— Důležitost, důkazy a poznámky

• **Neuron jako lineární klasifikátor**

$\xi < 0$ | $\xi = 0$ | $\xi > 0$



$\left| \dfrac{-w_0}{\|\vec{w}\|} \right|$

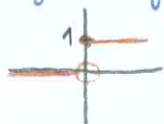• přímka ($\xi = 0$) je definovaná vektorem $\vec{w}$ (váhy neuronu)
• $x_1, x_2, x_3$ (směr, kam ukazuje vektor $\vec{w}$) mají $\xi > 0$; a $x_4, x_5, x_6$ mají $\xi < 0$
• posunutí přímky (lineárního klasifikátoru) je definováno $\left| \dfrac{-w_0}{\|\vec{w}\|} \right|$, kde $w_0$ je bias

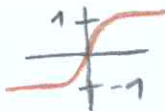• **Aktivační funkce**

 – Unit step function : $\sigma(\xi) = \begin{cases} 1 & \xi \geq 0 \\ 0 & \xi < 0 \end{cases}$

 

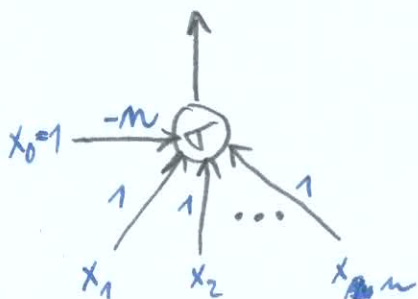 – Logistická sigmoida : $\sigma(\xi) = \dfrac{1}{1 + e^{-\lambda \cdot \xi}}$ ; $\lambda \in \mathbb{R}$ parametr strmosti

 

 – Hyperbolic tangens : $\sigma(\xi) = \dfrac{1 - e^{-\xi}}{1 + e^{-\xi}}$
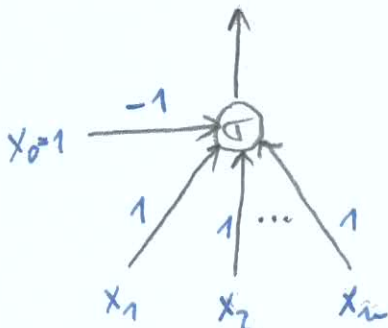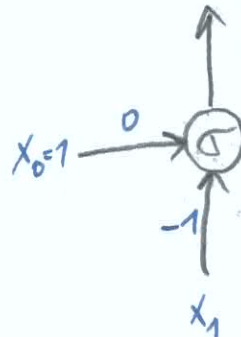
 

• **Logické funkce** (unit step function)

AND | OR | NOT

**ADALINE**

- Error function

$$E(\vec{w}) = \frac{1}{2}\sum_{k=1}^{p}\left(\vec{w}\cdot\vec{x}_k - d_k\right)^2 = \frac{1}{2}\sum_{k=1}^{p}\left(\sum_{i=0}^{n} w_i x_{ki} - d_k\right)^2$$

$\Rightarrow$ cíl: najít $\vec{w}$ takové, které minimalizuje $E(\vec{w})$

- **Gradient of the error function**

$$\nabla E(\vec{w}) = \left(\frac{\partial E}{\partial w_0}(\vec{w}), \ldots, \frac{\partial E}{\partial w_n}(\vec{w})\right)$$

- vektor v prostoru vah, který míří ve směru nejstrmějšího svahu chybové fce.
- vektory $\vec{x}_k$ jsou fixované $\leftarrow$ vstupy

$\Rightarrow \nabla E(\vec{w}) = \vec{0} = (0, \ldots, 0)$ znamená globální minimum chybové fce E

- 
$$\frac{\partial E}{\partial w_\ell}(\vec{w}) = \frac{1}{2}\sum_{k=1}^{p}\frac{\delta}{\delta w_\ell}\left(\sum_{i=0}^{n} w_i x_{ki} - d_k\right)^2$$

$$= \sum_{k=1}^{p}\left(\vec{w}\cdot\vec{x}_k - d_k\right) x_{k\ell}$$

$$\Rightarrow \boxed{\nabla E(\vec{w}) = \left(\frac{\partial E}{\partial w_0}(\vec{w}), \ldots, \frac{\partial E}{\partial w_n}(\vec{w})\right) = \sum_{k=1}^{p}\left(\vec{w}\cdot\vec{x}_k - d_k\right)\vec{x}_k}$$

- Learning algorithm

(Batch): $\boxed{\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon\cdot\nabla E(\vec{w}^{(t)}) = \vec{w}^{(t)} - \varepsilon\cdot\sum_{k=1}^{p}\left(\vec{w}^{(t)}\cdot\vec{x}_k - d_k\right)\cdot\vec{x}_k}$

(Online): 
$$\boxed{\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t)\cdot\left(\vec{w}^{(t)}\cdot\vec{x}_k - d_k\right)\cdot\vec{x}_k}$$

(2)

# MULTI LAYER PERCEPTRON

- **Notace :**
  - $X$ ..... set vstupních neuronů (inputs)
  - $Y$ ..... set výstupních neuronů (outputs)
  - $Z$ ..... set všech neuronů ($X, Y \subseteq Z$)
  - $i, j, ...$ indexy neuronů :
    - $\xi_j$ ..... vnitřní potenciál neuronu $j$ $\qquad \left( \sum\limits_{i \in j^{\leftarrow}} w_{ji} \, y_i \right)$
    - $y_j$ ..... output neuronu $j$ $\left( \sigma(\xi_j) \right)$
  - $w_{ji}$ ..... váha z neuronu $i$ do neuronu $j$
  - $j^{\leftarrow}$ ..... set všech neuronů, které míří do $j$
  - $j^{\rightarrow}$ ..... set všech neuronů, na které míří $j$

- **Error function**

$$\left\{ (\vec{x}_\ell, \vec{d}_\ell) \mid \ell = 1, \dots, p \right\} - \text{Tréninková sada}$$

$$E(\vec{w}) = \sum_{\ell=1}^{p} E_\ell(\vec{w}) \quad \dots\dots \text{ suma přes všechny tréninkové příklady}$$

$$E_\ell(\vec{w}) = \frac{1}{2} \sum_{j \in Y} \left( y_j(\vec{w}, \vec{x}_\ell) - d_{\ell j} \right)^2$$

- **Learning algorithm (Batch)**

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

komponenta gradientu $\nabla E \quad \Rightarrow$

$$\Rightarrow \vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t) \cdot \nabla E(\vec{w}^{(t)})$$

