



Pravidlo k -nejbližších sousedů

Jana Štanclová

jana.stanclova@ruk.cuni.cz

LS 2/0 Zk

www.cuni.cz/~stancloj/



Pravidlo k -nejbližších sousedů

- obsah
 - pravidlo nejbližšího souseda
 - vlastnosti
 - pravidlo k -nejbližších sousedů
 - zlepšení výpočetní složitosti
 - metriky



Pravidlo nejbližšího souseda



Pravidlo nejbližšího souseda

- pravidlo nejbližšího souseda
 - vstup
 - $D^n = \{x_1, \dots, x_n\}$... „označkové prototypy“
 - x ... neznámý vzor → klasifikován
 - výstup
 - x' ... prototyp z D^n nejbližší k x → klasifikace x do třídy, kam patří 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?

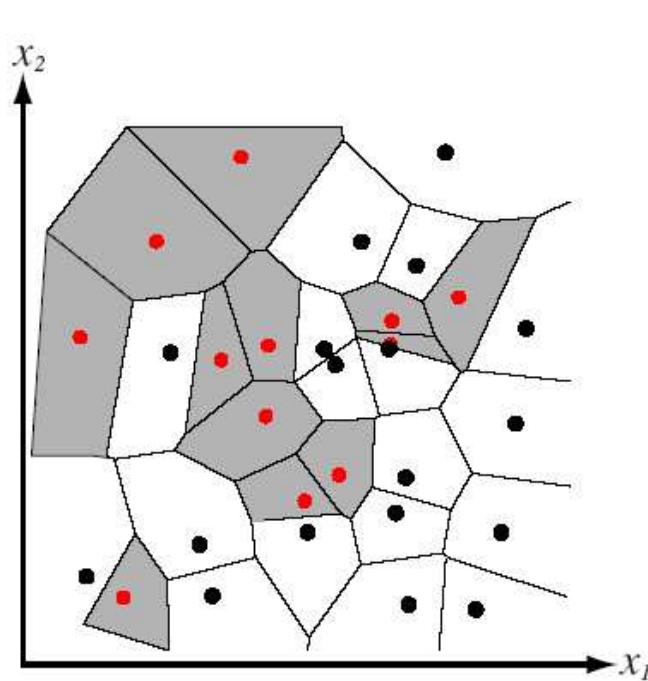
- pravidlo nejbližšího souseda
 - θ' ... značka třídy, kam patří prototyp \mathbf{x}'
 - pravděpodobnost $\theta' = \omega_i$
 - aposteriorní pravděpodobnost $P(\omega_i / \mathbf{x}')$
 - velký počet vzorů
 - lze předpokládat, že \mathbf{x}' dostatečně blízko k \mathbf{x}
 - $P(\omega_i / \mathbf{x}') \approx P(\omega_i / \mathbf{x})$
- 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
 - ??



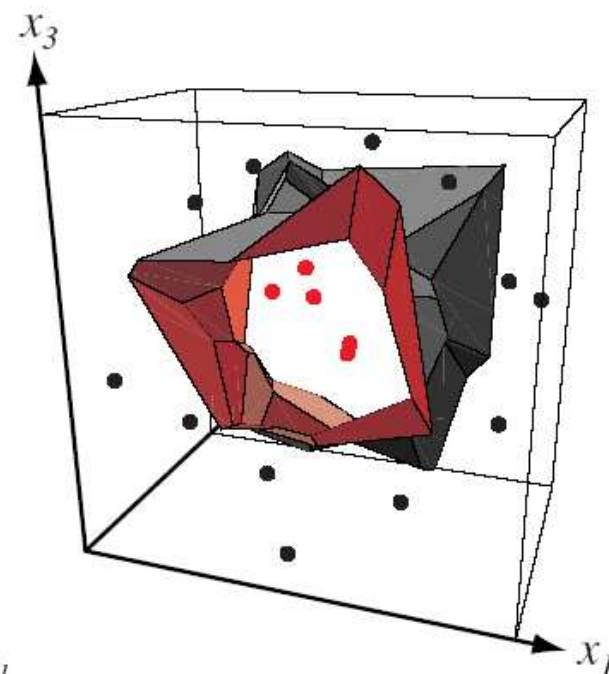
Pravidlo nejbližšího souseda – proč funguje?

