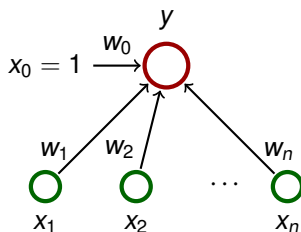


## Perceptron a ADALINE

- ▶ Perceptron
- ▶ Perceptronové učící pravidlo
- ▶ Konvergence učení perceptronu
- ▶ ADALINE
- ▶ Učení ADALINE

## Organizační dynamika:



$\vec{w} = (w_0, w_1, \dots, w_n)$  a  $\vec{x} = (x_0, x_1, \dots, x_n)$  kde  $x_0 = 1$ .

## Aktivní dynamika:

- ▶ vnitřní potenciál:  $\xi = w_0 + \sum_{i=1}^n w_i x_i = \sum_{i=0}^n w_i x_i = \vec{w} \cdot \vec{x}$
- ▶ aktivační funkce:  $\sigma(\xi) = \begin{cases} 1 & \xi \geq 0; \\ 0 & \xi < 0. \end{cases}$
- ▶ funkce sítě:  $y[\vec{w}](\vec{x}) = \sigma(\xi) = \sigma(\vec{w} \cdot \vec{x})$

## Adaptivní dynamika:

- ▶ Dána množina **tréninkových vzorů**

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Zde  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , je vstup  $k$ -tého vzoru a  $d_k \in \{0, 1\}$  je očekávaný výstup.

( $d_k$  určuje, do které ze dvou kategorií dané  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn})$  patří).

- ▶ Vektor vah  $\vec{w} \in \mathbb{R}^{n+1}$  je **konzistentní s  $\mathcal{T}$**  pokud  $y[\vec{w}](\vec{x}_k) = \sigma(\vec{w} \cdot \vec{x}_k) = d_k$  pro každé  $k = 1, \dots, p$ .  
Množina  $\mathcal{T}$  je **vnitřně konzistentní** pokud existuje vektor  $\vec{w}$ , který je s ní konzistentní.
- ▶ Cílem je nalézt vektor  $\vec{w}$ , který je konzistentní s  $\mathcal{T}$  za předpokladu, že  $\mathcal{T}$  je vnitřně konzistentní.

## Online učící algoritmus:

Idea: Cyklicky prochází vzory a adaptuje podle nich váhy, tj. otáčí dělící nadrovinu tak, aby se zmenšila vzdálenost špatně klasifikovaného vzoru od jeho příslušného poloprostoru.

Prakticky počítá posloupnost vektorů vah  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ váhy v  $\vec{w}^{(0)}$  jsou inicializovány náhodně blízko 0
- ▶ v kroku  $t + 1$  je  $\vec{w}^{(t+1)}$  vypočteno takto:

$$\begin{aligned}\vec{w}^{(t+1)} &= \vec{w}^{(t)} - \varepsilon \cdot (y[\vec{w}^{(t)}](\vec{x}_k) - d_k) \cdot \vec{x}_k \\ &= \vec{w}^{(t)} - \varepsilon \cdot (\sigma(\vec{w}^{(t)} \cdot \vec{x}_k) - d_k) \cdot \vec{x}_k\end{aligned}$$

Zde  $k = (t \bmod p) + 1$  (tj. cyklické procházení vzorů) a  $0 < \varepsilon \leq 1$  je **rychlost učení**.

## Věta (Rosenblatt)

*Jestliže je  $\mathcal{T}$  vnitřně konzistentní, pak existuje  $t^*$  takové, že  $\vec{w}^{(t^*)}$  je konzistentní s  $\mathcal{T}$ .*

# Důkaz Rosenblattovy věty

Pro zjednodušení budeme dále předpokládat, že  $\varepsilon = 1$ .

Nejprve si algoritmus přepíšeme do méně kompaktní formy:

- ▶ váhy v  $\vec{w}^{(0)}$  jsou inicializovány náhodně blízko 0
- ▶ v kroku  $t + 1$  je  $\vec{w}^{(t+1)}$  vypočteno takto:
  - ▶ **Jestliže**  $\sigma(\vec{w}^{(t)} \cdot \vec{x}_k) = d_k$ , **pak**  $\vec{w}^{(t+1)} = \vec{w}^{(t)}$
  - ▶ **Jestliže**  $\sigma(\vec{w}^{(t)} \cdot \vec{x}_k) \neq d_k$ , **pak**
    - ▶  $\vec{w}^{(t+1)} = \vec{w}^{(t)} + \vec{x}_k$  pro  $d_k = 1$
    - ▶  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \vec{x}_k$  pro  $d_k = 0$

(Řekneme, že nastala korekce.)

kde  $k = (t \bmod p) + 1$ .

