

Neparametrické metody

Neparametrické metody

- obsah
 - odhad funkce hustoty
 - metoda Parzenova okénka
 - pravděpodobnostní neuronové sítě
 - metoda k_n -nejbližších sousedů
 - odhad a posteriorní pravděpodobnosti

 - pravidlo k-nejbližších sousedů
 - vlastnosti
 - zlepšení výpočetní složitosti
 - metriky

Odhady hustoty

Neparametrické metody

- hustota rozdělení
 - 1. známá
 - 2. neznámá, ale známý tvar funkce
 - parametry lze odhadnout
 - pro „běžné“ funkce hustoty (např. normální rozdělení, ...)
 - 3. neznámá → nejčastěji
 - problém: hustota reálných dat málokdy „pasuje“ na některou z běžných hustot

→ neparametrické metody

- použití: pro libovolné rozdělení (+)
 - není třeba předpoklad o tvaru hustoty
- potřeba více trénovacích dat (–)

Neparametrické metody

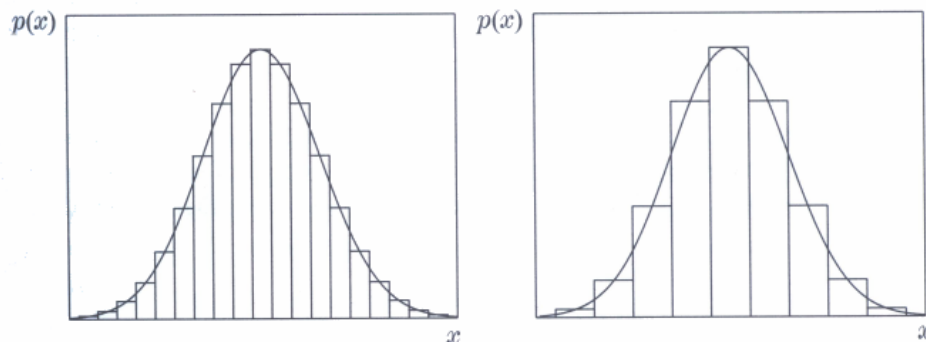
- rozdělení metod
 - odhady funkce hustoty $p(\mathbf{x}|\omega_j)$
 - odhad hustoty na základě trénovacích dat
 - uspokojivé odhady → použity jako skutečné hodnoty do klasifikátoru
 - proces nutné provést pro každou třídu ω_j

→ metoda histogramu, Parzenovo okénko
 - odhady apost. pravděpodobnosti $P(\omega_j|\mathbf{x})$
 - přímý odhad rozhodovacího pravidla

→ metody nejbližšího souseda

Odhady hustoty – myšlenka (1)

- princip metody
 - aproximace hustoty pomocí histogramu



- P ... pravděpodobnost, že vzor padne do přihrádky o rozměru R

$$P = \int_R p(\mathbf{x}') d\mathbf{x}' \quad (1)$$

neznáme

- obráceně
 - když se odhadne $P \rightarrow$ lze odhadnout i $p(\mathbf{x})$

Odhady hustoty – myšlenka (2)

- trénovací vzory $\mathbf{x}_1, \dots, \mathbf{x}_n$ vybrané nezávisle podle $p(\mathbf{x})$
 - k počet vzorů padnoucí do oblasti R
 - n počet vzorů celkem
- P_k ... pravděpodobnost, že k z n vzorů spadne do oblasti R

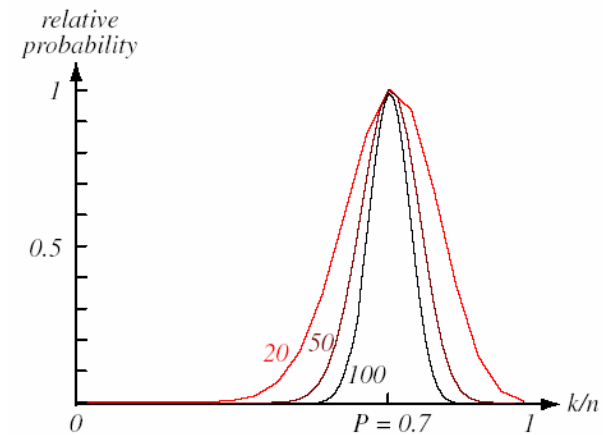
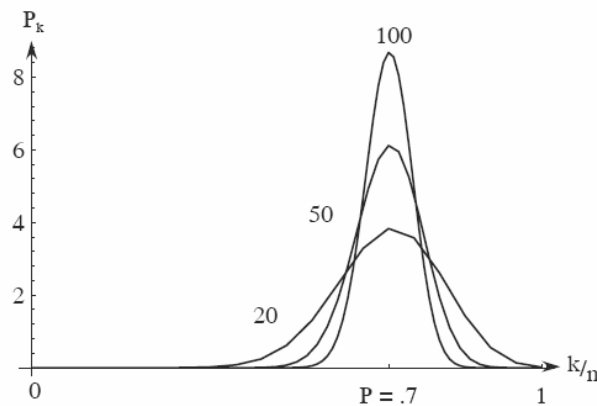
$$P_k = \binom{n}{k} P^k (1-P)^{n-k}$$

→ P_k má binomické rozdělení a střední hodnota pro k dat je

$$E(k) = n \cdot P$$

Odhady hustoty – ilustrace

- závislost P_k na hodnotě k/n
 - skutečná hodnota hustoty rozdělení, z něhož se v bodě x vybíralo
 $P = 0.7$
 - každá křivka je binomická

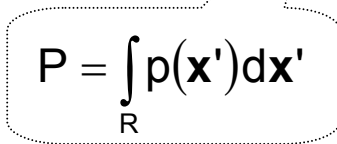


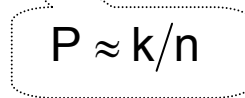
- rostoucí n ... křivka s ostřejším vrcholem
- $n \rightarrow \infty$... Diracova delta funkce

Odhady hustoty – myšlenka (3)

- binomické rozdělení P_k
 - velmi špičaté okolo své střední hodnoty
 - $P \approx k/n$ (2)
 - dobrý odhad zejména pro velké n
- kombinace vztahů (1) a (2)
 - pravděpodobnost, že \mathbf{x} padne do oblasti R

$$P = \int_R p(\mathbf{x}') d\mathbf{x}' \approx \frac{k}{n} \quad (3)$$


$$P = \int_R p(\mathbf{x}') d\mathbf{x}'$$


$$P \approx k/n$$

Odhady hustoty – myšlenka (4)

- když oblast R malá a p se v oblasti R příliš nemění
→ oblast lze aproximovat obdélníkem

- V ... „šířka“ oblasti R

$$V = \int_R 1 d\mathbf{x}'$$

- odhad pravděpodobnosti P

$$P = \int_R p(\mathbf{x}') d\mathbf{x}' \approx p(\mathbf{x}) \cdot V \quad (4)$$

kde \mathbf{x} je bod v R

- odhad hustoty $p(\mathbf{x})$ v bodě \mathbf{x}

$$p(\mathbf{x}) \approx \frac{k/n}{V}$$

(4)

$$P = \int_R p(\mathbf{x}') d\mathbf{x}' \approx p(\mathbf{x}) \cdot V$$

(3)

$$P = \int_R p(\mathbf{x}') d\mathbf{x}' \approx \frac{k}{n}$$

Odhady hustoty – praktický problém

- n fixní a fixní velikost V

- k/n konverguje k $p(\mathbf{x}) \cdot V$

prostorově zprůměrovaný
odhad $p(\mathbf{x})$

- pro získání $p(\mathbf{x})$ je třeba $V \rightarrow 0$

$$p(\mathbf{x}) \approx \frac{k/n}{V}$$

- n fixní a $V \rightarrow 0$

- velikost oblasti se zmenšuje „do nekonečna“

- do oblasti nepadne žádný vzor → odhad $p(\mathbf{x}) \approx 0$

- do oblasti náhodou padne 1 nebo více vzorů → odhad $p(\mathbf{x}) \approx \infty$

nepoužitelné odhady

- odhad vždy vyhlazen

Odhady hustoty – jak obejít problém

- teoreticky
 - nekonečně vzorů pro dané rozdělení
- prakticky
 - pousloupnosti oblastí $R_1, R_2 \dots$ různé velikosti kolem \mathbf{x}
 - každá R_i obsahuje vzor \mathbf{x}
 - R_1 pracuje s 1 vzorem
 - R_2 pracuje se 2 vzory
 -
 - $V_n \dots$ „šířka“ oblasti R_n
 - n -tý odhad hustoty $p(\mathbf{x})$ v bodě $\mathbf{x} \dots p_n(\mathbf{x}) \cong \frac{k_n/n}{V_n}$
 - $k_n \dots$ počet vzorů, které padnou do R_n

Odhady hustoty – konvergence

- konvergence $p_n(\mathbf{x}) \rightarrow p(\mathbf{x})$

$$\lim_{n \rightarrow \infty} V_n = 0$$

$$\lim_{n \rightarrow \infty} k_n = \infty, \quad p(\mathbf{x}) \neq 0$$

$$\lim_{n \rightarrow \infty} \frac{k_n}{n} = 0$$

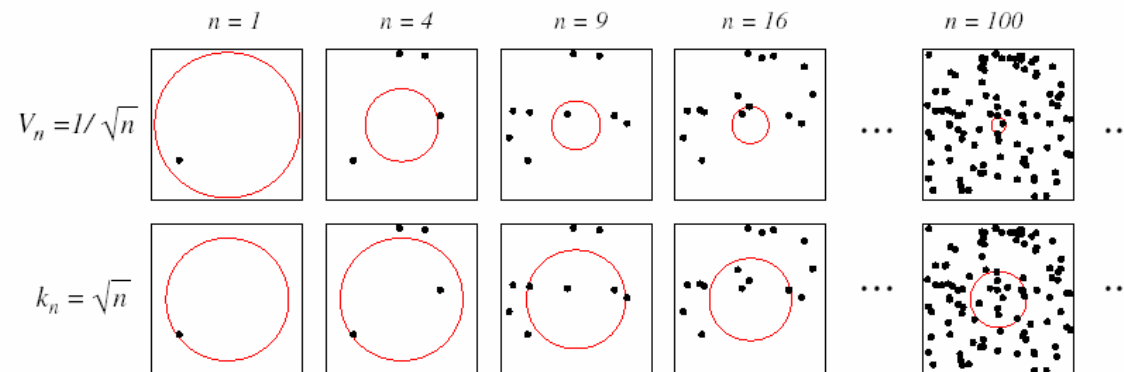
oblasti se zmenšují

z nekonečně vzorků jich nekonečně mnoho padne do oblasti R_n

ačkoliv velké množství vzorků padne do $R_n \rightarrow$
z celkového počtu vzorků n je to zanedbatelné množství