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E}{\partial w_{ji}} (\vec{w}^{(t)}) \quad \dots\dots \text{ update } w_{ji} \text{ v } t+1 \text{ kroku}$$

- **Error function gradient**

  - pro každé $w_{ji}$: $\qquad \dfrac{\partial E}{\partial w_{ji}} = \sum\limits_{\ell=1}^{p} \dfrac{\partial E_\ell}{\partial w_{ji}}$
    (váhy)

  - a pro každé $j \in Z \setminus X$ :
    (neurony)

  $$\frac{\partial E_\ell}{\partial y_j} = y_j - d_{\ell j} \qquad \text{pro } j \in Y$$

  - zde pro každé $\ell$ :
    (examples)

  $$\frac{\partial E_\ell}{\partial w_{ji}} = \frac{\partial E_\ell}{\partial y_j} \cdot \sigma_j'(\xi_j) \cdot y_i$$

  $$\frac{\partial E_\ell}{\partial y_j} = \sum_{n \in j^{\rightarrow}} \frac{\partial E_\ell}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot w_{nj}$$

  pro $j \in Z \setminus (Y \cup X)$

③

... slovy :

- Změna chyby podle dané váhy je rovna součtu chyb všech tréninkových příkladů (podle dané váhy).
- Změna chyby příkladu (example) podle váhy je rovna změně chyby podle výstupu daného neuronu (do kterého váha vchází), vynásobení vnitřním potenciálem, na který je aplikována derivovaná aktivační fce, a výstupem neuronu, ~~po které váha o tesy~~ ze kterého daná daná váha vychází (o úroveň ~~níž~~ níž)
- Změna chyby podle výstupu je rovna :
  - pro výstupní neurony : Rozdíl reálného a očekávaného výstupu
  - pro ostatní neurony : • suma přes všechny neurony $\hat{r}$, do kterých daný neur. $\check{j}$ míří (o 1 vrstvu výš) :
    - • Změna chyby podle výstupu daného neuronu o vrstvu ~~níž~~ výš, krát váha mezi těmito 2 neurony, krát vnitřní potenciál neuronu r (vrchní vrstva), na který je aplikována derivovaná aplikační fce

- <u>Derivace aktivačních funkcí</u>

  Logistická sigmoida :  $\sigma'_j(\xi_j) = \lambda_j y_j (1 - y_j)$

  Hyperbolický tangens :  $\sigma_j(\xi_j) = a \cdot \tanh(b \cdot \xi_j)$   • modifikující tvar (strmost, roztočení) křivky

  $\sigma'_j(\xi_j) = \frac{b}{a}(a - y_j)(a + y_j)$

- <u>Gradient Descent algoritmus</u>

  $\mathcal{E}_{ji} = 0$   (na konci $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

  for $(\ell = 1, \dots, p)$ :

  1) <u>forward pass</u> : - spočítej $y_j = y_j(\vec{w}; \vec{x}_\ell)$ pro ~~ka~~ všechny $j \in Z$
  2) <u>backward pass</u> : - spočítej $\frac{\partial E_\ell}{\partial y_j}$ pro všechny $j \in Z$ pomocí [Backpropagation]
  3) spočítej $\frac{\partial E_\ell}{\partial w_{ji}} = \frac{\partial E_\ell}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$   pro všechny $w_{ji}$
  4) $\mathcal{E}_{ji} = \mathcal{E}_{ji} + \frac{\partial E_\ell}{\partial w_{ji}}$

  Výsledná $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$

- **Backpropagation**

  Spočítej $\dfrac{\partial E_k}{\partial y_j}$ pro všechna $j \in Z$ :

  - $j \in Y \Rightarrow \dfrac{\partial E_k}{\partial y_j} = y_j - d_{kj}$

  - $j \in Z \setminus (Y \cup X) \Rightarrow$ předpokl., že $\dfrac{\partial E_k}{\partial y_n}$ již známe:

  $$\frac{\partial E_k}{\partial y_j} = \sum_{n \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot w_{nj}$$

- **Chain rule**

  $$\frac{\partial E_k}{\partial w_{12}} = \frac{\partial E_k}{\partial y_1} \cdot \frac{\partial y_1}{\partial \xi_1} \cdot \frac{\partial \xi_1}{\partial w_{12}} = \frac{\partial E_k}{\partial y_1} \cdot \sigma_1'(\xi_1) \cdot y_2$$

  $$\frac{\partial E_k}{\partial w_{23}} = \ldots = \frac{\partial E_k}{\partial y_2} \cdot \sigma_2'(\xi_2) \cdot y_3$$

  $$\frac{\partial E_k}{\partial w_{34}} = \ldots = \frac{\partial E_k}{\partial y_3} \cdot \sigma_3'(\xi_3) \cdot y_4$$
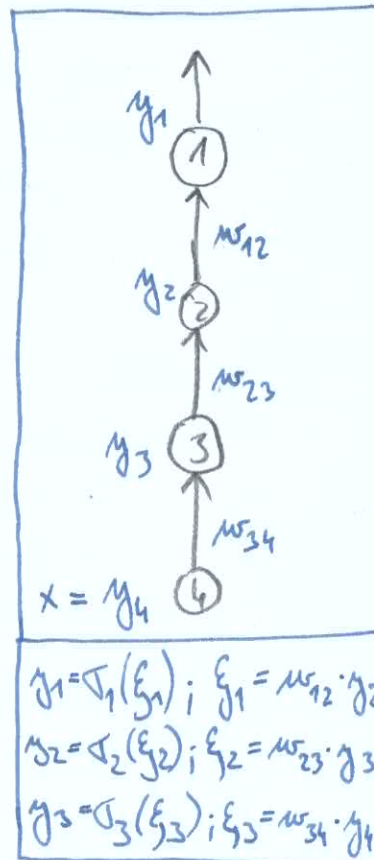
- **Error functions**

  - Square error: $E(\vec{w}) = \displaystyle\sum_{k=1}^{p} E_k(\vec{w})$ , kde

  $$E_k(\vec{w}) = \frac{1}{2} \sum_{j \in Y} \left( y_j(\vec{w}, \vec{x}_k) - d_{kj} \right)^2$$

  - Mean square error (mse):

  $$E(\vec{w}) = \frac{1}{p} \sum_{k=1}^{p} E_k(\vec{w}) \quad \ldots\ldots \text{ průměrná chyba objektu v daném}$$
  $$\text{batchi}$$



