tanh(h)$)

$\Rightarrow 2^{N_v}$ combinations per hidden layer $v$

$\Rightarrow$ overall $\prod\limits_{v=1}^{L} N_v! \, 2^{N_v}$ equivalent solution

(also for global or local minima true)

$\rightarrow$ no unique optimal parameters! but at last $\prod\limits_{v=1}^{L} N_v! \, 2^{N_v}$

equivalent solutions (same cost)