Vytvoření posloupností oblastí

- „šířka“ oblasti V_n funkcí n
 - oblast se zmenšuje určením V_n jako funkci n
 - např. $V_n = 1 / \sqrt{n}$
 - metoda: Parzenovo okénko
- počet vzorků k_n funkcí n
 - oblast je tak velká, aby obsahovala k_n sousedů vzoru x
 - např. $k_n = \sqrt{n}$
 - metoda: k_n -nejbližších sousedů



Metoda Parzenova okénka

- oblast R_n ... d -dimenzionální hyperkostka

- střed hyperkostky v bodě \mathbf{x}
- objem hyperkostky se stranou h_n

$$V_n = h_n^d$$

→ počet vzorů, které padnou do hyperkostky, se definuje pomocí okénka

- okénko ... jednotková hyperkostka se středem v počátku

$$\varphi(\mathbf{u}) = \begin{cases} 1 & |u_j| \leq 1/2 \text{ pro } j = 1, \dots, d \\ 0 & \text{jinak} \end{cases}$$

vzor \mathbf{u} padne do okénka

vzor \mathbf{u} nepadne do okénka

- okénko ... hyperkostka o objemu V_n vycentrovaná v \mathbf{x}

$$\varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) = \begin{cases} 1 & \text{když vzor } \mathbf{x}_i \text{ padne do hyperkostky} \\ & \text{o objemu } V_n \text{ vycentrované v } \mathbf{x} \\ 0 & \text{jinak} \end{cases}$$

Metoda Parzenova okénka

- počet vzorů k_n v hyperkostce

$$k_n = \sum_{i=1}^n \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

$$p_n(\mathbf{x}) \cong \frac{k_n/n}{V_n}$$

- n -tý odhad hustoty $p_n(\mathbf{x})$

$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

- spojitá funkce $p(\mathbf{x}) \rightarrow$ aproximace nespojitou skokovou funkcí $\varphi(\mathbf{u})$
→ odhad zatížen chybou

- **Parzen:** odhad pomocí hladké funkce $\varphi(\mathbf{u})$

- okénko $\varphi(\mathbf{u})$

$$\varphi(\mathbf{u}) \geq 0 \quad \text{a} \quad \int \varphi(\mathbf{u}) d\mathbf{u} = 1$$

- pro takové okénko $\varphi(\mathbf{u}) \rightarrow$ odhad $p_n(\mathbf{x})$ je funkcí hustoty

Parzenovo okénko – vliv h_n

- vliv šířky okénka h_n na odhad $p_n(\mathbf{x})$
 - funkce $\delta_n(\mathbf{x})$ Diracův puls

$$\delta_n(\mathbf{x}) = \frac{1}{V_n} \varphi\left(\frac{\mathbf{x}}{h_n}\right)$$

- $p_n(\mathbf{x})$ pomocí $\delta_n(\mathbf{x})$

$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \delta_n(\mathbf{x} - \mathbf{x}_i)$$

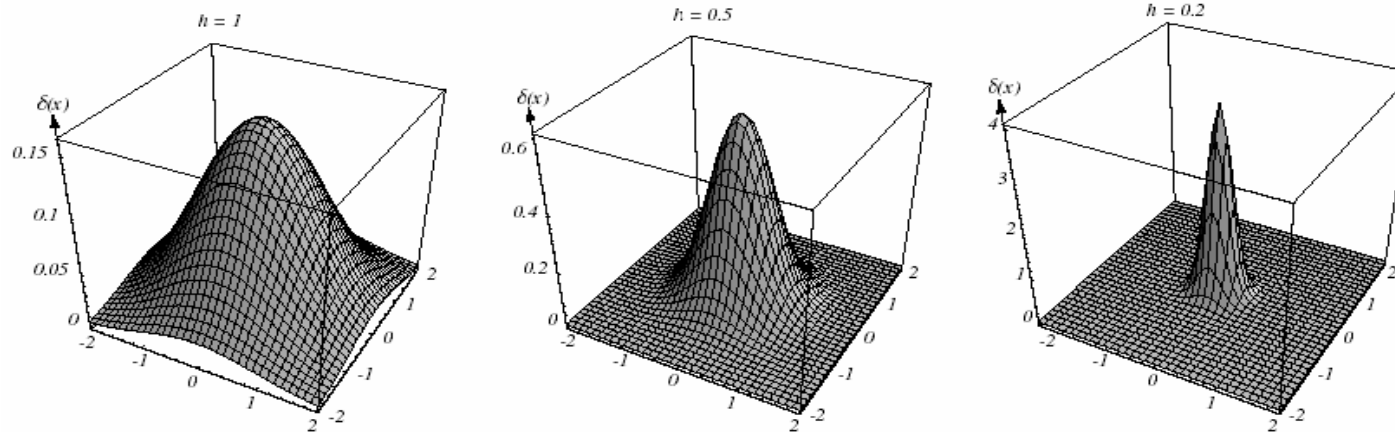
$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

$$\int \delta_n(\mathbf{x}) d\mathbf{x} = \int \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) d\mathbf{x} = \int \varphi(\mathbf{u}) d\mathbf{u} = 1$$

- Diracův puls δ_n
 - δ_n má konstantní objem
 - $V_n = h_n^d \rightarrow h_n$ má vliv na amplitudu (výšku) i šířku funkce $\delta_n(\mathbf{x})$

Parzenovo okénko – příklady

- příklad ($dim=2$)
 - Parzenovo okénko ... 2D-rotační symetrická normální funkce
 - δ_n $h_n=1$ a $h_n=0,5$ a $h_n=0,2$



$\delta_n(\mathbf{x})$ normalizována → na svislé ose různá škálování

Vlastnosti Parzenova okénka

- volba šířky h_n

- h_n velmi velké

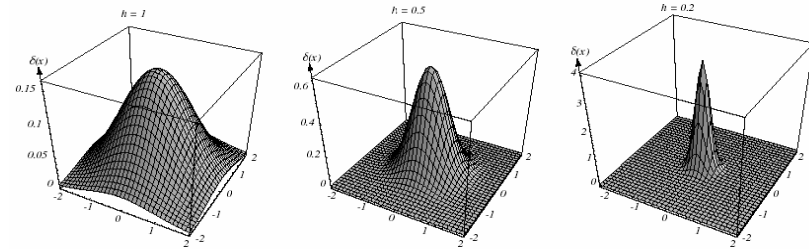
- malá amplituda δ_n (nízký vrchol)
- \mathbf{x} daleko od \mathbf{x}_i aby se změnilo $\delta_n(\mathbf{x}-\mathbf{x}_i)$ od $\delta_n(\mathbf{0})$

- h_n velmi malé

- vysoký vrchol $\delta_n(\mathbf{x}-\mathbf{x}_i)$
- vrchol blízko $\mathbf{x}=\mathbf{x}_i$

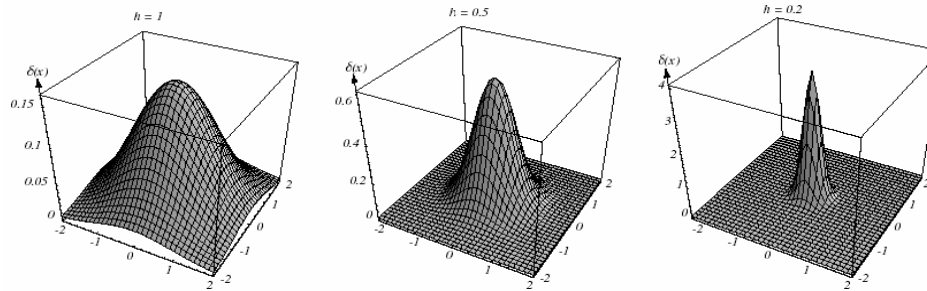
- h_n dosáhne 0

- $\delta_n(\mathbf{x}-\mathbf{x}_i)$ Diracova delta funkce vycentrovaná v \mathbf{x}_i

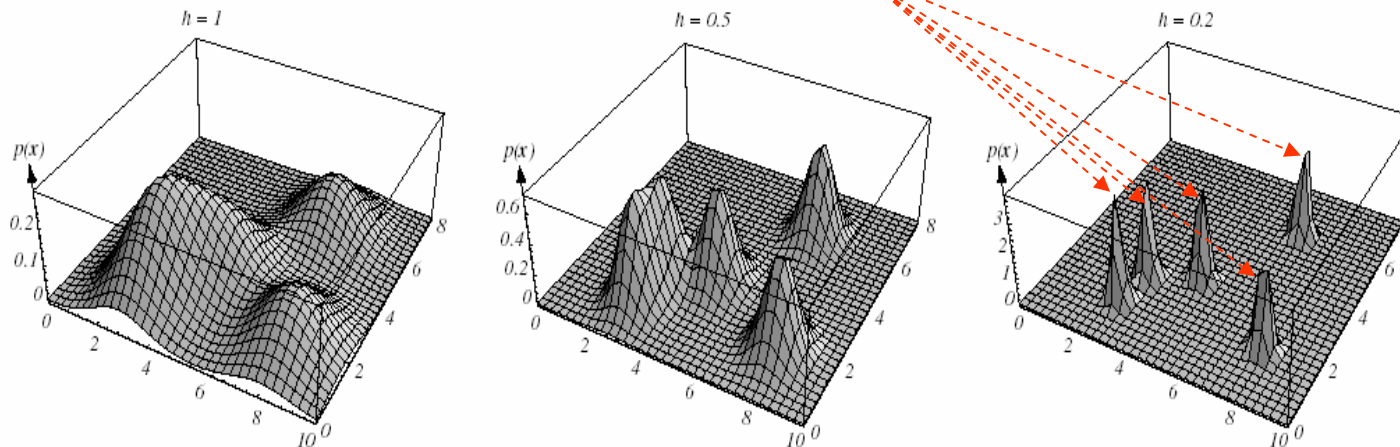


Parzenovo okénko – příklady

- Parzenova okénka – tvar δ_n



- odhad $p_n(\mathbf{x})$ pomocí 5 trénovacích vzorů



na svislé ose různá škálování

Ilustrační příklad 1 (dim=1)

- normální rozdělení $p(\mathbf{x}) \sim N(0, 1)$
 - okénko ... Gaussova funkce

$$\varphi(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$$

- „šířka“ kostky h_n

$$h_n = \frac{h_1}{\sqrt{n}}$$

h_1 ... parametr ovlivňující šířku okénka (uživatelsky definovaný)