$y_1 = \sigma_1(\xi_1) ; \ \xi_1 = w_{12} \cdot y_2$
$y_2 = \sigma_2(\xi_2) ; \ \xi_2 = w_{23} \cdot y_3$
$y_3 = \sigma_3(\xi_3) ; \ \xi_3 = w_{34} \cdot y_4$

- **Momentum**

  Problém: Moc velké kroky mohou po sebelehčí údolí (správným směrem) vyběhnout na druhé straně opět nahoru. ~~...~~

  Řešení: V každém kroku přičteme směr provedený minulým krokem (vynásobený faktorem $\alpha$) ⟶ **HYBNOST**

  $$\Delta \vec{w}(\Delta) = -\varepsilon(\Delta) \cdot \sum_{\mathcal{R} \in T} \nabla E_{\mathcal{R}}(\vec{w}^{(\Delta)}) + \boxed{\alpha \cdot \Delta w_{j}^{(\Delta-1)}}$$

- **Adaptující se rychlost učení (learning rate)**

  - začnout s vyšším (např. $\varepsilon = 0.1$)
  - později snižovat (dohledávka lokálního minima)

  ~~...~~ Případně:

  - ~~...~~ Error ↓ ⟹ $\varepsilon$ ↑        (chyba se může krátkodobě zvedat a nám
  - Error ↑ ⟹ $\varepsilon$ ↓          to nemusí vadit - minimum je schované
                                                na konci)

- **Ada Grad**

  - každá váha má svůj learning rate ($\varepsilon$), který se updatuje samostatně:
    - málo měnící se váhy ⟹ ~~velký~~ velké $\varepsilon$      } všechny váhy se
    - hodně měnící se váhy ⟹ malé $\varepsilon$          } mění ± stejně
  - problém: akumulace během celého procesu ⟹ RMSProp minimalizuje historii
    <br>          ↑ updatů                                              (exponenciálně)

1) diferencovatelná .... většinou (ne nutně vždy - ReLU) ⟶ Gradient Descent
<br>                                                         počítá s derivací

~~2) non-linea~~

2) nelineární .... i více - vrstevná síť (lineární) je ekvivalentní jedno - vrstevné

3) monotónnost ..... bez lokálních extrémů

4) "lineární" ..... jednodušeji se učí ⟶ tradeoff
<br>               ..... dnes se používají co nejvíce lineární

- **Input preprocessing**
  - velké hodnoty mají výrazně větší vliv na trénink než malé
  => Standardizace : • průměr = 0 (odečtení průměru)
                     • odchylka = 1 (podělení standardní odchylkou)
                       (variance)                    "standard deviation"
  - redukce dimenze : PCA (Principal Component Analysis)

- **Inicializace vah**   interval $[-w; w]$; průměr = 0; odchylka = 1; $d$ = počet vstupů
                                              (rozptyl)        do neuronu
  - příliš malá čísla => příliš lineární
  - příliš velká čísla => plochá oblasti (rozseknuté)

  <u>Řešení:</u> vybrat $w$ tokové, aby standardní odchylka je blízko přelomu mezi
  lineární a nelineární části fce. ($\xi$ se herz rozptýlí do ~~tohoto~~
  tohoto intervalu)

  $$\Rightarrow \boxed{w = \frac{\sqrt{3}}{\sqrt{d}}} \Rightarrow \left[ -\frac{\sqrt{3}}{\sqrt{d}} ; \frac{\sqrt{3}}{\sqrt{d}} \right]$$   herz rozptýlí vnitřní potenciál do intervalu $[-1; 1]$

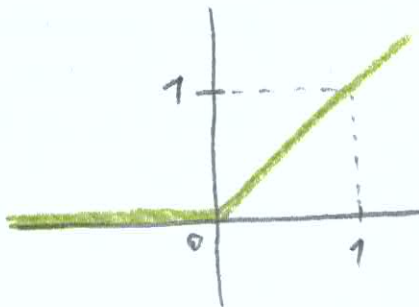- **Glorot & Bengio inicializace**
  - pro lineární modely ideální
  - lépe zohledňuje forward pass (řeší i počet neuronů nahoře, ne jen dole)
  $$\boxed{w = \sqrt{\frac{6}{m+n}}}$$

- **Moderní aktivační funkce**
  - Hidden layers : <u>ReLU</u>
  $$\sigma(\xi_j) = \max\{0, \xi_j\}$$

  

  - Output neurons:
    regrese => lineární / sigmoida
    klasifikace => binární -> logistická sigmoida
                           => tanh
                => n - ární => <u>Softmax</u>
                   (multi-class)
    $$\boxed{\sigma_j(\xi_j) = \frac{e^{\xi_j}}{\sum_{i \in Y} e^{\xi_{ki}}}}$$

    ... spolu se softmaxem se nejčastěji používá <u>Cross-entropy</u>
        jako error function

⑦

- **Cross-entropy**

$$E = -\frac{1}{l} \sum_{k=1}^{l} \sum_{j \in Y} \left[ d_{kj} \ln(y_j) + (1 - d_{kj}) \ln(1 - y_j) \right]$$

.... minimalizováním Cross-entropy maximalizujeme likelihood

"Maximum likelihood principle"

- **Proč je Cross-entropy lepší než MSE (mean square error)?**

Slovně: • Když se MSE plete, zůstane relativně zaseklá ve špatném rozhodnutí ⇒ přeúčení je pomalý.

• Cross-entropy dělá sama od sebe velké kroky směrem od špatného rozhodnutí.

Formálně: binární klasifikace - $\{0;1\}$; 1 output neuron → logistická sigmoida

$\sigma' = \sigma(1-\sigma)$



**MSE :** $E_k^{mse} = \left( \sigma(w x_k) - d_k \right)^2$

$$\frac{\partial E^{mse}}{\partial w} = \left( \sigma(w x_k) - d_k \right) \cdot \sigma'(w x_k) \cdot x_k =$$

$$= \left( \sigma(w x_k) - d_k \right) \cdot \boxed{\sigma(w x_k)} \left( 1 - \sigma(w x_k) \right) \cdot x_k$$

$\cdot d_k = 1$ :

→ odečítáme tuto hodnotu od stávajících vah

• pokud jsme ~~ověření špatě~~ ~~jijí 0~~ je ovšem vnitřní potenciál záporný (a my jej chceme dostat do kladných čísel), ~~~~

$\frac{\partial E^{mse}}{\partial w}$ je malé číslo. Člen ● ($w x_k$ je vnitřní potenciál ⇒ $\sigma(w x_k)$ je blízko 0).

Moc se nepohneme ⇒ Z tohoto špatného stavu se nedostaneme.