- 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 – Voronoiova mozaika



n -úhelníkové oblasti



3D-krystaly



Pravidlo nejblížešího souseda – vlastnosti 1

- vlastnosti pravidla

- $\omega_i \dots$ nejpravděpodobnější třída \mathbf{x}

$$P(\omega_i | \mathbf{x}) = \max_k P(\omega_k | \mathbf{x})$$

- $P(\omega_i | \mathbf{x}) \sim 1$

- pravidlo nejblížešího souseda \approx Bayesovský výběr
- minimum chybné pravděpodobnosti malé \rightarrow pravděpodobnost chyby při metodě nejblížešího souseda také malá

- $P(\omega_i | \mathbf{x}) \sim 1/c$

- pravidlo nejblížešího souseda \neq Bayesovský výběr
- pravděpodobnosti chyby přibližně stejná u obou metod

Pravidlo nejblížešího 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

○ $n < \infty$

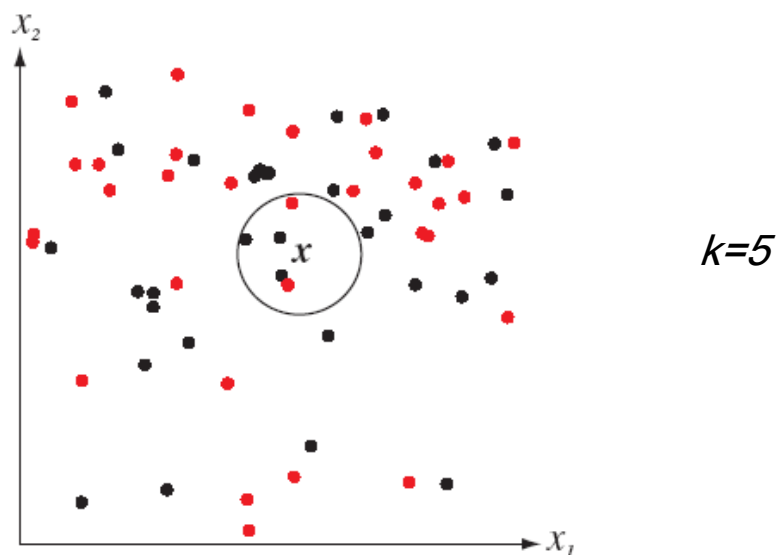
- obecný případ \rightarrow není znám žádný vztah
 - konvergence může být libovolně pomalá
 - chyba P_n se nemusí monotónně snižovat
- obecně nelze analyticky spočítat
 - nutné další předpoklady o pravděpodobnostní struktuře problému



Pravidlo k -nejbližších sousedů

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 $d_r(x, x')$ prvních r příznaků
 - $d_r(x, x') > d(x, x_{\text{dosud_nejbližší}})$ → výpočet končí

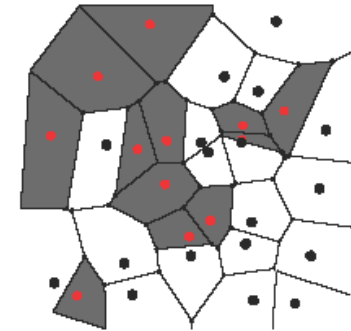
celá vzdálenost

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



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

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

- Minkowského metrika

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

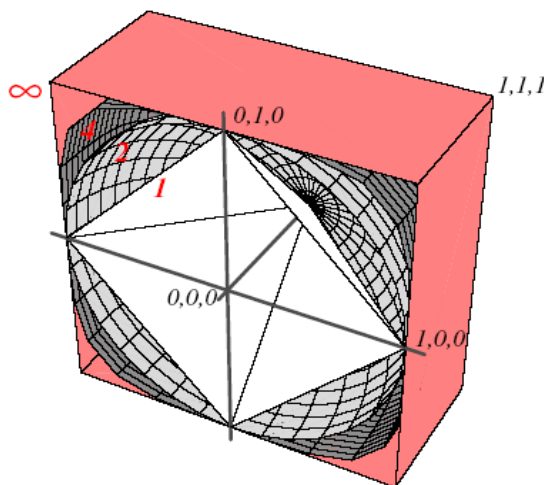


Metriky pro metody nejbližších sousedů – 2

- Minkowského metrika $L_k(\mathbf{a}, \mathbf{b}) = \left(\sum_{i=1}^d |a_i - b_i|^k \right)^{1/k}$
 - $k=1$... 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 $k=1$, $k=2$, $k=\infty$?