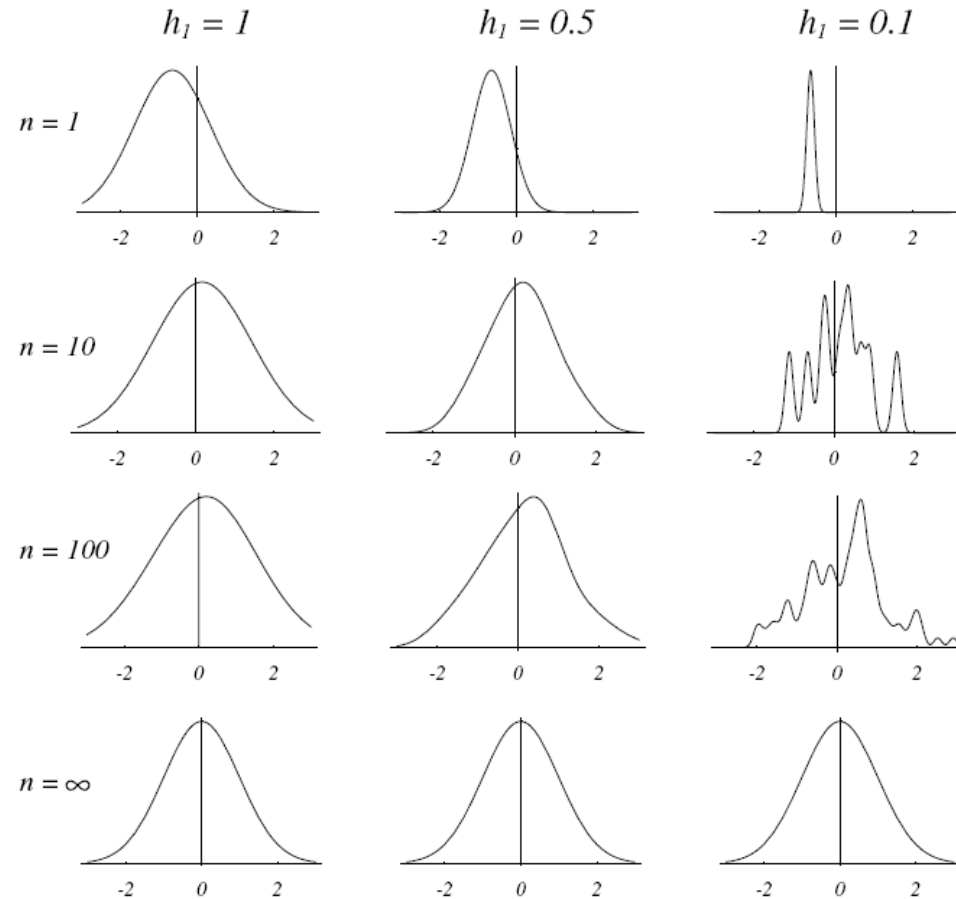
→ $p_n(\mathbf{x})$

$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

$d=1 \rightarrow V_n = h_n$

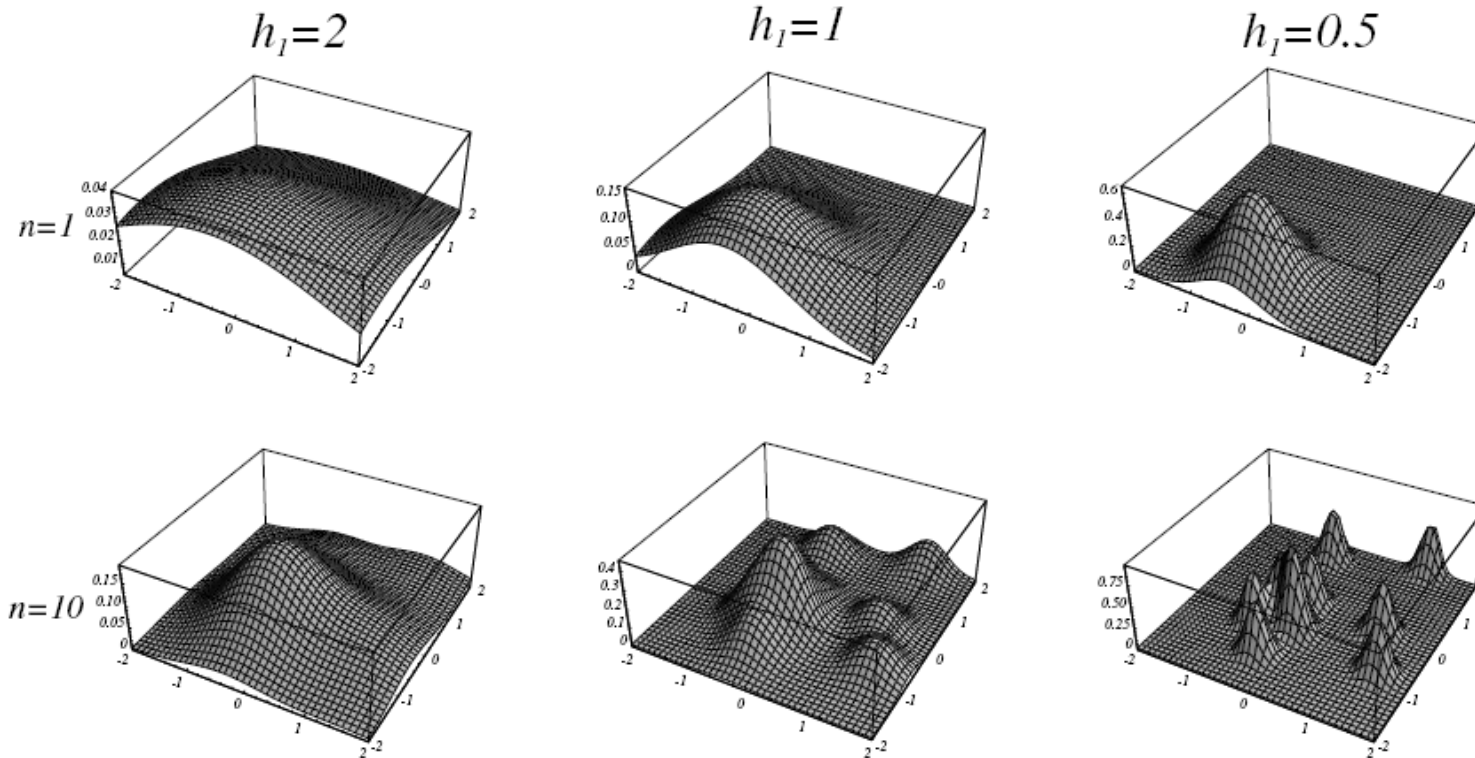
Ilustrační příklad 1 – výsledky

- odhad neznámé $p(\mathbf{x}) \sim N(0, 1)$
 - iniciální šířka okénka h_1 a počet vzorů n



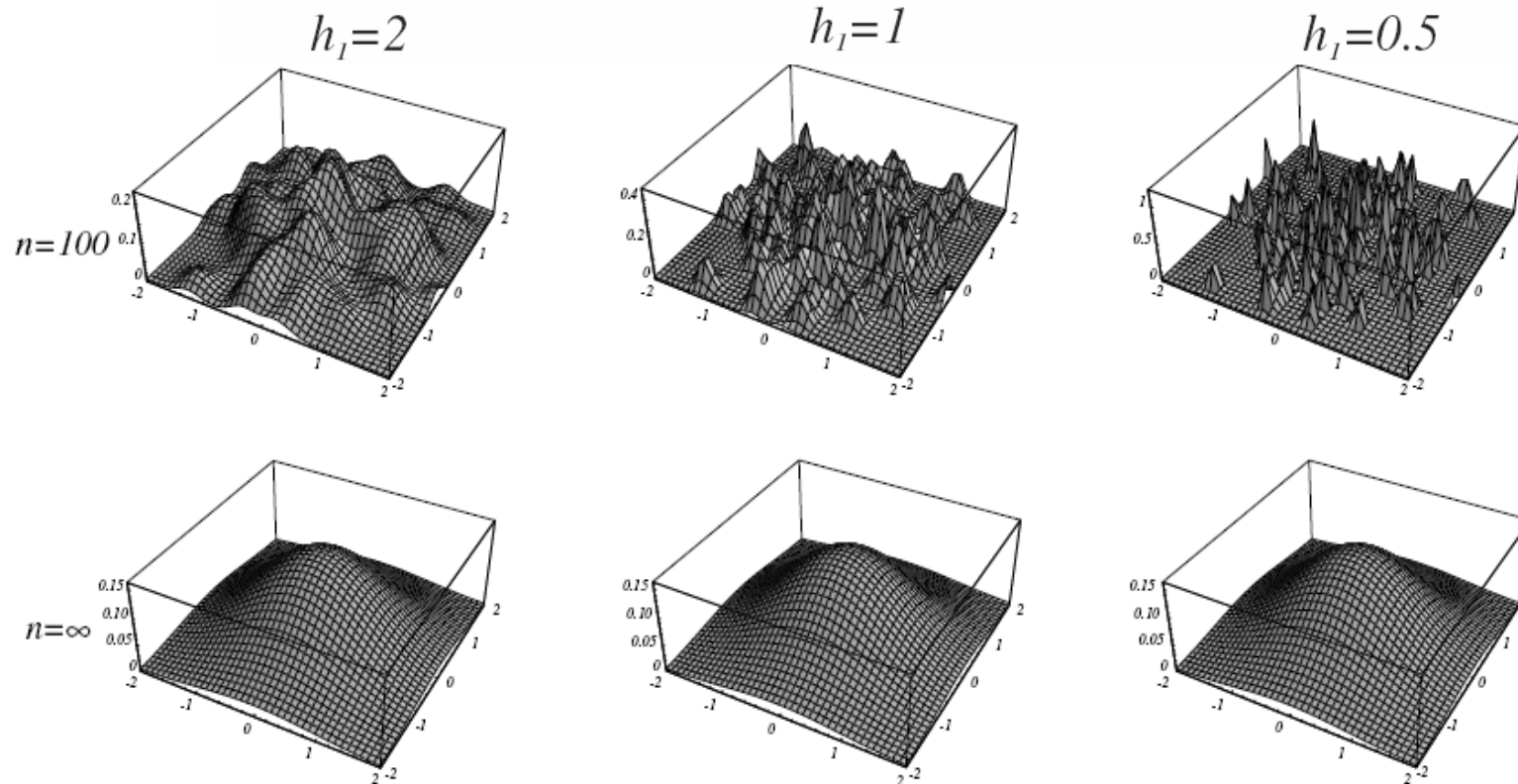
Ilustrační příklad 2 (dim=2)

- odhad $p(\mathbf{x})$... dvourozměrné normální rozdělení
 - iniciální šířky okénka h_1 a počet vzorů n ($n=1$ a $n=10$)