# Důkaz Rosenblattovy věty

(Pro daný vektor  $\vec{a} = (a_0, \dots, a_n)$  označme  $\|\vec{a}\|$  jeho eukleidovskou normu  $\sqrt{\vec{a} \cdot \vec{a}} = \sqrt{\sum_{i=0}^n a_i^2}$ )

Idea:

- ▶ Uvážíme *hodně dlouhý vektor* (spočítáme jak dlouhý)  $\vec{w}^*$ , který je konzistentní s  $\mathcal{T}$ .
- ▶ Ukážeme, že pokud došlo v kroku  $t + 1$  ke korekci vah (tedy buď  $\vec{w}^{(t+1)} = \vec{w}^{(t)} + \vec{x}_k$  nebo  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \vec{x}_k$ ), pak

$$\|\vec{w}^{(t+1)} - \vec{w}^*\|^2 \leq \|\vec{w}^{(t)} - \vec{w}^*\|^2 - \max_i \|\vec{x}_i\|^2$$

Všimněte si, že  $\max_i \|\vec{x}_i\| > 0$  *nezávisí* na  $t$ .

- ▶ Z toho plyne, že algoritmus nemůže udělat nekonečně mnoho korekcí.

## Dávkový učící algoritmus:

Vypočte posloupnost  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$  váhových vektorů.

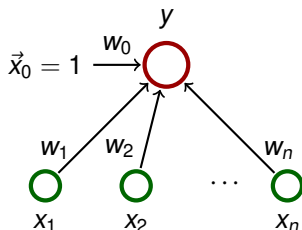
- ▶ váhy v  $\vec{w}^{(0)}$  jsou inicializovány náhodně blízko 0
- ▶ v kroku  $t + 1$  je  $\vec{w}^{(t+1)}$  vypočteno takto:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^p (\sigma(\vec{w}^{(t)} \cdot \vec{x}_k) - d_k) \cdot \vec{x}_k$$

Zde  $k = (t \bmod p) + 1$

a  $0 < \varepsilon \leq 1$  je rychlost učení.

## Organizační dynamika:



$\vec{w} = (w_0, w_1, \dots, w_n)$  a  $\vec{x} = (x_0, x_1, \dots, x_n)$  kde  $x_0 = 1$ .

## Aktivní dynamika:

- ▶ vnitřní potenciál:  $\xi = w_0 + \sum_{i=1}^n w_i x_i = \sum_{i=0}^n w_i x_i = \vec{w} \cdot \vec{x}$
- ▶ aktivační funkce:  $\sigma(\xi) = \xi$
- ▶ funkce sítě:  $y[\vec{w}](\vec{x}) = \sigma(\xi) = \vec{w} \cdot \vec{x}$



**Adaptivní dynamika:**

- ▶ Dána množina **tréninkových vzorů**

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

Zde  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$ ,  $x_{k0} = 1$ , je vstup  $k$ -tého vzoru a  $d_k \in \mathbb{R}$  je očekávaný výstup.

Intuice: chceme, aby síť počítala afinní aproximaci funkce, jejíž (některé) hodnoty nám předepisuje tréninková množina.

- ▶ **Chybová funkce:**

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

- ▶ Cílem je nalézt  $\vec{w}$ , které minimalizuje  $E(\vec{w})$ .

# Gradient chybové funkce

Uvažme **gradient** chybové funkce:

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

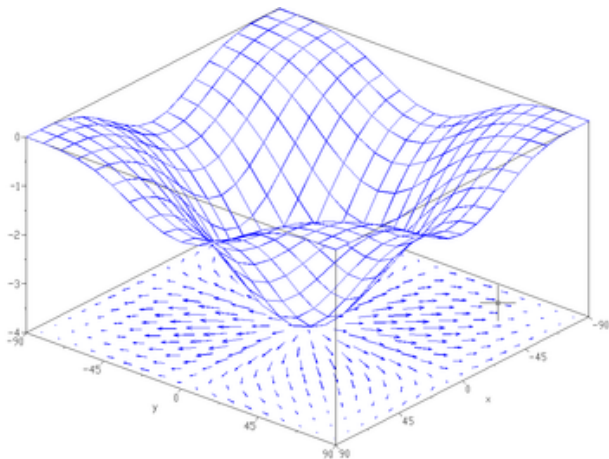
Intuice:  $\nabla E(\vec{w})$  je vektor ve **váhovém prostoru**, který ukazuje směrem nejstrmějšího „růstu“ funkce  $E(\vec{w})$ . Vektory  $\vec{x}_k$  zde slouží pouze jako parametry funkce  $E(\vec{w})$  a jsou tedy fixní!