**CROSS-ENTROPY :**

$$E_k^{cross}(w) = -d_k \ln(\sigma(w x_k)) + (1 - d_k) \underbrace{\ln(1 - \sigma(w x_k))}_{=0}$$

$d_k = 1$ : $\dfrac{\partial E_k^{cross}}{\partial w} = -\dfrac{1}{\sigma(w x_k)} \cdot \sigma'(w x_k) \cdot x_k = -\dfrac{1}{\sigma(w x_k)} \cdot \sigma(w x_k)(1 - \sigma(w x_k)) \cdot x_k =$

velké záporné číslo → ⇒ odečítám od vah → posunene "doprava" → ⇒ zlepšíme se

$$= -\left( 1 - \sigma(w x_k) \right) \cdot x_k = \boxed{\left( \sigma(w x_k) - d_k \right) \cdot x_k}$$

(pro $d_k = 0$ získáme to stejné)

⑧

- ~~Cross~~ ~~validation~~ Early stopping (kdy skončit s učením)

    - rozdělit dataset na 3 části:
        • training set (60%) – k trénování použít pouze tyto data
        • validation (20%) – ke zjištění, kdy zastavit učení
        • test set (20%) – k otestování sítě (porovnání s jinými modely, ...)
    → validation set se taky „všímá" ⇒ na opravdové testování musíme
    ◊ VŽDY ◊ použít čerstvý, nedotknutý dataset

- Dropout
    - v každém kroku gradient descentu vybereme náhodně jen podmnožinu neuronů, které se použijí do učení. Zbytek se pro toto kolo ignoruje
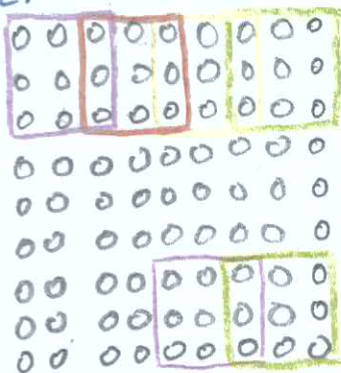
- Weight decay
    - při každém učeném kole, vynásobíme koeficientem $0 < \alpha < 1$, a tím snížíme všechny váhy
    - silné (důležité) přežijí a slabé (nedůležité) spadnou na 0

$$\boxed{\text{KONVOLUČNÍ SÍTĚ}}$$

- Zpracování obrazu používá téměř výlučně konvoluční sítě                    $k \times k$
- každý neuron ve skryté vrstvě (hidden layer) je napojen na čtverec vstupních neuronů. Toto okno se s každým dalším neuronem ve skryté vrstvě posouvá
  (o hodnotu ~~shift~~) a tvoří tzv. feature mapy.

stride    stride
~~shift~~ = 2:
k = 3:



feature map :
    - všechny neurony se v jedné mapě spolu sdílí váhy
    - mají pouze $k \times k$ (9) vah
    - každá vrstva může mít několik takovýchto feature map
    - každá feature mapa „hledá" v obrázku nějaký vzor"

vstupní obrázek
(~~pixely~~) (např. pixely)

Pooling
    - redukce dimenze (každý neuron v pooling vrstvě je opět napojen
    - nad feature mapou     na $n \times n$ čtverec ve feature mapě a
    → Max-pooling: bere maximum (abstrahuje je do 1 výstupu)
    → L2-pooling: odmocnina ze součtu 2. mocnin
    → Avg-pooling: ~~va~~ průměr
       ⋮
                - nepřekrývají se (narozdíl od konvol. vrstvy)

(9)

- Příklad



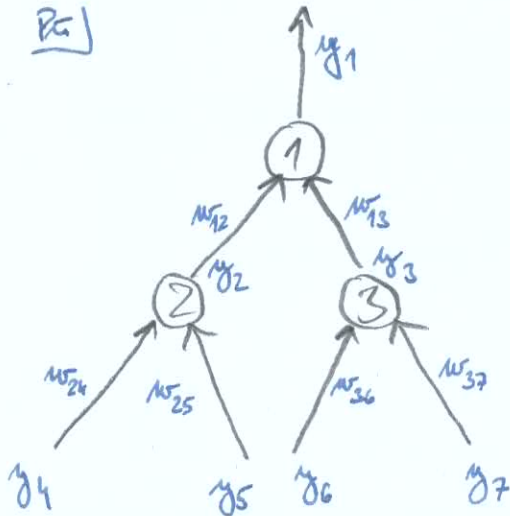| Vstupní obrázek (input) | → | 3 feature mapy | → | Max-pooling vrstva pro každou feat. m. | ⟹ | 10 výstupních neuronů |
|---|---|---|---|---|---|---|
| $28 \times 28$ | | $3 \times 24 \times 24$ | | $3 \times 12 \times 12$ | | |

→ úplné propojení: každý výstupní neuron je napojen na úplně všechny neurony z pooling vrstev

⟹ Gradient descent pro trénink

- Jak vypočítat $\dfrac{\partial E_k}{\partial w_{ji}}$ ?

Př.]



$$w_{24} = w_{36} = w_a$$
$$w_{25} = w_{37} = w_b$$

$$y_1 = \sigma_1(w_{12}\, y_2 + w_{13}\, y_3)$$
$$y_2 = \sigma_2(w_{24}\, y_4 + w_{25}\, y_5)$$
$$y_3 = \sigma_3(w_{36}\, y_6 + w_{37}\, y_7)$$

1 feature mapa:
②,③

1 input vrstva:
$y_4, y_5, y_6, y_7$ (strive 2)

- polimetí s $w_a$ má vliv na $y_2$ i $y_3$
⟹ počítá se přes všechny neurony, které to ovlivní (které sdílí váhy)

$$\frac{\partial E_k}{\partial w_a} = \frac{\partial E_k}{\partial y_2} \cdot \frac{\partial y_2}{\partial w_a} + \frac{\partial E_k}{\partial y_3} \cdot \frac{\partial y_3}{\partial w_a} =$$

$$= \frac{\partial E_k}{\partial y_2} \cdot \sigma_2' \cdot y_4 + \frac{\partial E_k}{\partial y_3} \cdot \sigma_3' \cdot y_6$$

$$\frac{\partial E_k}{\partial w_{ji}} = \sum_{n\in j\ shore} \frac{\partial E_k}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot y_n$$

Opravdu tu má být $y_n$ ?
- ve slidech to tak je (str. 189), ale příklad v Schoboule ukazuje něco jiného

- Backpropagation

- pro $j\in Y$: $\dfrac{\partial E_k}{\partial y_j} = y_j - d_{kj}$ (pro MSE)