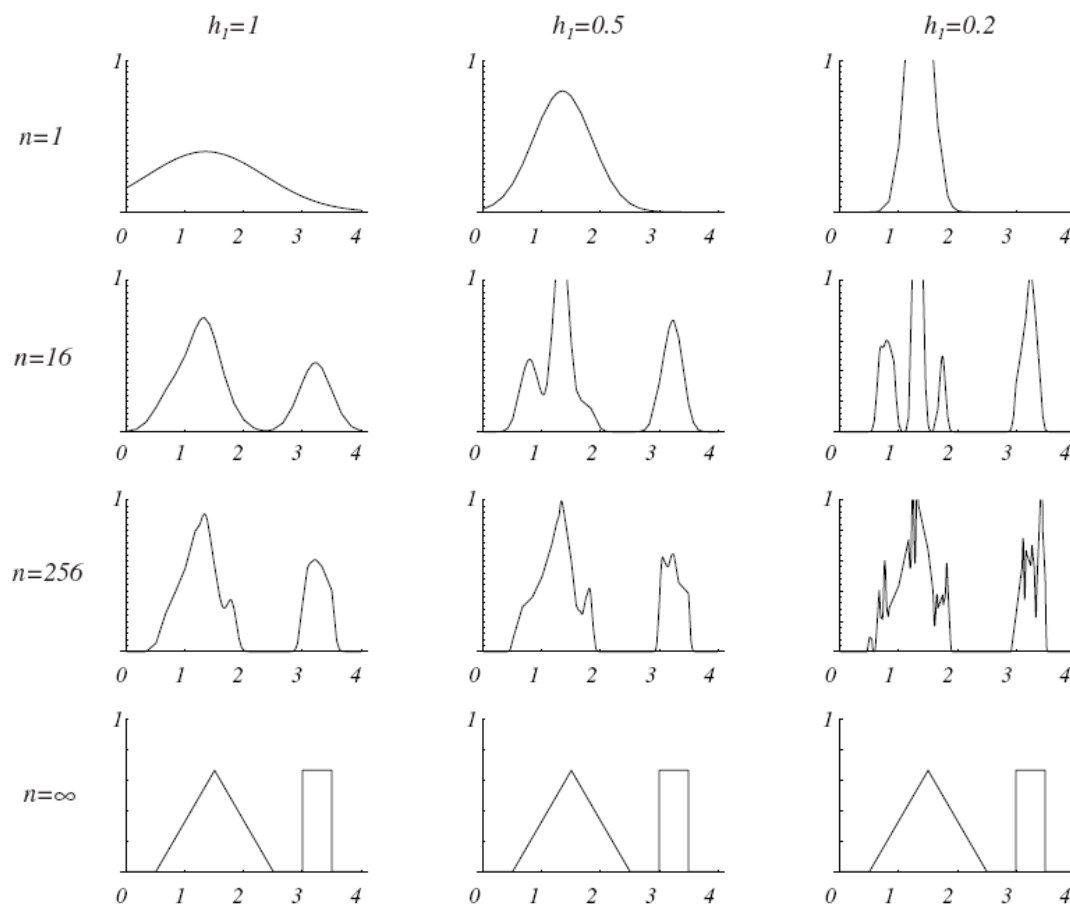
Ilustrační příklad 2 (dim=2)

- odhad $p(\mathbf{x})$... dvourozměrné normální rozdělení
 - iniciální šířky okénka h_1 a počet vzorů n ($n=100$ a $n=\infty$)



Ilustrační příklad 3 (dim=1)

- neznámá hustota: směs uniformní a trojúhelníkové hustoty



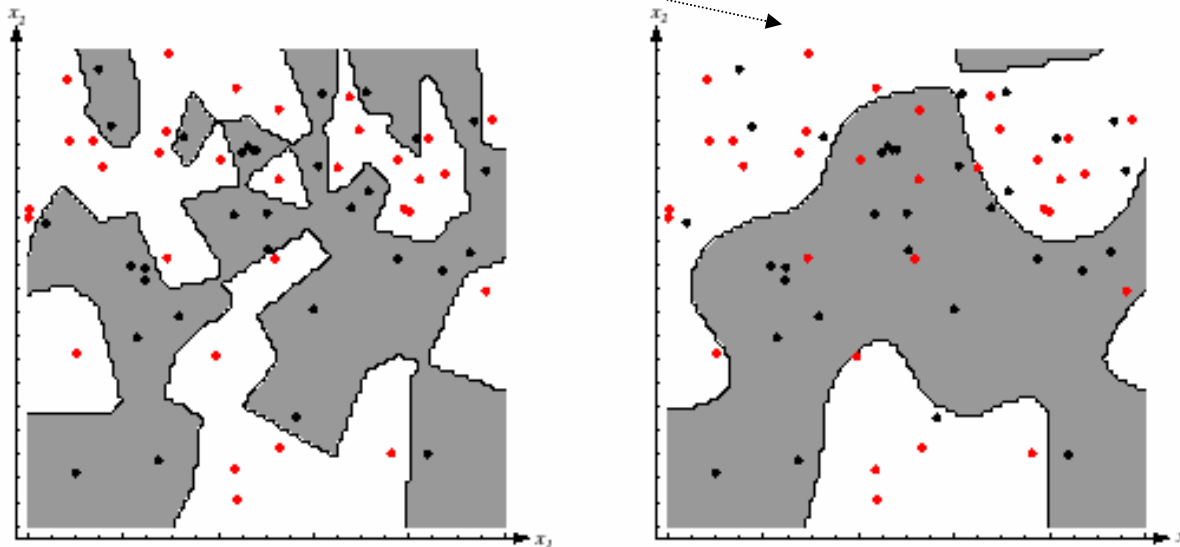
Klasifikace – Parzenovo okénko

- klasifikátor založený na Parzenově okénku
 - 1. předložení vzoru \mathbf{x}
 - 2. vypočten odhad hustoty v bodě \mathbf{x} pro každou třídu
 - 3. klasifikace vzoru \mathbf{x} podle minimální ztráty
- klasifikátor pro 2 třídy
 - \mathbf{x} do ω_1 (resp. ω_2) když

$$\frac{\frac{1}{N_1} \sum_{i=1}^{N_1} \frac{1}{h} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right)}{\frac{1}{N_2} \sum_{i=1}^{N_2} \frac{1}{h} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right)} > (<) \frac{P(\omega_2) \lambda_{21} - \lambda_{22}}{P(\omega_1) \lambda_{12} - \lambda_{11}}$$

Klasifikace – rozhodovací hranice

- rozhodovací hranice vliv okénkové funkce φ a šířky h
 - h malé \rightarrow komplikované oblasti
 - h velké \rightarrow jednodušší oblasti

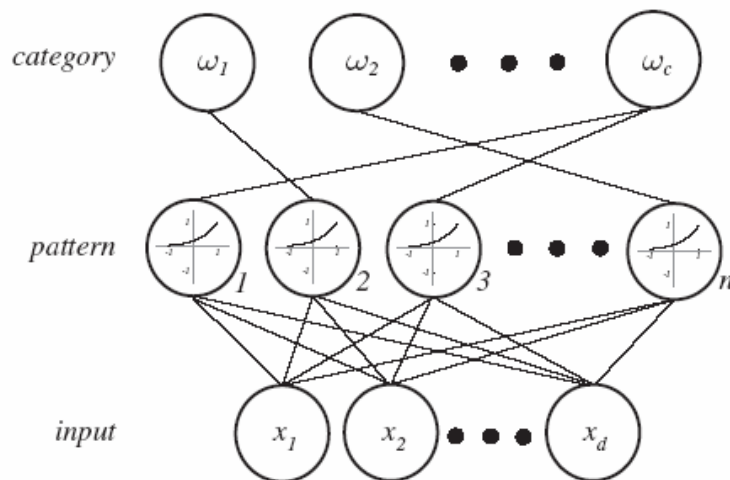


Parzenovo okénko – závěr

- klasifikátor
 - okénko úzké → malá chyba na trénovacích vzorech → problém přeučení
- metoda Parzenova okénka
 - obecné (+)
 - velké množství trénovacích dat (–)
 - časová i prostorová náročnost (–)
 - počet trénovacích dat roste exponenciálně s dimenzí dat
→ „prokletí dimenze“
 - mnohodimenzionální funkce typicky mnohonásobně složitější než nízko-dimenzionální funkce

Pravděpodobnostní neuronové sítě (PNS)

- paralelní implementace Parzenova okénka
 - neuronové sítě (NS)
- vstup
 - n d -dimenzionálních vzorů v c třídách



- d vstupních neuronů
- n skrytých neuronů
- c výstupních neuronů
- vstupní neuron → spojen s **každým** skrytým neuronem
- skrytý neuron → spojen s **právě jedním** výstupním neuronem

Pravděpodobnostní neuronové sítě (PNS)

- vstupní vrstva ↔ skrytá vrstva
 - váhy, které je nutné naučit

- skrytý neuron

- skalární součin váhového vektoru \mathbf{w} a normalizovaného předloženého vzoru \mathbf{x}

$$\text{net} = \mathbf{w}^T \mathbf{x}$$

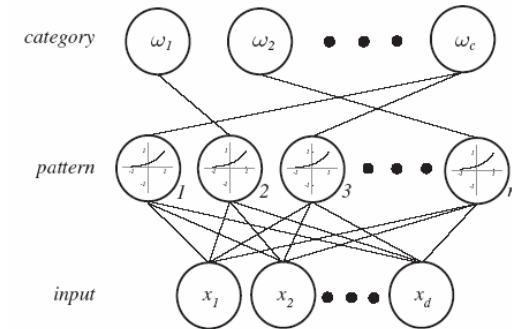
- skalární součin → aplikována funkce

$$\exp [(\text{net}-1)/\sigma^2]$$

- σ ... uživatelsky definovaný parametr

- výstupní neuron

- součet výstupů ze skrytých neuronů, se kterými je spojen



PNS – učení

- 1. normalizace trénovacích vzorů

$$\sum_{i=1}^d x_i^2 = 1$$

- 2. předložení trénovacího vzoru \mathbf{x}_k
 - I) aktualizace vah: vstupní neurony \leftrightarrow k -tý skrytý neuron
$$\mathbf{w}_k = \mathbf{x}_k$$
 - II) vytvoření vazby: k -tý skrytý neuron \leftrightarrow výstupní neuron i pro třídu ω_i vzoru \mathbf{x}_k
$$a_{ki} = 1$$
- 3. proces zopakován pro všechny trénovací vzory

PNS – algoritmus učení

- značení

- $\mathbf{x}_j = (x_{j1}, \dots, x_{jn}) \quad j=1, \dots, n$
- $\mathbf{w}_j = (w_{j1}, \dots, w_{jn}) \quad j=1, \dots, n$

- algoritmus

```
1. begin
2.   initialize  $j \leftarrow 0, n, a_{ij} \leftarrow 0 \quad j=1, \dots, n \quad i=1, \dots, c$ 
3.   do  $j \leftarrow j+1$ 
4.      $\mathbf{x}_{jk} \leftarrow \mathbf{x}_{jk} / (\sum \mathbf{x}_{ji}^2)^{1/2} \quad // \text{normalizace}$ 
5.      $\mathbf{w}_{jk} \leftarrow \mathbf{x}_{jk} \quad // \text{učení}$ 
6.     if ( $\mathbf{x}$  in  $\omega_c$ ) then  $a_{jc} \leftarrow 1$ 
7.   until  $j=n$ 
8. end
```


