

## Neuron

- dendrites (inputs)
- axon (single output)
- synapses (branching)
- $8,6 \times 10^{10}$  neurons ( $8 \times 10^5 / \text{mm}^3$ )
- short-memory = signal circulation
- long-term memory = weights updating

## Standardization

- interval scaling

$$\begin{cases} [-1, 1] \\ [0, 1] \end{cases}$$

$$V_{ij} = \sum_h (x_{hi} - \mu_i) \cdot (x_{hj} - \mu_j)$$

- decimal scaling: dělím nejmenším výsledkem 10, aby  $\forall x : x \in [-1, 1]$

$$\left. \begin{array}{l} \text{- MAD - scale: "mean absolute deviation"} \\ \frac{1}{n} \sum |x - \bar{x}| \end{array} \right\} = \frac{x - \bar{x}}{\text{deviations}}$$

$$\left. \begin{array}{l} \text{- MSD - scale: "mean standard deviation"} \\ \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} \end{array} \right\}$$

## All Transform n PCA

$$X^m \rightarrow X^p : m \gg p$$

$$X = \begin{pmatrix} \dots & x_{1j} & \dots \\ \vdots & \vdots & \vdots \\ \dots & x_{uj} & \dots \end{pmatrix}$$

Vždy itemji jsou sloupcy,  $\mu_j = \frac{1}{n} \sum^u x_{ij}$   $M = (\dots \mu_j \dots \dots)$

Přitom se sčítá

2. jistit, které ughadit. Kovariaci matice:  $W_{ij} = W_{ji} = \frac{1}{n} \sum^u (x_{ki} - \mu_i) \cdot (x_{kj} - \mu_j)$

nejzájemnější vlastní čísla a jejich vlastní vektory odpovídají

nejzájemnějším směrem

redukované vektory  $y = V^T X$

## Normalní rozdělení:

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$\mathbb{E}(x) = \mu, \text{ resp. } \sigma^2$$

## Hypotheses testing

$$\text{Error}_S(h) = P_{x \in S} [f(x) \neq h(x)] = \frac{r}{|S|}$$

- past, že jsem se netrefil

mismatched

## k-Fold Validation:

- po jednotlivých podmnožinách průchodu error dvanácti hypotez

$$\delta_i = \text{err}_{T_i}(A) - \text{err}_{T_i}(B)$$

$$\bar{\delta} = \frac{1}{k} \sum_i \delta_i$$

$$\text{Conf. int.: } \bar{\delta} \pm t_{N, k-1} \cdot \sqrt{\frac{\frac{1}{k-1} \sum_i (\delta_i - \bar{\delta})^2}{N}}$$

estimated std. dev.

When having binomial distr.:

95% conf. int.:

$$\text{error}_S(h) \pm 1,96 \sqrt{\frac{\text{err}_S(h) \cdot (1 - \text{err}_S(h))}{n}}$$

Can compare two hyps.:

$$\Delta = \text{err}_{S_1}(h_1) - \text{err}_{S_2}(h_2)$$

$$\text{Conf. int.: } \Delta \pm 2\sqrt{\frac{\text{err}_{S_1}(h_1) \cdot (1 - \text{err}_{S_1}(h_1)) + \text{err}_{S_2}(h_2) \cdot (1 - \text{err}_{S_2}(h_2))}{n}}$$

↑ Requires independent testsets. ↑

Must be on identical test sets.

## Perception

$$x^T w - \gamma \geq 0 \Rightarrow y=1$$

$$x^T w - \gamma < 0 \Rightarrow y=0$$

## Lineární segmentabilita:

-> pokud je perceptronum

segmentovatelný mimo A, B.

neuron

active	: $\sigma(w^T x - \gamma) = 1$
silent	: $\sigma(w^T x - \gamma) = \frac{1}{2}$
passive	: $\sigma(w^T x - \gamma) = 0$

Pokud  $t_i = 0$ , ale  $y_i = 1$ :  $w_i = w_i - \alpha x_i$

$t_i = 1$ , ale  $y_i = 0$ :  $w_i = w_i + \alpha x_i$

$$w' = (w^{**}, -\gamma) \\ x = (x^{**}, 1) \rightarrow y=1 \Leftrightarrow x^T w \geq 0$$