- pro $j\in Z-Y$, kde $j^{\rightarrow}$ je konvoluční vrstva nebo dense vrstva: (úplně propojená)

$$\frac{\partial E_k}{\partial y_j} = \sum_{n\in j^{\rightarrow}} \frac{\partial E_k}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot w_{nj}$$

- pro $j \in Z \setminus Y$, kde $j^{\rightarrow}$ je max-pooling vrstva, pak $j^{\rightarrow} = \{i\}$ pro každý 'max' neuron:

$$\frac{\partial E_k}{\partial y_j} = \begin{cases} \dfrac{\partial E_k}{\partial y_{ji}} & \text{pokud } j \text{ je ten maximální neuron} \\ \\ 0 & \text{jinak} \end{cases}$$

## DEEP MLP

PRO: - lepší výsledky
- těžší problémy
- 1 vrstva může být neefektivní

PROTI: - vanishing gradient
- rychlé přeučení (overfit)

· <u>Vanishing gradient</u>

$$\frac{\partial E_k}{\partial y_j} = \sum_{M \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_M} \cdot \boxed{\sigma'_M(\xi_M)} \cdot w_{Mj}$$

pro standardní aktivační fce (logistická sigmoida tanh,...) je menší než 1 ⟹ opakovaným násobením gradient stále snižujeme ⟹ ⟹ v nižších vrstvách to není vůbec patrný

· <u>Předtrénování MLP pomocí RBM</u>

- každé 2 sousední vrstvy MLP můžeme považovat za RBM (Restricted Boltzman Machine) ..... $B_i$ odpovídá vrstvě $i$ a $i-1$ ($B_1 = $ input + 1. skrytá)

- 1. fáze učení (unsupervised pretraining):        $L = $ počet vrstev
  · bereme jen vstupy (ignorujeme očekávané výstupy)
  · na náhodné podmnožině trénivkové sady natrénujeme $B_i$ (pro $i = 1, ..., L$) pomocí RBM algoritmu
  · tímto vlastně inicializujeme váhy MLP tak, že již v sobě obsahují jistou dávku informace (vzorů) z dat. - „nejdříve trochu podgříme doménu"
  · postupujeme od spodu - každou předtrénovanou podčást sítě (1. vrstva, 1. dvě vrstvy,...) použijeme k transformaci dat pro následující vrstvu (2., 3.,...)

# REKURENTNÍ SÍTĚ

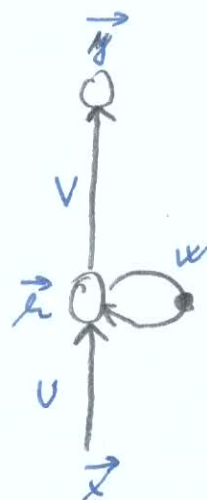- **Rekurentní MLP** (se smyčkou):

    - vstupy (input): $\vec{x}$
    - stav: $\vec{h}$
    - output: $\vec{y}$
    - matice: $U, W, V$
      váhy

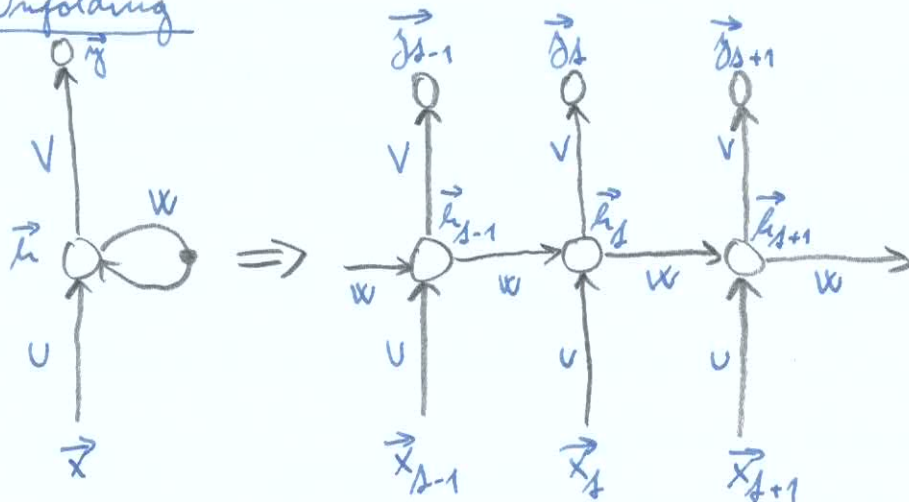$$\vec{h} = \sigma(U\vec{x} + W\vec{h})$$ .... sigmoidal /ReLU

z minulého kroku – staré hodnoty

$$\vec{y} = \sigma(V\vec{h})$$ .......... softmax / sigmoidal

=> ke zpracování sekvence vstupů => časové řady atd. (zpracování řeči)

- **Unfolding**



+1 dimenze
vstup = sekvence vektorů
výstup = sekvence vektorů

$\vec{h}_t$ je paměť sítě

$$\vec{h}_t = \sigma(U\vec{x}_t + W\vec{h}_{t-1})$$
$$\vec{y}_t = \sigma(V\vec{h}_t)$$

matice $U, V, W$ jsou pro všechny stejné

- **Learning**

    - training example (1 případ): $(X, d)$
    
    $$X = \vec{x}_1, \ldots, \vec{x}_T \quad \text{(sekvence vstupů)}$$
    
    $$\vec{x}_t = (x_{t1}, \ldots, x_{tN})$$
    
    (to stejné pro $d$ - výstupy)
    
    - rozbalení (unfolding) pro dané $X$ vytvoří sekvenci vnitřních stavů:
    
    $$\vec{h}_0, \vec{h}_1, \vec{h}_2, \ldots, \vec{h}_T \quad \text{kde } \vec{h}_t = (h_{t1}, \ldots, h_{tH})$$
    
    $$\vec{h}_0 = (0, \ldots, 0)$$
    
    - sekvence výstupů: $\vec{y}_1, \ldots, \vec{y}_T \quad \text{kde } \vec{y}_t = (y_{t1}, \ldots, y_{tM})$
    
    Chybová fce:
    
    $$E(x, d) = \sum_{t=1}^{T} \sum_{k=1}^{M} (y_{tj} - d_{tj})^2$$
    
    → online verze (při minibatchi můžeme přidat ještě jeden stupeň)

- **Backpropagation**

$\nabla E_{(x,d)}\left(\vec{w}^{(t)}\right)$ ..... vektor všech parciálních derivací $\dfrac{\partial E_{(x,d)}}{\partial w_{ji}}$