PNS – klasifikace

- PNN klasifikace

- 1. **normalizace** testovaného vzoru \mathbf{x}
- 2. **předložení** vzoru \mathbf{x} vstupní vrstvě

- 3a. každý skrytý **neuron** k vypočte skalární součin net_k

$$net_k = \mathbf{w}_k^T \cdot \mathbf{x}$$

- 3b. na net_k aplikována **aktivační** funkce

$$\exp\left(\frac{net_k - 1}{\sigma^2}\right)$$

σ ... uživatelsky definovaný parametr → šířka Gaussova okénka

- 4. každý **výstupní neuron** sečte příspěvky od **skrytých neuronů**, se kterými je spojen

PNS – klasifikace

- **aktivační funkce** → exponenciála
 - odvození

$$\varphi\left(\frac{\mathbf{x} - \mathbf{w}_k}{h_n}\right) \propto e^{-\frac{(\mathbf{x} - \mathbf{w}_k)^T (\mathbf{x} - \mathbf{w}_k)}{2\sigma^2}} = e^{-\frac{(\mathbf{x}^T \mathbf{x} + \mathbf{w}_k^T \mathbf{w}_k - 2\mathbf{x}^T \mathbf{w}_k)}{2\sigma^2}} = e^{\frac{\text{net}_k - 1}{\sigma^2}}$$

Parzenovo okénko
= příslušný Gaussián

$$\mathbf{w}_k^T \cdot \mathbf{w}_k = \mathbf{x}^T \cdot \mathbf{x} = 1, h_n(\sigma) \text{ konstanta}$$

- **výstup skrytého neuronu**
 - pravděpodobnost, že vzor \mathbf{x} vygenerován Gaussovou umístěnou ve středu trénovacího vzoru
- **výstupní neuron**
 - součet lokálních odhadů skrytých neuronů → diskriminační funkce $g_i(\mathbf{x})$
 - $\max\{g_i(\mathbf{x})\}$ → hledaná třída pro vzor \mathbf{x}

PNS – algoritmus klasifikace

- algoritmus

```
1. begin
2.   initialize  $k \leftarrow 0$ ,  $\mathbf{x} \leftarrow \text{test\_pattern}$ 
3.   do  $k \leftarrow k+1$ 
4.      $\text{net}_k \leftarrow \mathbf{w}_k^T \cdot \mathbf{x}$ 
5.     if ( $a_{ki}=1$ ) then  $g_i \leftarrow g_i + \exp((\text{net}_k - 1) / \sigma^2)$ 
6.   until  $k=n$ 
7.   return  $\text{class} \leftarrow \text{argmax}_i(g_i(\mathbf{x}))$ 
8. end
```

- vlastnosti PNS

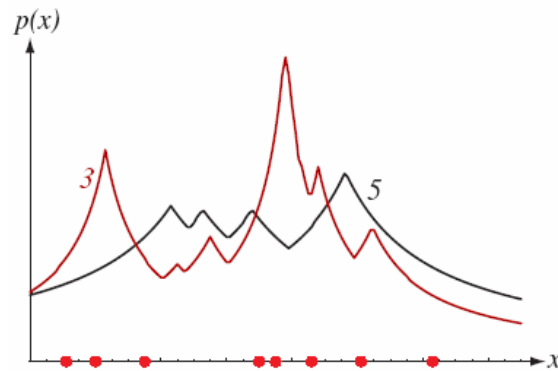
- rychlost učení
 - jeden průchod trénovací množinou
- prostorová složitost
 - úměrná počtu vazeb $O((n+1)d)$
- lze použít online učení
 - nový vzor lze snadno zabudovat do již naučeného modelu

Metoda k_n -nejbližších sousedů

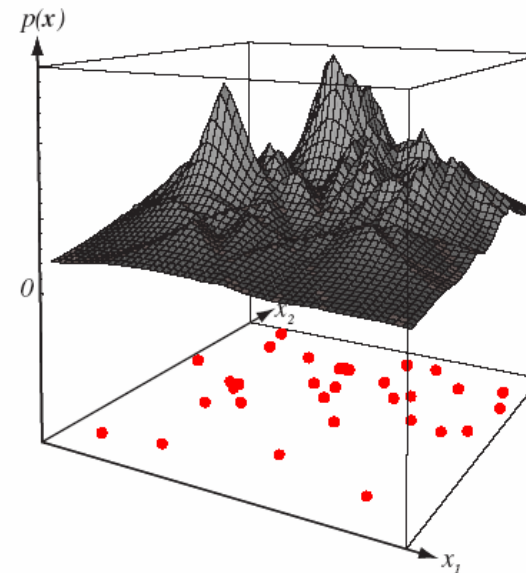
- odhad hustoty $p_n(\mathbf{x})$ z n trénovacích vzorů
 - okénko okolo \mathbf{x}
 - okénko se zvětšuje, dokud v něm není k_n vzorů
 - k_n -nejbližších sousedů k bodu \mathbf{x}
- velikost okénka → závislá na trénovacích datech
 - velká hustota okolo \mathbf{x} → malé okénko → dobré rozlišení
 - malá hustota okolo \mathbf{x} → široké okénko (zastaví se, až narazí na nějakou oblast s velkou hustotou)
- odhad hustoty $p_n(\mathbf{x}) \cong \frac{k_n/n}{V_n}$

k_n -nejbližších sousedů – postřehy

- postřehy
 - $p_n(\mathbf{x})$ spojitá \rightarrow gradient není spojitý
 - body nespojitosti gradientu \rightarrow nemusí být na pozici trénovacích vzorů



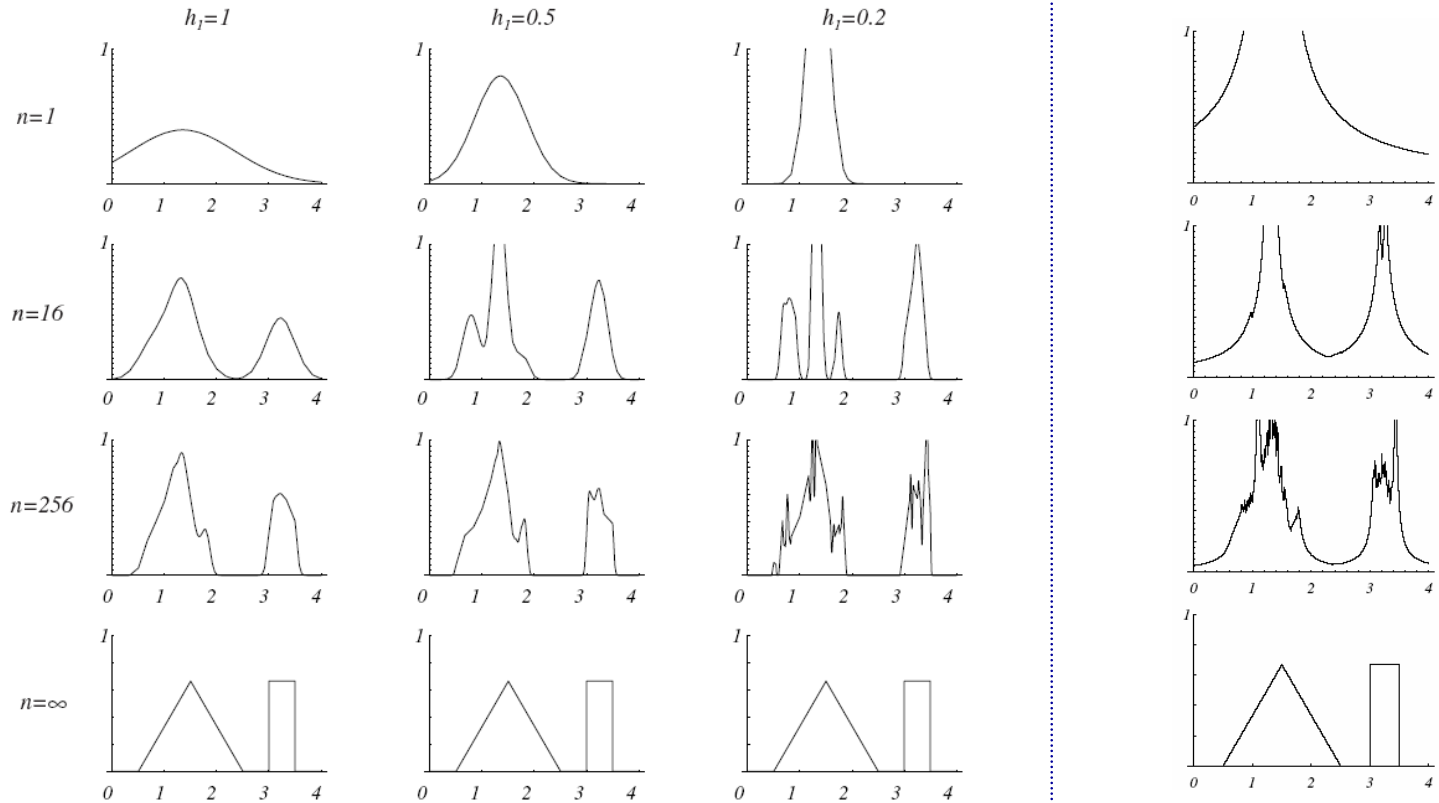
odhad hustoty pomocí 3- a 5- nejbližších sousedů
(dim=1)



odhad hustoty pomocí 5-nejbližších sousedů
(dim=2)

k_n -nejbližších sousedů a Parzenovo okénko

- porovnání metod odhadu
 - uniformní a trojúhelníková hustota pro $n = 1, 16, 256, \infty$