## Fakt

*Pokud  $\nabla E(\vec{w}) = \vec{0} = (0, \dots, 0)$ , pak  $\vec{w}$  je globální minimum funkce  $E$ .*

Námi uvažovaná chybová funkce  $E(\vec{w})$  má globální minimum, protože je konvexním paraboloidem.

# Gradient - ilustrace



Pozor! Tento obrázek pouze ilustruje pojem gradientu, nezobrazuje chybovou funkci  $E(\vec{w})$

## Dávkový algoritmus (gradientní sestup):

- ▶ váhy v  $\vec{w}^{(0)}$  jsou inicializovány náhodně blízko 0
- ▶ v kroku  $t + 1$  je  $\vec{w}^{(t+1)}$  vypočteno takto:

$$\begin{aligned}\vec{w}^{(t+1)} &= \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)}) \\ &= \vec{w}^{(t-1)} - \varepsilon \cdot \sum_{k=1}^p (\vec{w}^{(t)} \cdot \vec{x}_k - d_k) \cdot \vec{x}_k\end{aligned}$$

Zde  $k = (t \bmod p) + 1$

a  $0 < \varepsilon \leq 1$  je rychlost učení.

(Všimněte si, že tento algoritmus je téměř stejný jako pro perceptron, protože  $\vec{w}^{(t)} \cdot \vec{x}_k$  je hodnota funkce sítě (tedy  $\sigma(\vec{w}^{(t)} \cdot \vec{x}_k)$  kde  $\sigma(\xi) = \xi$ .)

## Tvrzení

*Pro dostatečně malé  $\varepsilon > 0$  posloupnost  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$  konverguje (po složkách) ke globálnímu minimu funkce  $E$  (tedy k vektoru  $\vec{w}$ , který splňuje  $\nabla E(\vec{w}) = \vec{0}$ ).*

## Online algoritmus (Delta-rule, Widrow-Hoff rule):

- ▶ váhy v  $\vec{w}^{(0)}$  jsou inicializovány náhodně blízko 0
- ▶ v kroku  $t + 1$  je  $\vec{w}^{(t)}$  vypočteno takto:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t) \cdot (\vec{w}^{(t)} \cdot \vec{x}_k - d_k) \cdot \vec{x}_k$$

kde  $k = (t + 1) \bmod p$

a  $0 < \varepsilon(t) \leq 1$  je rychlost učení v kroku  $t + 1$ .

Všimněte si, že tento algoritmus nepracuje s celým gradientem, ale jenom s jeho částí, která přísluší právě zpracovávanému vzoru!

## Věta (Widrow & Hoff)

*Pokud  $\varepsilon(t) = \frac{1}{t}$  pak  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$  konverguje ke globálnímu minimu chybové funkce  $E$ .*

- ▶ Množina **tréninkových vzorů** je

$$\mathcal{T} = \{(\vec{x}_1, d_1), (\vec{x}_2, d_2), \dots, (\vec{x}_p, d_p)\}$$

kde  $\vec{x}_k = (x_{k0}, x_{k1}, \dots, x_{kn}) \in \mathbb{R}^{n+1}$  a  $d_k \in \{1, -1\}$ .

- ▶ Síť se natrénuje ADALINE algoritmem.
- ▶ Očekáváme, že bude platit následující:
  - ▶ jestliže  $d_k = 1$ , pak  $\vec{w} \cdot \vec{x}_k \geq 0$
  - ▶ jestliže  $d_k = -1$ , pak  $\vec{w} \cdot \vec{x}_k < 0$
- ▶ To nemusí vždy platit, ale často platí. Výhoda je, že se ADALINE algoritmus postupně stabilizuje i v neseparabilním případě (na rozdíl od perceptronového algoritmu).