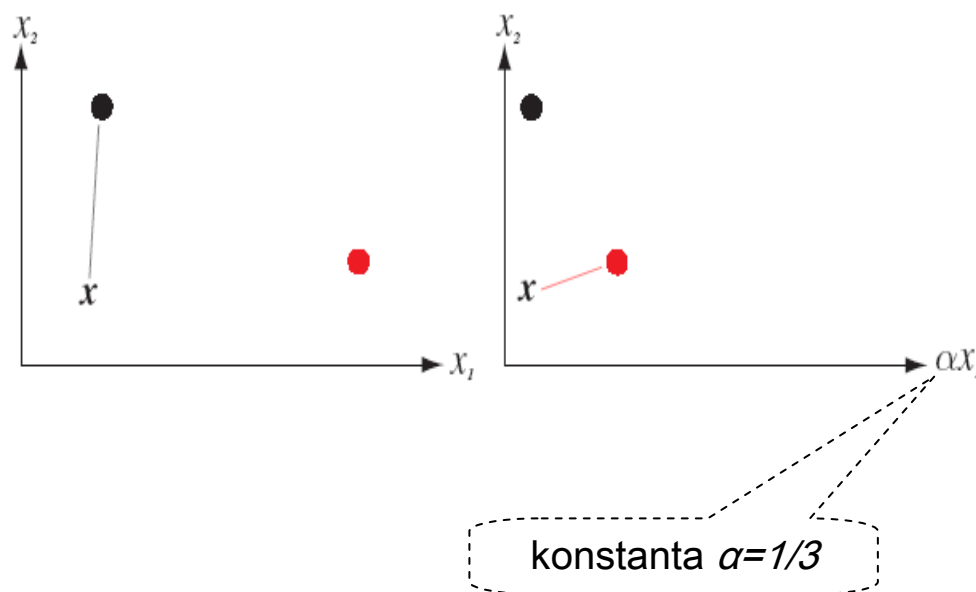
Metriky pro metody nejbližších sousedů – 2

- Minkowského metrika $L_k(\mathbf{a}, \mathbf{b}) = \left(\sum_{i=1}^d |a_i - b_i|^k \right)^{1/k}$
 - $k=1$... 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 $k=1$, $k=2$, $k=\infty$?



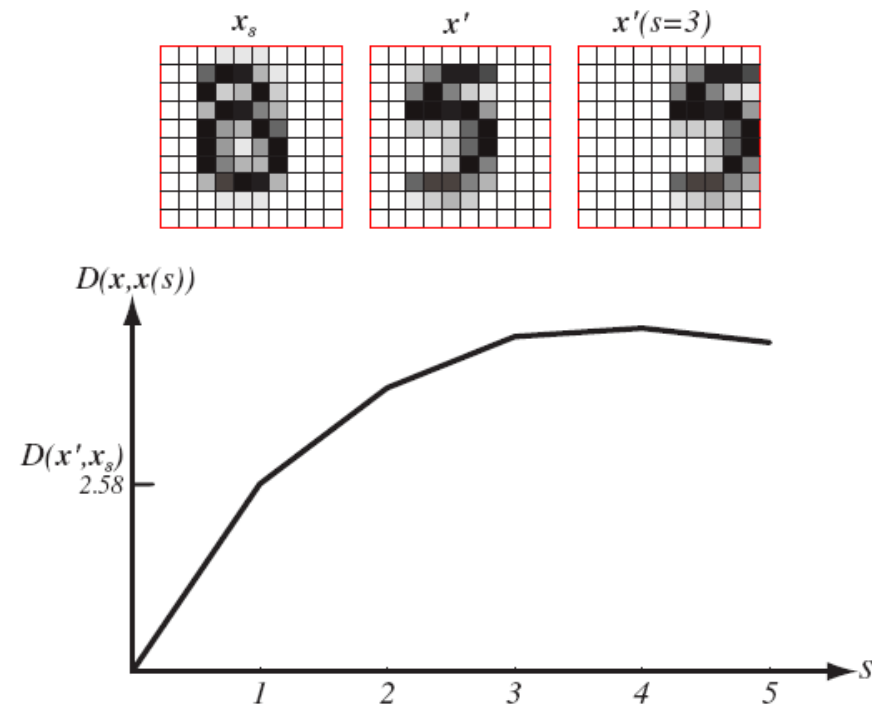
Metriky pro metody nejbližších sousedů – 3