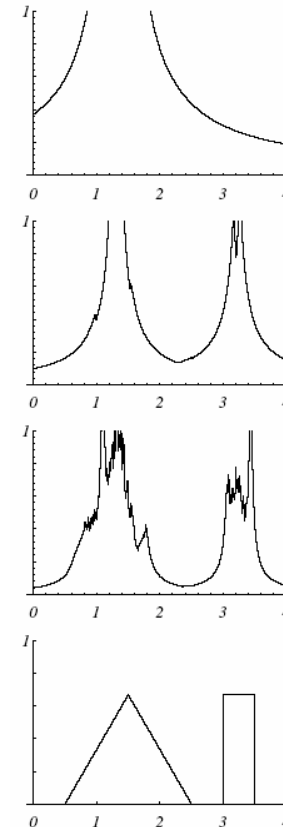
$k_n = \sqrt{n}$

Parzenovo okénko

k_n -nejbližších sousedů

k_n -nejbližších sousedů – postřehy

- $n=1$ a $k_n=\sqrt{n}=1$
 - odhad hustoty $p_n(x) = 1/(2|x-x_1|)$
→ špatný odhad
- rostoucí n
 - odhad se zlepšuje ale $\int p_n(x) = \infty$ (ale má být 1)
 - **výhoda:** $p_n(x)$ se nikdy nepřiblíží k 0
 - **použití:** mnohorozměrné prostory



k_n -nejbližších sousedů – klasifikace

- klasifikace do 2 tříd ω_1 a ω_2
 - 1. předložení vzoru \mathbf{x}
 - 2. N_1 a N_2 ... počet trénovacích vzorů pro ω_1 a ω_2 ($N = N_1 + N_2$)
 - 3. nalezení poloměrů (hyper)koulí
 - r_1 ... poloměr (hyper)koule se středem v \mathbf{x} , která obsahuje k_n trénovacích bodů z ω_1
 - r_2 ... poloměr (hyper)koule se středem v \mathbf{x} , která obsahuje k_n trénovacích bodů z ω_2
 - 4. spočtení objemů V_1 a V_2 (hyper)koulí
 - 5. klasifikace vzoru \mathbf{x} do ω_1 (ω_2)

$$\frac{k_n N_2 V_2}{k_n N_1 V_1} > (<) \frac{P(\omega_2) \lambda_{21} - \lambda_{22}}{P(\omega_1) \lambda_{12} - \lambda_{11}}$$
$$\frac{V_2}{V_1} > (<) \frac{N_1 P(\omega_2) \lambda_{21} - \lambda_{22}}{N_2 P(\omega_1) \lambda_{12} - \lambda_{11}}$$

- poznámky
 - hodnota k_n ... obecně různá pro jednotlivé třídy
 - hledání nejbližších sousedů
 - Eukleidovská vzdálenost → (hyper)koule
 - Mahalanobisova vzdálenost → (hyper)elipsoidy

Odhad apost. pravděpodobnosti

- odhad a posteriori pravděpodobnosti $P(\omega_i|\mathbf{x})$
 - vstup
 - n trénovacích vzorů
 - umístění okénka
 - okénko o velikosti V okolo \mathbf{x} , aby mělo k trénovacích vzorů
 - odhad $P_n(\omega_i|\mathbf{x})$
 - k_i z k trénovacích vzorů patří do $\omega_i \rightarrow$ odhad $P_n(\omega_i|\mathbf{x}) = k_i / k$
 - klasifikace
 - třída, která je nejvíce reprezentována v daném okénku
- \rightarrow dostatek trénovacích vzorů + dostatečně malé okénko \rightarrow dobré výsledky

Metody odhadu hustoty – shrnutí

- metody na odhad hustoty

- Parzenovo okénko

- definuje se, jak se má zvětšovat objem V_n okénka
 - V_n je funkcí n

např. $V_n = 1/\sqrt{n}$

- k_n -nejbližších sousedů

- definuje se, kolik vzorů k_n má okénko obsahovat $\rightarrow V_n$ se zvětšuje, dokud nemá k_n vzorů
 - k_n je funkcí n

např. $k_n = \sqrt{n}$

Pravidlo k -nejbližších sousedů

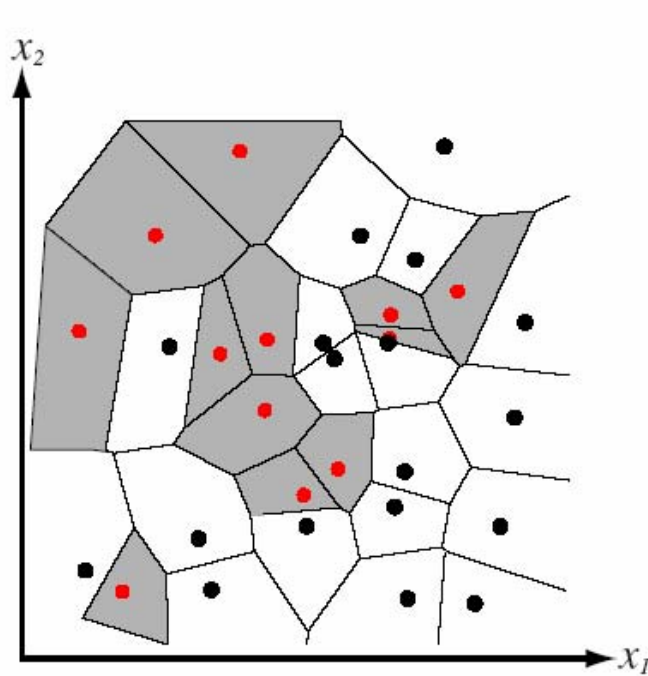
Pravidlo nejbližšího souseda

- pravidlo nejbližšího souseda
 - vstup
 - $D^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$... „označkové prototypy“
 - \mathbf{x} ... neznámý vzor \rightarrow klasifikován
 - výstup
 - \mathbf{x}' ... prototyp z D^n nejbližší k $\mathbf{x} \rightarrow$ klasifikace \mathbf{x} do třídy, kam patří \mathbf{x}'
- vlastnosti
 - suboptimální
 - typicky
 - větší chyba pravděpodobnosti než minimální možná chyba
 - lze ukázat
 - pro $n \rightarrow \infty$: $\text{pravděpodobnost_chyby} \leq 2 * \text{Bayesovská_pravděpodobnost_chyby}$

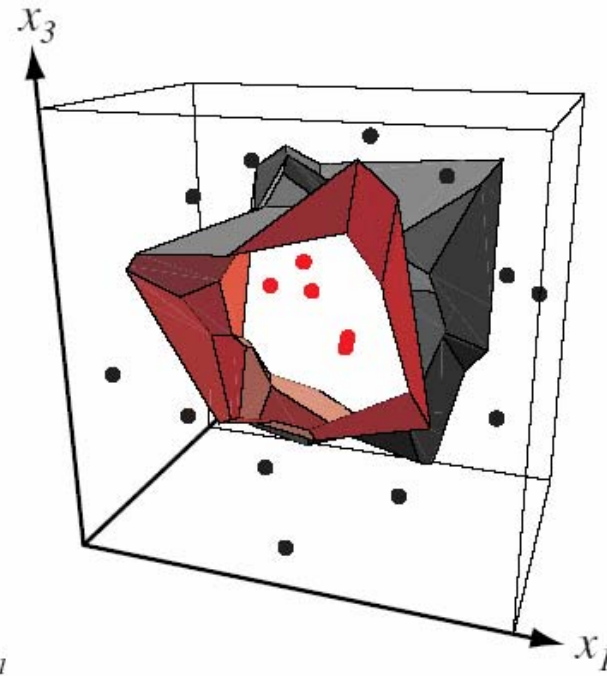
Pravidlo nejbližšího souseda – proč funguje?

- pravidlo nejbližšího souseda
 - θ' ... značka třídy, kam patří prototyp \mathbf{x}'
 - pravděpodobnost $\theta' = \omega_i$
 - aposteriorní pravděpodobnost $P(\omega_i | \mathbf{x}')$
 - velký počet vzorů
 - lze předpokládat, že \mathbf{x}' dostatečně blízko k \mathbf{x}
 - $P(\omega_i | \mathbf{x}') \approx P(\omega_i | \mathbf{x})$
- rozdělení příznakového prostoru do oblastí
 - oblast $R_i = \{ \mathbf{x} ; d(\mathbf{x}, \mathbf{x}_i) \leq d(\mathbf{x}, \mathbf{x}_k), i \neq k \}$
 - body v oblasti označovány třídou „jejich“ trénovacího vzoru
 - Voronoiova mozaika

Pravidlo nejbližšího souseda – Voronoiiova mozaika



n -úhelníkové oblasti



3D-krystaly

Pravidlo nejbližšího souseda – vlastnosti 1

- vlastnosti pravidla
 - ω_i ... nejpravděpodobnější třída \mathbf{x}
$$P(\omega_i | \mathbf{x}) = \max_k P(\omega_k | \mathbf{x})$$
 - $P(\omega_i | \mathbf{x}) \sim 1$
 - pravidlo nejbližšího souseda \approx Bayesovský výběr
 - minimum chybné pravděpodobnosti malé \rightarrow pravděpodobnost chyby při metodě nejbližšího souseda také malá
 - $P(\omega_i | \mathbf{x}) \sim 1/c$
 - pravidlo nejbližšího souseda \neq Bayesovský výběr
 - pravděpodobnosti chyby přibližně stejná u obou metod

Pravidlo nejbližšího souseda – vlastnosti 2

- $n \rightarrow \infty$

- značení

$$P = \lim_{n \rightarrow \infty} P_n$$

průměrná pravděpodobnost chyby
na n trénovacích vzorech

- lze dokázat

$$P^* \leq P \leq P^* \left(2 - \frac{c}{c-1} P^* \right) \leq 2P^*$$

Bayesovská chyba

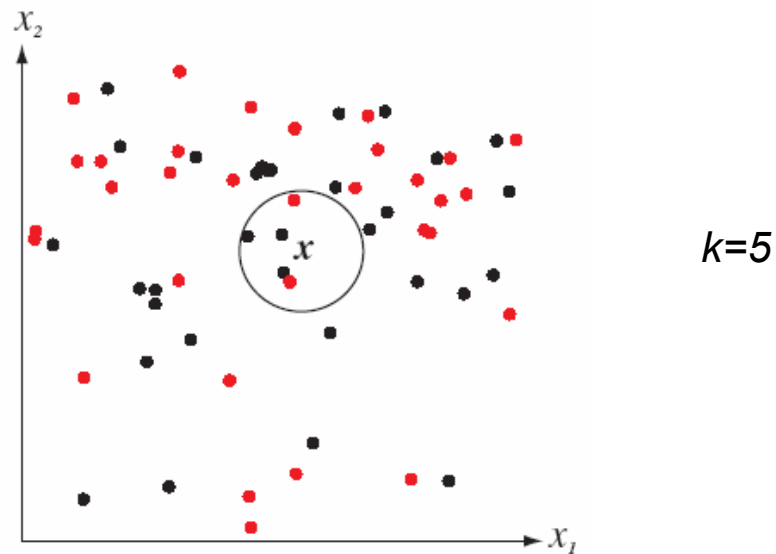
počet tříd

- $n < \infty$

- obecný případ \rightarrow není znám žádný vztah
 - konvergence může být libovolně pomalá
 - chyba P_n se nemusí monotónně snižovat
- obecně nelze analyticky spočítat
 - nutné další předpoklady o pravděpodobnostní struktuře problému

Pravidlo k -nejbližších sousedů

- zobecnění metody
 - vzor x klasifikován do třídy **nejvíce zastoupené** mezi k nejbližšími trénovacími vzory
 - nalezení k nejbližších sousedů (trénovacích vzorů)
 - zjištění tříd nejbližších sousedů
 - vítězství třídy s „největším počtem hlasů“