1) jak nahradit derivace podle $w_{ji}$ na derivace podle $y_{ji}$:

$$\frac{\partial E_{(x,d)}}{\partial w_{ji}} = \sum_{n \in j^{\rightarrow}} \frac{\partial E_{(x,d)}}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot y_i$$

2) pro $j \in Y$: (MSE)

$$\frac{\partial E_{(x,d)}}{\partial y_j} = y_j - d_j$$

3) pro $j \in Z \setminus Y$:

$$\frac{\partial E_{(x,d)}}{\partial y_j} = \sum_{n \in j^{\rightarrow}} \frac{\partial E_{(x,d)}}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot w_{nj}$$

RNN:

$$\cdot \ \frac{\partial E_{(x,d)}}{\partial h_{sk}} = \sum_{k'=1}^{N} \frac{\partial E_{(x,d)}}{\partial y_{sk'}} \cdot \sigma' \cdot V_{kk'} + \sum_{k'=1}^{H} \frac{\partial E_{(x,d)}}{\partial h_{(s+1)k'}} \cdot \boxed{\sigma' \cdot W_{kk'}}$$

$$\cdot \ \frac{\partial E_{(x,d)}}{\partial V_{kk'}} = \sum_{s=1}^{T} \frac{\partial E_{(x,d)}}{\partial y_{sk}} \cdot \sigma' \cdot h_{sk}$$

$$\cdot \ \frac{\partial E_{(x,d)}}{\partial W_{kk'}} = \sum_{s=1}^{T} \frac{\partial E_{(x,d)}}{\partial h_{sk}} \cdot \sigma' \cdot h_{(s-1)k}$$

$$\cdot \ \frac{\partial E_{(x,d)}}{\partial U_{kk'}} = \sum_{s=1}^{T} \frac{\partial E_{(x,d)}}{\partial h_{sk}} \cdot \sigma' \cdot x_{sk}$$

- špatly v čase = vanishing gradient nebo exploding

$\Rightarrow$ řešení: LTSM

chceme co nejblíže 1, jinak
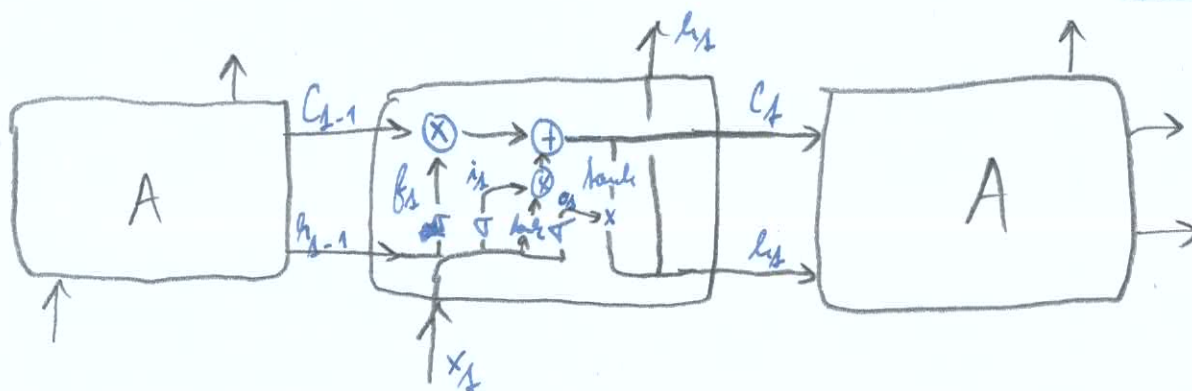
− každý neuron je nahrazen komplexní 'mašinou'

○...... hadamardův product (násoben po komponentách, ne dot product)

~~$(\vec{y}_t =) \vec{h}_t = \vec{\sigma}_t \odot \vec{c}_t \odot \vec{x}_t$~~

• $\underline{RNN}$ vs $\underline{LSTM}$
(tanh)



$\boxed{RNN}$

$\boxed{LSTM}$

$$(\vec{y}_t =) \; \vec{h}_t = \vec{\sigma}_t \odot \sigma_h(\vec{c}_t) \qquad \text{output}$$

$$\vec{c}_t = \vec{f}_t \odot \vec{c}_{t-1} + \vec{i}_t \odot \tilde{c}_t \qquad \text{memory}$$

$$\tilde{c}_t = \sigma_c(W_c \cdot \vec{h}_{t-1} + U_c \cdot x_t) \qquad \text{new memory contexts}$$

$$\vec{\sigma}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t) \qquad \text{output gate}$$

$$\vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t) \qquad \text{forget gate}$$

$$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t) \qquad \text{input gate}$$

− kontroluje, kolik informace má převzít z historie, kolik má zahodit, kolik má vzít z nových dat, kolik poslat dál, atd.

- **Vector quantization**

  - pro každý input $\vec{x}$, určíme nejbližší střed $\vec{w}_{c(\vec{x})}$:

$$C(\vec{x}) = \arg\min_{i=1,\dots,h} \left\{ \| \vec{x} - \vec{w}_i \| \right\}$$

  - následně minimalizujeme chybu:

$$E = \int \| \vec{x} - \vec{w}_{c(\vec{x})} \|^2 \, f(\vec{x}) \, d\vec{x}$$

  $\rightarrow$ chceme nahradit všechy objekty $\vec{x}$ středy $\vec{w}$ tak, abychom snížili dimenzi a zároveň dobře popsali dataset
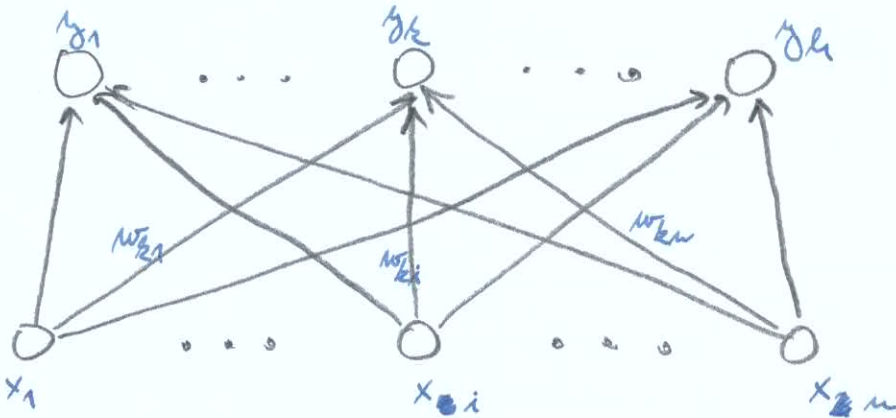
  $\rightarrow$ co nejvíc středů v hustých oblastech, co nejmíň v řídkých

  jak to vyřešit? $\Rightarrow$ k-means (Lloyd's algorithm)

