

Neuronové sítě

Shluková analýza

Metrický prostor

Metrický prostor $\langle \mathbb{M}, d \rangle$

$d(\bar{x}, \bar{y})$ – metrika

$\forall x, y, z \in \mathbb{M}$:

- 1 $d(\bar{x}, \bar{y}) \geq 0$
- 2 $d(\bar{x}, \bar{y}) = d(\bar{y}, \bar{x})$
- 3 $d(\bar{x}, \bar{y}) = 0 \Leftrightarrow \bar{x} = \bar{y}$
- 4 $d(\bar{x}, \bar{y}) \leq d(\bar{x}, \bar{z}) + d(\bar{z}, \bar{y})$

Semimetrický prostor

platí 1) – 3)

$d(\bar{x}, \bar{y})$ – semimetrika

Vzdálenosti

Minkowského vzdálenost

$\langle \mathbb{R}^n, d_1 \rangle$	$d_1(\bar{x}, \bar{y}) = \sum_{i=1}^n x_i - y_i $	Manhattanská vzdálenost
$\langle \mathbb{R}^n, d_2 \rangle$	$d_2(\bar{x}, \bar{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$	Euklidovská vzdálenost
$\langle \mathbb{R}^n, d_\infty \rangle$	$d_\infty(\bar{x}, \bar{y}) = \max_{i=1, \dots, n} x_i - y_i $	Čebyševova vzdálenost
$\langle \mathbb{R}^n, d_p \rangle$	$d_p(\bar{x}, \bar{y}) = \left(\sum_{i=1}^n x_i - y_i ^p \right)^{\frac{1}{p}}$	Minkowského vzdálenost
	$p \geq 1 \quad (0 < p < 1 - \text{semimetrika})$	

Vzdálenosti

Hammingova vzdálenost

$$\langle \mathbb{M}^n, d_H \rangle$$

$$d_H(\alpha, \beta) = |\{i \in \hat{n} \mid a_i \neq b_i\}|, \quad d_H(\alpha, \beta) \in \langle 0, n \rangle$$

Levenshteinova vzdálenost

$$\langle \mathbb{M}^*, d_L \rangle$$

$$d_L(\alpha, \beta) = n_I + n_D + l n_R, \quad l \in \langle 1, 2 \rangle$$

$$|n_1 - n_2| \leq d_L(\alpha, \beta) \leq \max(n_1, n_2) \quad |\alpha| = n_1, |\beta| = n_2$$

$$n_1 = n_2 \Rightarrow d_L(\alpha, \beta) \leq d_H(\alpha, \beta)$$

Střed shluku

Shluk (cluster) $\left\{ \begin{array}{l} \text{jak vybrat reprezentanta } x \in \mathbb{C} \text{ shluku } \mathbb{C} \\ \text{jak vybrat obecného reprezentanta shluku } \mathbb{C} \end{array} \right.$

Střed \mathbb{C} v $\langle \mathbb{M}, d \rangle$

$$f(x) = \sum_{k=1}^m d^2(x, x_k), \quad x \in \mathbb{M}, \quad x_k \in \mathbb{C} \subset \mathbb{M}$$

Každý prvek $x \in \mathbb{M}$, pro který platí $f(x) = \min$, nazýváme středem shluku, $x \in s(\mathbb{C})$.

Souvislosti v \mathbb{R}^n

- $\langle \mathbb{R}, d_1 \rangle$: $s(\mathbb{C}) = \{\mu\}$

$$\mu = \frac{1}{m} \sum_{k=1}^m x_k$$

- $\langle \mathbb{R}^n, d_2 \rangle$: $\bar{s}(\mathbb{C}) = \{\bar{\mu}\}$

$$\bar{\mu} = \frac{1}{m} \sum_{k=1}^m \bar{x}_k$$

Shlukování

shluk – skupina s malým rozptylem a velkou vzdáleností od ostatních shluků

shlukování – disjunktní rozdělení dat na shluky

- rozdělení m objektů do H shluků
- metrický prostor $\langle M, d \rangle$
- $C = \{x_1, x_2, \dots, x_m\} \subset M$
- $C = C_1 \cup C_2 \cup \dots \cup C_H, \quad i \neq j : C_i \cap C_j = \emptyset$

rozptyl shluku $C : \sum_{x \in C} d^2(x, s(C))$

součet přes všechny shluky : $J_e = \sum_{i=1}^H \sum_{x \in C_i} (d^2(x, s(C_i)))$