Zlepšení výpočetní složitosti nejbližšího souseda

- literatura
 - mnoho analýz ohledně výpočetní složitosti (pro $dim=1$ a $dim=2$)
- 2 základní techniky
 - 1. metoda částečné vzdálenosti
 - 2. eliminace zbytečných prototypů

Metoda částečné vzdálenosti – 1

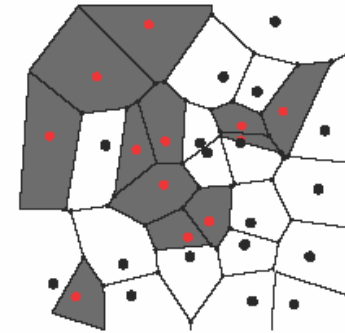
- vstup
 - n trénovacích vzorů (prototypů) dimenze d
 - konstanta r ($r < d$)
- výpočet
 - částečná vzdálenost $d_r(\mathbf{x}, \mathbf{x}')$ prvních r příznaků
 - $d_r(\mathbf{x}, \mathbf{x}') > d(\mathbf{x}, \mathbf{x}_{\text{dosud_nejbližší}}) \rightarrow$ výpočet končí



celá vzdálenost

Eliminace zbytečných prototypů – 2

- metoda
 - eliminace prototypů, v jejichž okolí jsou jen prototypy ze stejné třídy
 - rozhodovací hranice i celková chyba se nezmění



- algoritmus
 1. begin
 2. inicializace $j \leftarrow 0$, $D \leftarrow \text{trénovací_data}$, $n \leftarrow \text{počet_prototypů}$
 3. vytvoření úplného Voronoiova diagramu z D
 4. do $j \leftarrow j+1$
 5. nalezení Voronoiovy sousedy pro prototyp x_j'
 6. if (některý soused z jiné třídy než x_j') then označení x_j'
 7. until $j=n$
 8. odstranění všech neoznačených prototypů
 9. vytvoření Voronoiova diagramu ze zbylých (označených) prototypů
 10. end

Eliminace zbytečných prototypů – 2

- prototyp zůstává
 - přispívá-li k rozhodovací hranici → aspoň jeden z jeho sousedů patří k jiné třídě
- vlastnosti
 - negarantuje minimální množinu prototypů (–)
 - sníží výpočetní složitost (+)
 - bez změny přesnosti výpočtu
 - nelze dodatečně přidávat prototypy do „vyčištěného“ modelu (–)
 - k vyčištění je potřeba znalost **všech** trénovacích dat

Metriky pro metody nejbližších sousedů – 1

- metrika

- nezápornost: $D(\mathbf{a}, \mathbf{b}) \geq 0$
- reflexivita: $D(\mathbf{a}, \mathbf{b}) = 0$ jen když $\mathbf{a} = \mathbf{b}$
- symetrie: $D(\mathbf{a}, \mathbf{b}) = D(\mathbf{b}, \mathbf{a})$
- trojúhelníková nerovnost: $D(\mathbf{a}, \mathbf{b}) + D(\mathbf{b}, \mathbf{c}) \geq D(\mathbf{a}, \mathbf{c})$

- nejběžnější metriky

- Eukleidovská metrika

$$D(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^d (\mathbf{a}_k - \mathbf{b}_k)^2 \right)^{1/2}$$

- Minkowského metrika

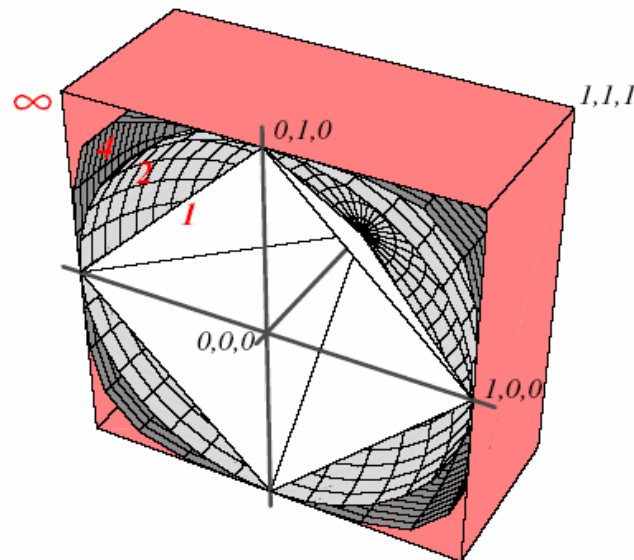
$$L_k(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^d |\mathbf{a}_k - \mathbf{b}_k|^k \right)^{1/k}$$

Metriky pro metody nejbližších sousedů – 2

- Minkowského metrika $L_k(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^d |a_k - b_k|^k \right)^{1/k}$
 - $k=1$... Manhatannovská vzdálenost
 - $k=2$... Eukleidovská vzdálenost
 - $k=\infty$... maximum z projekcí bodů na jednotlivé souřadnicové osy
 - jak vypadá množina bodů ve vzdálenosti 1 od počátku pro $k=1, k=2, k=\infty$?

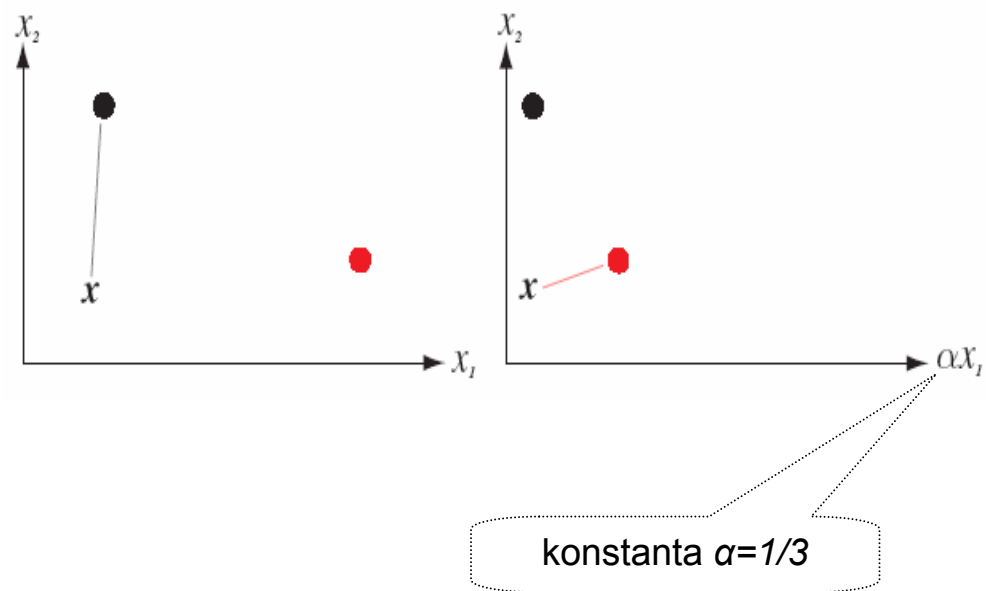
Metriky pro metody nejbližších sousedů – 2

- Minkowského metrika $L_k(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^d |a_k - b_k|^k \right)^{1/k}$
 - $k=1$... Manhatannovská vzdálenost
 - $k=2$... Eukleidovská vzdálenost
 - $k=\infty$... maximum z projekcí bodů na jednotlivé souřadnicové osy
 - jak vypadá množina bodů ve vzdálenosti 1 od počátku pro $k=1$, $k=2$, $k=\infty$?



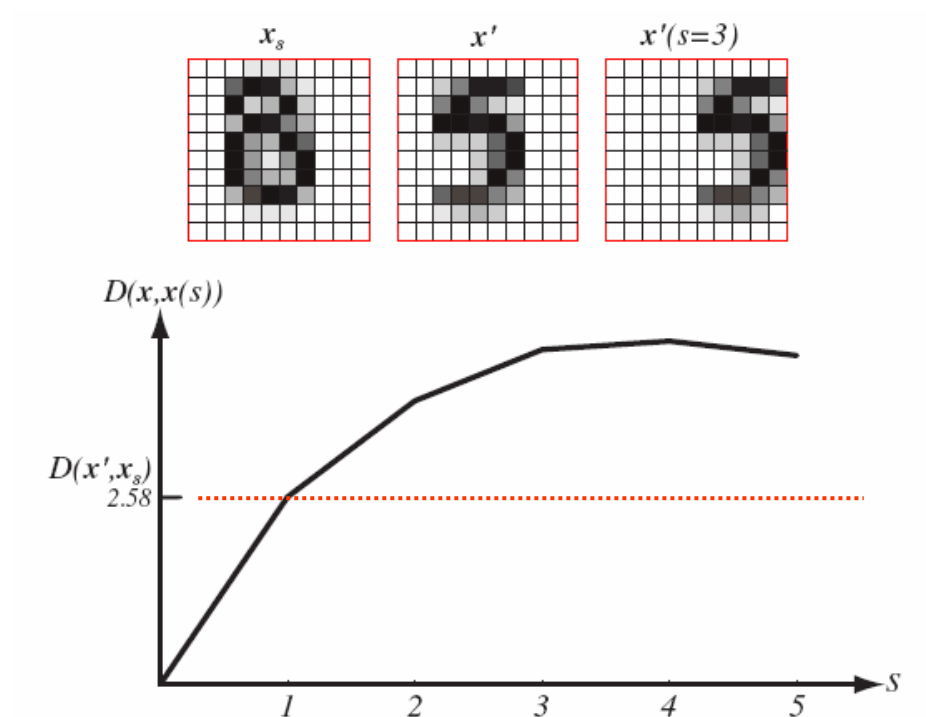
Metriky pro metody nejbližších sousedů – 3

- ne vždy se hodí Eukleidovská metrika
 - např. vynásobení každé souřadnicové osy konstantou
 - odlišné Eukleidovské vzdálenosti v transformovaném prostoru
 - vliv na nalezení nejbližšího souseda



Metriky pro metody nejbližších sousedů – 4

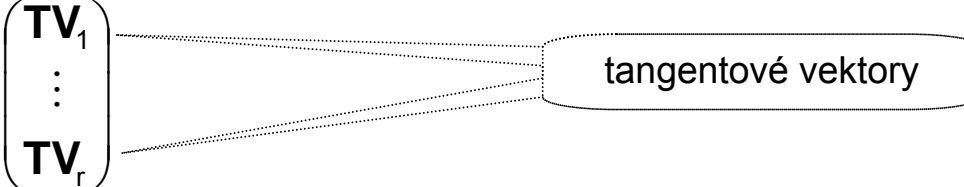
- Eukleidovská metrika – invariance vůči posunutí