  $\Rightarrow$ Kohonenovy mapy

| | ONLINE | |
|---|---|---|
| | NE :( | – potřebujeme celý dataset |
| | ANO :) | |

- **Z pohledu NN**

  1 vrstva:



vstupní objekty:
$x$
výsledná souřadnice středu:
$w$

neurony Kohonenovy mapy jsou spolu propojené, ale nezáleží na standardním propojování (váhy, ...)

pro vstup $\vec{x} \in \mathbb{R}^u$ a $k = 1, \dots, h$:

$$y_k \begin{cases} 1 & k = \arg\min_{i=1,\dots,h} \| \vec{x} - \vec{w}_i \| \\ 0 & \text{jinak} \end{cases}$$

Activity

(15)

- **Learning**
    - bereme v úvahu topologii sítě, která je „nad" klasickou strukturou
    - $\rightarrow$ propojení neuronů do mřížky

    $\underline{d(c, q)}$ = nejkratší cesta z $c$ do $q$

    topological neighbourhood neuronu $c$ o velikosti $s$:

$$N_s(c) = \{ q \mid d(c, q) \leq s \}$$

- **Krok učícího algoritmu**

$$\vec{w}_q^{(t)} = \begin{cases} \vec{w}_q^{(t-1)} + \theta \cdot \left( \vec{x}_t - \vec{w}_q^{(t-1)} \right) & \text{pro } q \in N_s(c(\vec{x}_t)) \\ \\ \vec{w}_q^{(t-1)} & \text{jinak} \end{cases}$$

kde $c(\vec{x}_t) = \arg \min\limits_{i = 1, \ldots, k} \left\| \vec{x}_t - \vec{w}_i^{(t-1)} \right\|$

$\left.\begin{array}{l} \theta \in \mathbb{R} \\ s \in \mathbb{N}_0 \end{array}\right\}$ parametry, které se během tréninku mohou měnit

- **V Praxi**

    2 fáze: 
    - **coarse phase**: 
        - cca 1000 kroků
        - learning rate $\theta$: od 0.1 do 0.01 (postupně)
        - sousedství neuronů: začínat s velkým okruhem, postupně zmenšovat
    - **fine tuning**
        - cca 500 × počet neuronů
        - $\theta$ blízko 0,01 (jinak je pravděpodobnost defektů)
        - jen malá sousedství neuronů

- **Klasifikace pomocí Kohonenových map**

    1) Natrénuj mapu na vstupních datech (ignoruj třídy)
    2) Označ neurony třídami: - každé třídě vyber jen neuron, který rozpoznává nejvíc z jejich instancí (obráceně)
    3) dotrénuj trénink (fine tune) za pomoci LVQ (později)

    Použití: - spočítej pro daný vstup $\vec{x}$ nejbližší střed a vrať jeho třídu

- **LVQ**

    pro každý trénindový příklad:

    - najdi nejbližší neuron:
        - pokud je ~~správná ta~~ třída správě, posuň střed blíže k novému examplu
        - pokud je třída špatně, posuň střed kousek od nového examplu

$$\vec{w_c}(t) = \begin{cases} \vec{w_c}(t-1) + \alpha(\vec{X_t} - \vec{w_c}(t-1)) & d_t = N_c \\[2em] \vec{w_c}(t-1) - \alpha(\vec{X_t} - \vec{w_c}(t-1)) & d_t \neq N_c \end{cases}$$

## HOPFIELD NETWORK

- Klika – kompletní graf
- všechny neurony jsou zároveň vstupní i výstupní
- váhy jsou symetrické $w_{ji} = w_{ij}$
- Pokud síť dostane naučený vstup, měla by zůstavit a dál se neměnit.
- Pokud síť dostane nenaučený vstup, měla by se dostat do stavu odpovídajícímu nejbližšímu naučenému příkladu (vstupu)
- Training set je jen vektorů (množina) $\vec{x_k} \in \{-1; 1\}^m$
- Hebb's rule: _____ věime síť, které dvojice neuronů nabývají často stejnou hodnotu

$$w_{ji} = \sum_{k=1}^{t} X_{kj} X_{ki} \qquad 1 \leq j \neq i \leq m$$

.... pokud jsou $X_{kj}$ a $X_{ki}$ stejná (oba 1 nebo oba -1), vyjde veliké positivní číslo => jejich spojení by mělo být silné.

.... pokud jsou rozdílná => vyjde veliké negativní číslo => jejich spojení by mělo být slabé

- **Activity:** – nejprve nastavíme neurony na vstupy sítě

    – cyklicky upravujeme stavy neuronů:

    1) vnitřní potenciál:

    $$\xi_j^{(t)} = \sum_{i=1}^{m} w_{ji} y_i^{(t)}$$

    2)
    $$y_j^{(t+1)} = \begin{cases} 1 & \xi_j^{(t)} > 0 \\ y_j^{(t)} & \xi_j^{(t)} = 0 \\ -1 & \xi_j^{(t)} < 0 \end{cases}$$

    → tento proces zůstaví na korelém vstupu

    $$y_j^{(t^*+n)} = y_j^{(t^*)} \qquad (j = 1, \ldots, m)$$

- <u>Energie</u> (potenciální) pro každý stav $\vec{y} \in \{-1; 1\}^n$

$$E(\vec{y}) = -\frac{1}{2} \sum_{j=1}^{n} \sum_{i=1}^{n} w_{ji} \, y_j \, y_i$$

    – stavy s malou energií jsou stabilní (málo neuronů „chce změnit" svůj stav)

    – stavy s velikou energií jsou málo stabilní $\Rightarrow$ large (positive) $w_{ji} \, y_j \, y_i$ je stabilní

               small (negative) $w_{ji} \, y_j \, y_i$ není stabilní

- Energie se během procesu nezvyšuje

- <u>Phantoms</u>

    = lokální minimum energetické fce, které nekorespondije s trénínkovým příkladem

    – mohou být "odnaučeni" (inverzní Hebb's rule):

$$w_{ji}' = w_{ji} - x_i \, x_j$$

- <u>Kapacita Hopfieldovy sítě</u>

$$p \leq n / (4 \log n)$$

           $p$ = počet tréninkových příkladů schopných se naučit

           $n$ = počet neuronů

- <u>Boltzman activity</u>

    – místo chození cyklicky dokola, v každém kroku vybereme náhodné neuron.