- 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



$$L_2(x', x'(s)) \gg L_2(x', x_s)$$

→ pravidlo nejbližšího souseda s Eukleid. vzdáleností vede k velkým chybám

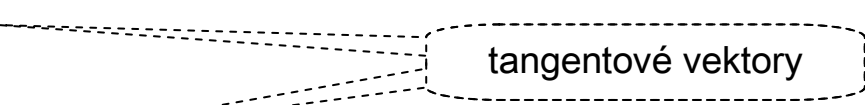
Tangentová metrika – konstrukce klasifikátoru

- 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}'

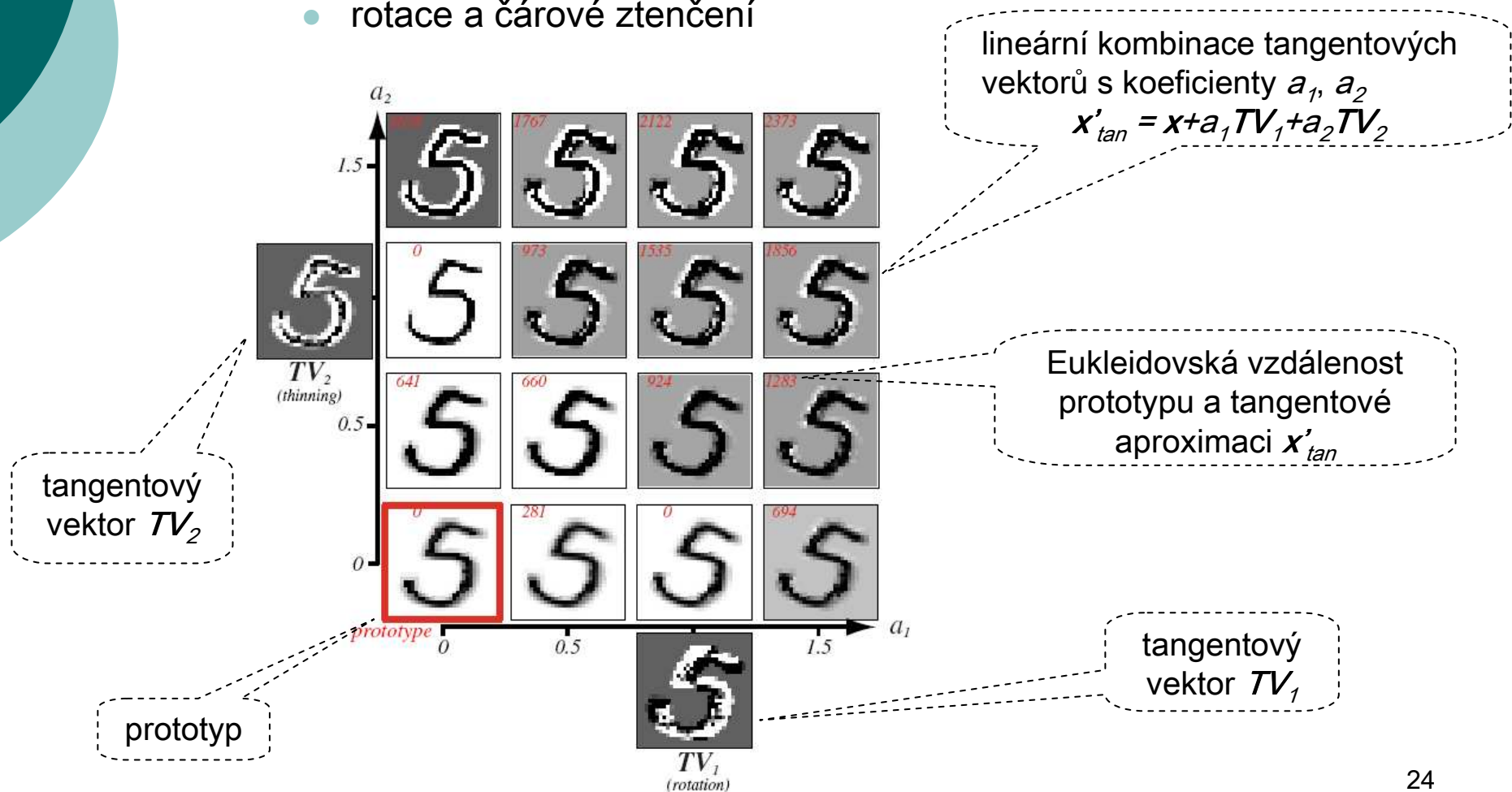
$$D_{\text{tan}}(\mathbf{x}', \mathbf{x}) = \min_a \{ \|(\mathbf{x}' + T(\mathbf{x}')\mathbf{a}) - \mathbf{x}\| \}$$