## MLP

**error function:**

$$\frac{1}{2} \sum_p \sum_j (y_{j,p} - d_{j,p})^2$$

**adjustment:**

$$w_{ij} = w_{ij} + \alpha \Delta_E w_{ij}$$

$\rightarrow$  minor best abstractions

o moment

$$\alpha_m (E(w(t)) - E(w(t+1)))$$

$$= y_i \cdot (\sigma_j)$$

$$(d_j - y_j) \cdot y_j \cdot (1 - y_j)$$

output neuron

$$\left( \sum_h \sigma_h w_{jh} \right) \cdot y_j \cdot (1 - y_j)$$

hidden neurons

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

$$\text{Softmax } S(x)_j = \frac{e^{-x_j}}{\sum_i e^{-x_i}}$$

$$\frac{\partial y_j}{\partial x_i} = y_j (1 - y_j)$$

$$\frac{\partial y_j}{\partial x_n} = -y_j \cdot y_n$$

$$\Delta w(i+1)$$

**Momentum - learning:**

$$w(i+1) = w(i) - \alpha \frac{\partial E}{\partial w(i)} + \alpha_m (w(i) - w(i-1))$$

keeping the previous direction

to not oscillate too much

where would I get using the weights from previous iteration.

**Silva & Amida:**

$$\text{minimizing } C \cdot w_i^2 + h_1 w_1 + h_2$$

accelerate / decelerate

nothing changed

derivatives changed

$$\alpha_i^{(h+1)} = \begin{cases} \alpha_i^{(h)} u & \text{if } \nabla E_i^{(h)} > 0 \\ \alpha_i^{(h)} d & \text{if } \nabla E_i^{(h)} < 0 \\ \alpha_i^{(h)} & \text{else} \end{cases}$$

$$d < 1 < u$$

**Delta-Bar-Delta**

$$\alpha_i^{(h+1)} = \begin{cases} \alpha_i^{(h)} + u & : \nabla_i E^{(h)} \cdot \bar{\sigma}^{(h-1)} > 0 \\ \alpha_i^{(h)} - d & : \nabla_i E^{(h)} \cdot \bar{\sigma}^{(h-1)} < 0 \\ \alpha_i^{(h)} & \text{else} \end{cases}$$

$$\bar{\sigma}_i^h = (1 - \phi) \nabla_i E^L + \phi \bar{\sigma}_{i-1}^h$$

where  $\bar{\sigma}_i^h$

## (Super) SAB

Def:  $\alpha^+$  increasing ( $= 1, 0.5$ )  
 $\alpha^-$  decreasing ( $= 2$ )  
 $\alpha_{\text{START}}$  initial ( $= 1, 2$ )  
 $\alpha_m$  momentum ( $= 0.3$ )

If error sign didn't change, increase:  $\alpha_{ij}(t+1) = \alpha^+ \cdot \alpha_{ij}(t)$

else:

$$\alpha_{ij}(t+1) = \alpha_{ij}(t) / \alpha^-$$

$$\Delta w_{ij}(t+1) = 0$$

## AdaGrad

$$A_i = \frac{\partial E}{\partial w_i} \quad \rightarrow \text{Vlastní mi drž RMS shadem JE}$$

$$A_i(t+1) = A_i(t) + \left( \frac{\partial E}{\partial w_i} \right)^2$$

$$w_i(t+1) = w_i(t) - \frac{\alpha}{\sqrt{A_i}} \cdot \frac{\partial E}{\partial w_i}$$

tohle zůstává stejné

## Adam

$$N_i(t+1) = \beta_1 N_i(t) + (1-\beta_1) \cdot \frac{\partial E}{\partial w_i} \quad \rightarrow \text{first order}$$

$$A_i(t+1) = \beta_2 N_i(t) + (1-\beta_2) \left( \frac{\partial E}{\partial w_i} \right)^2 \quad \rightarrow \text{second order}$$

$$\tilde{N}_i(t+1) = N_i(t+1) / (1 - \beta_1^{t+1})$$

$$\tilde{A}_i(t+1) = A_i(t+1) / (1 - \beta_2^{t+1})$$

$$w_i(t+1) = w_i(t) - \alpha \frac{\tilde{N}_i(t+1)}{\sqrt{\tilde{A}_i(t+1)} + c}$$