    ~~neuron~~ inner potenciál :

$$\xi_j^{(t)} = \sum_{i=1}^{n} w_{ji} \, y_i^{(t)}$$

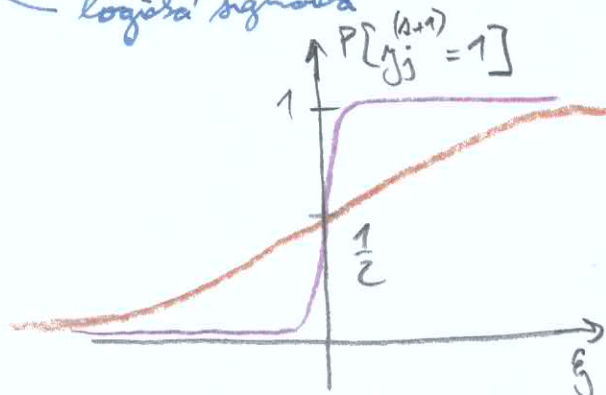poté vyber náhodně 1 nebo -1    ( $y_j^{(t+1)} \in \{-1, 1\}$ )

$$P\left[ y_j^{(t+1)} = 1 \right] = \sigma\left( \xi_j^{(t)} \right)$$

kde

$$\sigma(\xi) = \frac{1}{1 + e^{-2\xi / T(t)}}$$

                   $T(t)$ je teplota v čase $t$

                   logická sigmoida
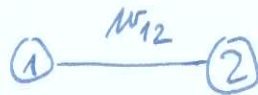
- vysoká teplota $\Rightarrow$ náhodné chování
- nízká teplota $\Rightarrow$ stabilní chování

- Boltzman Activity = Hopfield network + random noise
- Energie může poskočit na vyšší úroveň, v závislosti na teplotě

- Hopfield networks – příklad

  $\boxed{\text{2 neurony :}}$

  $\boxed{\text{4 stavy :}}$

  $(1) \xrightarrow{\;w_{12}\;} (2)$

- pravděpodobnost přechodu $a$:

  $P[a] = \boxed{\dfrac{1}{2}} \cdot P_1^a \qquad$ pst., že vybereme neuron 1

  $\xi_1^{(t)} = w_{12} \cdot 1 \Rightarrow P_1^a = \dfrac{1}{1+e^{-2w_{12}}}$

  teplota je 1



$(1,1) \xleftarrow{\;\;a\;\;} (-1,1)$

$(1,-1) \xleftarrow{\quad} (-1,-1)$

$P_1^a$ je pravděpodobnost 1 v neuronu 1

$\Rightarrow$ velké váhy $\rightarrow$ neurony chtějí mít stejné hodnoty

$\Rightarrow$ malé váhy $\rightarrow$ neurony chtějí mít různé hodnoty

- Simulated annealing

  – ke lepšímu dosažení glob. minima energie

  - začít s vysokou teplotou a postupně ji snižovat

  $\boxed{\text{BOLTZMAN MACHINE}}$

- Architektura

  – NN s cykly a neorientovanými vahami

  – jinak vlastně stejné jako Hopfield N.

- Aktivita:

  – v každém kroku $t+1$:

  1) vyber náhodně neuron

  2) spočítej vnitřní potenciál

  $$\xi_j^{(t)} = \sum_{i \in j \leftarrow}^{n} w_{ji}\, y_i^{(t)}$$

  3) vyber $y_j^{(t+1)} \in \{-1; 1\}$ náhodně:

  $$P[y_j^{(t+1)} = 1] = \sigma\left(\xi_j^{(t)}\right) \quad \cdots \quad T = \boxed{1}$$

$\Rightarrow$ definuje pravděpodobnostní rozložení vektorů $\{-1;1\}^{|N|}$

$\rightarrow$ Když necháme BM běžet dostatečně dlouho, relativní četnost jednotlivých stavů bude nezávislá na iniciálním stavu

$\rightarrow$ tyto četnosti bereme jako pravděpodobnosti jednotlivých stavů

$\Rightarrow$ během učení dostáváme pravděpod. rozložení stavů $\{-1;1\}^{|N|}$ a podle nich adaptujeme váhy tak, aby odpovídaly daným pravděp.

- Equilibrium
  - fixovaná teplota
  - dál tomu vůbec nerozumím ?

- Hidden neurons
  - rozdělíme N na 2 množiny:
    - visible V
    - hidden H

$\alpha \in \{-1;1\}^{|V|}$ :

$\not{p_N} \; p_V(\alpha) = \sum_{\beta \in \{-1;1\}^{|H|}} p_N(\alpha, \beta)$

pravděp., že stav visib. neuronů v serem. equil. je $\alpha$

  - ~~a je pst., že stav viditelných neuronů je v termálním equilibrium~~

given target prob.

$$E(W) = \sum_{\alpha \in \{-1;1\}^{|V|}} p_d(\alpha) \; \ln \frac{p_d(\alpha)}{p_V(\alpha)}$$

current learned prob.

- Maximum likelyhood

$$p_d(\alpha) = \#(\not{} \alpha, T) / m \quad , \; kde \; \#(a, T) \; je \; počet \; výskytů \; \not{} \alpha \; v \; T$$
$$m = |T|$$

naším cílem je najít konfiguraci (obarvan), že $p_V \approx p_d$

hledá $W$

likelyhood:

$$L(T) = \prod_{i=1}^{m} P_V(\vec{x_i}) \quad \leftarrow \; toto \; chceme \; maximalizovat$$

- Gradient descent

<div style="border:1px solid green; display:inline-block">RESTRICTED BOLTZMAN MACHINE</div>

- kompletní bipartitní graf



V

H

- stavy nejsou $-1$ a $1$, ale $\vec{y} \in \{0,1\}^{|N|}$
- ~~bias~~ : $w_{j0} \; \cdots \not{} \; y_0 \; je \; stále \; 0$
- Activity
  - v každém ~~sudém~~ sudém kroku : náhodně vyber hodnoty všech hidden neuronů
  - v každém lichém kroku : náhodně vyber hodnoty všech visible neuronů

- pro každé $j \in H$:

$$P\left[ y_j^{(k+1)} = 1 \right] = \frac{1}{1 + \exp\left(-w_{j0} - \sum_{i \in V} w_{ji} y_i^{(k)}\right)}$$

- pro $j \in V$ :

  to stejné, ale $\displaystyle\sum_{i \in H}$