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

Tangentová metrika – příklad

- 2 transformace
 - rotace a čárové ztenčení



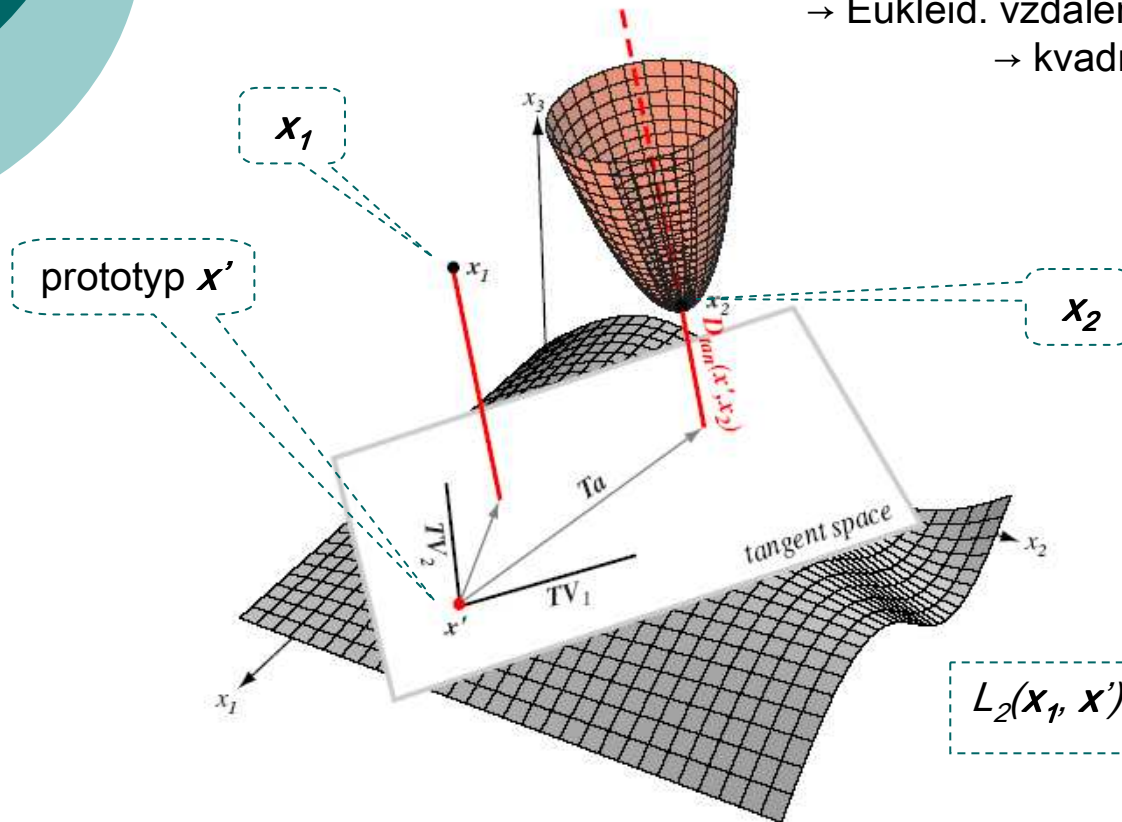
Tangentová metrika – příklad

předloženy vzory x_1 a x_2

tangentová vzdálenost $D_{tan}(x_2, x') = \min_a \{ \| (x' + T(x')a) - x_2 \| \}$

→ nejmenší Eukleidovská vzdálenost bodu x_2
k tangentovému prostoru bodu x'

→ Eukleid. vzdálenost x_2 k tangentovému prostoru x'
→ kvadratická funkce proměnné a



$$L_2(x_1, x') < L_2(x_2, x') \text{ ale } D_{tan}(x_1, x') > D_{tan}(x_2, x')$$