- scale invariant

## Second-order-algorithms

$$\nabla E(w) = \begin{pmatrix} \vdots \\ \frac{\partial E(w)}{\partial w_i} \\ \vdots \end{pmatrix} \quad \nabla^2 E(w) = \begin{pmatrix} \vdots & & & \\ \cdots & \frac{\partial^2 E(w)}{\partial w_i \partial w_j} & \cdots & \vdots \\ \vdots & & & \vdots \end{pmatrix}$$

$$w_{(k+1)} = w_{(k)} - (\nabla^2 E(w))^{-1} \cdot \nabla E(w)$$

## Pseudo-Newton method

- only diagonal is computed

$$w_{(k+1)} = w_{(k)} - \frac{\nabla E(w)}{\frac{\partial^2 E(w)}{\partial w_i^2}}$$

## Quicprop

$$w_i(k+1) = w_i(k) + \nabla_{w_i} w_i$$

$$\nabla_i E(w)$$

$$\nabla_{w_i} w_i = \nabla_{w_{i-1}} w_i \cdot \frac{\nabla_i E(w) - \nabla_i E(w_{-1})}{\nabla_i E(w)}$$

$$\nabla_i E(w)$$

$$= \frac{\nabla_i E(w) - \nabla_i E(w_{-1})}{\nabla_{w_{i-1}} w_i}$$

$$\left. \begin{array}{l} \text{tohle je} \\ \text{odhad} \\ \frac{\partial^2 E(w)}{\partial w^2} \end{array} \right\}$$

## Maxout activation

$$\max \{ x_i^\top w_1, x_i^\top w_2 \}$$

$$\nabla w_i = -\alpha \frac{E(w') - E(w)}{\beta}$$

## Weight perturbation

## Batch normalization

- upravují výběr fenomén

↪ hodnota všechno.

$$\alpha_i = \gamma_i \frac{y_i - \text{mean}}{\sqrt{\text{var}}} + \beta_i \quad \begin{matrix} \nearrow \text{mean of the} \\ \nearrow \text{lin. output} \end{matrix}$$

$L^2$ -reg.

$$L = \sum_{(y, t)} (t - y)^2 + \lambda \cdot \sum_i w_i^2$$

$$w_i = w_i - \alpha \lambda s_i - \alpha \frac{\partial L}{\partial w_i}$$

$$s_i = \begin{cases} -1 & w_i < 0 \\ +1 & w_i > 0 \end{cases}$$

## Space partitioning

$R(m, n) := \# \text{ of regions defined}$

by  $m$  hyperplanes in  
 $n$ -dimensional space.

then:  $\# \text{ threshold func.} \leq \# \text{ regions} \leq \# \text{ boolean func.}$

$$R(1, n) = 2, R(0, m) = 0 : n, m \geq 1$$

$$R(m, n) = R(m-1, n) + R(m-1, n-1)$$

$$= 2 \cdot \sum_{i=0}^{n-1} \binom{m-1}{i}$$

$\hookrightarrow$  if the dimension of vector  
is too big, we don't have suitable  
threshold func. for good approx.

## VC-dimension

"když užijete množství funkcií rovněž vystupují"

## Minimum number of patterns

$$P \geq \left( \frac{W}{\varepsilon} \right) \log \left( \frac{N}{\varepsilon} \right)$$

$N$  ... neurons

$W$  ... weight

$P$  ... patterns

### Definition:

Let  $C = \{f_i\}$  be a set of functions (concept class).

The set of  $m$  training patterns  $\{t_k\}_{k=1,\dots,m}$  can be shattered by means of  $C$ , if for each of the  $2^m$  possible labelings of these patterns with 1/0, there exists at least one function, that satisfies this labeling.

### Definition:

The VC-dimension  $V$  of a set of functions  $C$  is defined as the biggest  $m$ , for which a set of  $m$  training patterns exists that can be shattered.

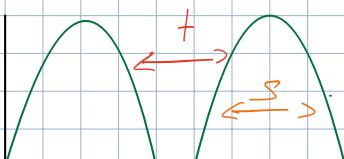
Předpoklad možnosti  $C$  obdržet funkce pro lib.  $m$ ,  
je také problém nenujatelný.

## The optimal error for enforcing condensed features

$$F = y^s (1-y)^s \cdot (y - 0.5)^+$$

Odm. výšší  $t$ , tím větší dim. uprostřed

Odm. výšší  $s$ , tím větší rozložení hory



## Pruning BP

redundant layer:

- no uniform neuron  $\times$

- no identical neuron repr.  $\times$

- no inverse neuron repr.  $\times$

$\rightarrow$  this is reduced.  $\rightarrow$  there is one  
for each BP network