$s > 1 \rightarrow L_2(\mathbf{x}', \mathbf{x}'(s)) > L_2(\mathbf{x}', \mathbf{x}_s) \rightarrow$ chybná klasifikace

Tangentová metrika – konstrukce klasifikátoru

- tangentová metrika
 - výpočetně náročná (–)
 - invariantní vůči základním transformacím (+)
 - r transformací
 - např. horizontální a vertikální posun, rotace, škálování, čárové ztenčení,
- konstrukce klasifikátoru
 - pro každý prototyp \mathbf{x}'
 1. aplikace jednotlivých transformací $F_i(\mathbf{x}', \alpha_i)$
 2. vytvoření tangentového vektoru \mathbf{TV}_i pro každou transformaci i
$$\mathbf{TV}_i = F_i(\mathbf{x}', \alpha_i) - \mathbf{x}'$$
 3. uspořádání tangentových vektorů \mathbf{TV}_i do matice $\mathbf{T}(\mathbf{x}')$

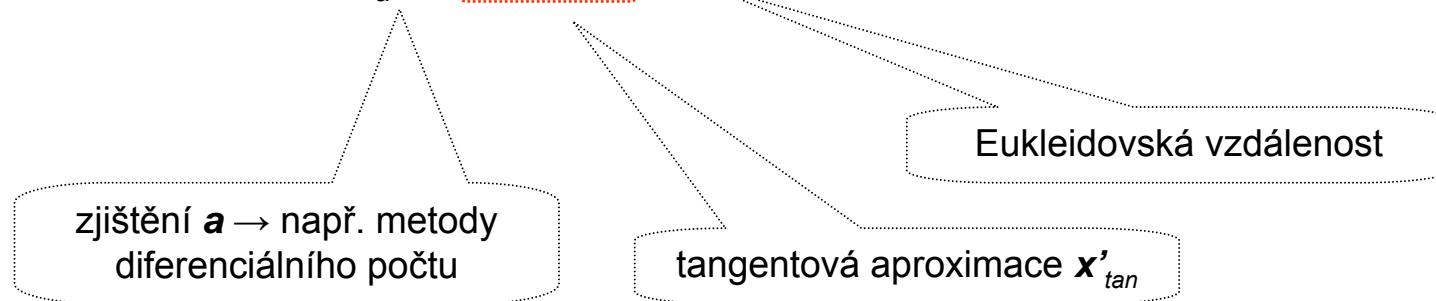
$$\mathbf{T}(\mathbf{x}') = \begin{pmatrix} \mathbf{TV}_1 \\ \vdots \\ \mathbf{TV}_r \end{pmatrix}$$


tangentové vektory

Tangentová metrika – klasifikace

- klasifikace neznámého vektoru \mathbf{x}
 - 1. spočtení tangentové vzdálenosti pro každý prototyp \mathbf{x}'

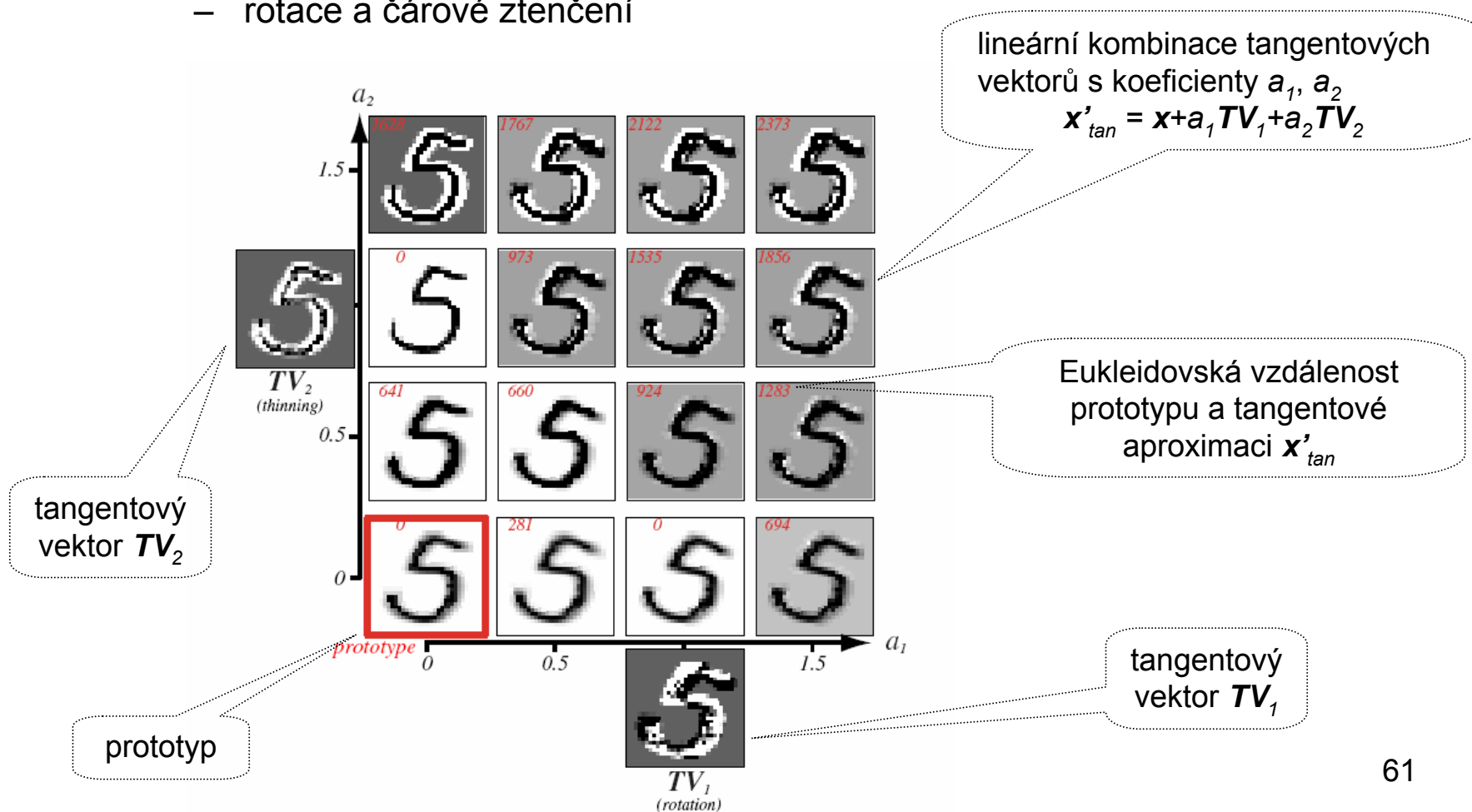
$$D_{\text{tan}}(\mathbf{x}', \mathbf{x}) = \min_a \left\{ \left\| (\mathbf{x}' + T(\mathbf{x}')\mathbf{a}) - \mathbf{x} \right\| \right\}$$



- 2. nalezení prototypu s minimální tangentovou vzdáleností od \mathbf{x}

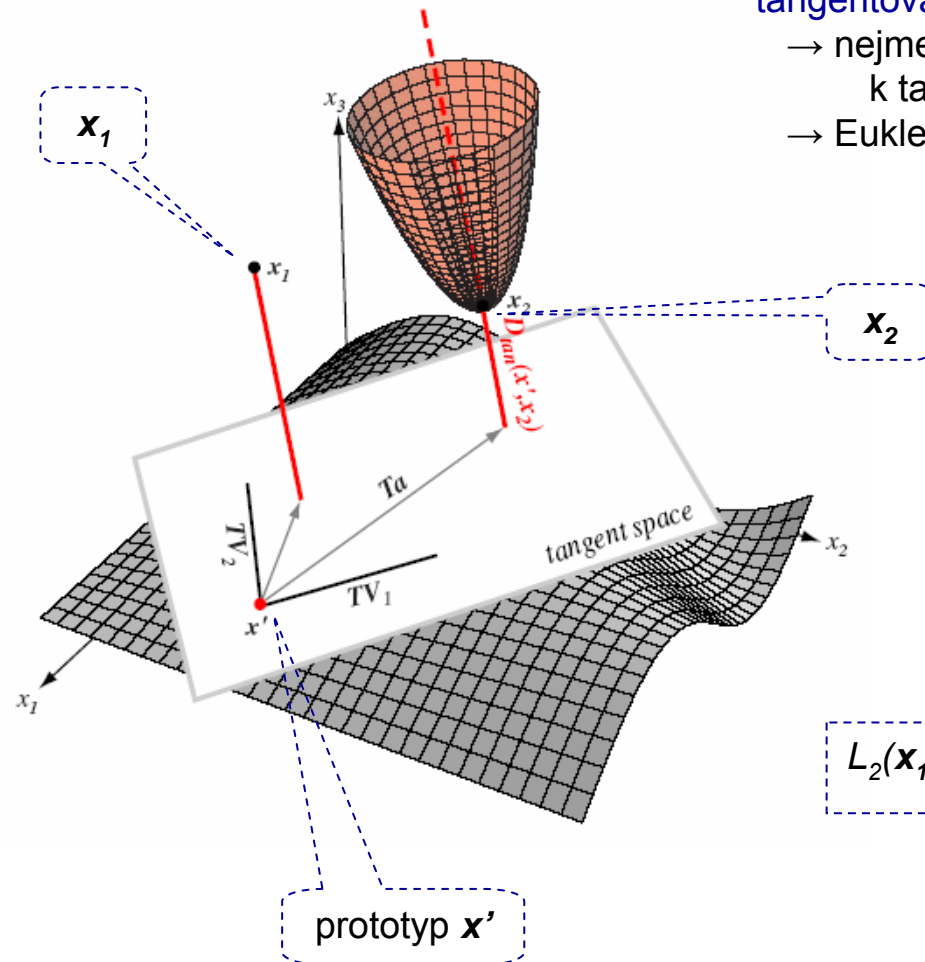
Tangentová metrika – příklad

- 2 transformace
 - rotace a čárové ztenčení



Tangentová metrika – příklad

předloženy vzory \mathbf{x}_1 a \mathbf{x}_2



tangentová vzdálenost $D_{tan}(\mathbf{x}_2, \mathbf{x}') = \min_a \{ \| (\mathbf{x}' + T(\mathbf{x}')\mathbf{a}) - \mathbf{x}_2 \| \}$

→ nejmenší Eukleidovská vzdálenost bodu \mathbf{x}_2
k tangentovému prostoru bodu \mathbf{x}'

→ Eukleid. vzdálenost \mathbf{x}_2 k tangentovému prostoru \mathbf{x}'
→ kvadratická funkce proměnné \mathbf{a}

$$L_2(\mathbf{x}_1, \mathbf{x}') < L_2(\mathbf{x}_2, \mathbf{x}') \text{ ale } D_{tan}(\mathbf{x}_1, \mathbf{x}') > D_{tan}(\mathbf{x}_2, \mathbf{x}')$$