- kritérium při shlukování, chybí zde ale podmínka velkých vzdáleností od ostatních shluků
- používá se při znalosti počtu shluků H

Algoritmy shlukové analýzy

k -means

- 1 Zvolíme k náhodných bodů (objektů) z M (nebo z C) za středy shluků.
- 2 Přiřadíme každý objekt k nejbližšímu středu shluku.
- 3 Přepočítáme středy shluků.
- 4 Pokud se příslušnost objektů ke shlukům změnila, opakujeme druhý a třetí krok.

Poznámka: V našem případě $k = H$.

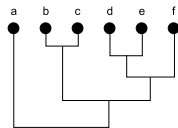
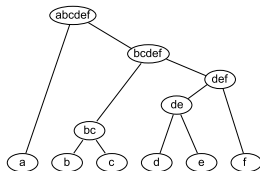
Cíl: Minimalizace kritéria J_e , které určuje kvalitu shlukování.

Algoritmy shlukové analýzy

Hierarchické shlukování

- Není nutné znát předem předpokládaný počet shluků, které data vytvářejí.
- Vycházíme z matice vzdáleností všech objektů.
- Na počátku každý objekt vytváří jeden shluk.
- V každém kroku dochází ke sloučení dvou nejbližších shluků.

dendrogram



Algoritmy shlukové analýzy

Vzdálenost mezi shluky $\mathbb{C}_i, \mathbb{C}_j \subset \mathbb{R}^n$ můžeme určit:

- metodou nejbližšího souseda

$$d(\mathbb{C}_i, \mathbb{C}_j) = \min\{d_2(\bar{x}, \bar{x}') \mid \bar{x} \in \mathbb{C}_i, \bar{x}' \in \mathbb{C}_j, i \neq j\},$$

- metodou nejvzdálenějšího souseda

$$d(\mathbb{C}_i, \mathbb{C}_j) = \max\{d_2(\bar{x}, \bar{x}') \mid \bar{x} \in \mathbb{C}_i, \bar{x}' \in \mathbb{C}_j, i \neq j\},$$

- průměrnou vzdáleností shluků \mathbb{C}_i a \mathbb{C}_j

$$d(\mathbb{C}_i, \mathbb{C}_j) = \frac{1}{m_i m_j} \sum_{\bar{x} \in \mathbb{C}_i} \sum_{\bar{x}' \in \mathbb{C}_j} d_2(\bar{x}, \bar{x}'),$$

- vzdáleností středů shluků \mathbb{C}_i a \mathbb{C}_j

$$d(\mathbb{C}_i, \mathbb{C}_j) = d_2(\bar{\mu}_i, \bar{\mu}_j),$$

- Wardovou metodou minimálního rozptylu

$$d(\mathbb{C}_i, \mathbb{C}_j) = \frac{m_i m_j}{m_i + m_j} d_2^2(\bar{\mu}_i, \bar{\mu}_j).$$

Poznámka: obecně lze vzít místo d_2 jinou metriku, na volbě metriky při shlukování záleží!

Algoritmy shlukové analýzy

Vzdálenost mezi shluky $\mathbb{C}_i, \mathbb{C}_j \subset \mathbb{M}^*$ můžeme určit:

- metodou nejbližšího souseda

$$d(\mathbb{C}_i, \mathbb{C}_j) = \min\{d_L(x, x') \mid x \in \mathbb{C}_i, x' \in \mathbb{C}_j, i \neq j\},$$

- metodou nejvzdálenějšího souseda

$$d(\mathbb{C}_i, \mathbb{C}_j) = \max\{d_L(x, x') \mid x \in \mathbb{C}_i, x' \in \mathbb{C}_j, i \neq j\},$$

- průměrnou vzdáleností shluků \mathbb{C}_i a \mathbb{C}_j

$$d(\mathbb{C}_i, \mathbb{C}_j) = \frac{1}{m_i m_j} \sum_{x \in \mathbb{C}_i} \sum_{x' \in \mathbb{C}_j} d_L(x, x'),$$

Poznámka: V případě $\mathbb{C}_i, \mathbb{C}_j \subset \mathbb{M}^n$ lze použít také Hammingovu vzdálenost.