

GEOINFORMATIKA

Shluková analýza

Algoritmy neřízené klasifikace

Ing. Markéta Potůčková, Ph.D.
Katedra aplikované geoinformatiky a kartografie, Přf UK

Obsah přednášky

- Shlukování
- ***Opakování - Metody neřízené klasifikace v DPZ***
- Používané přístupy ke shlukování
 - Přehled
 - Míry podobnosti
- Příklady algoritmů
 - Sekvenční algoritmy
 - Hierarchické shlukování (dendrogram)
 - k-means
 - DBSCAN

Literatura:

Koutroumbas, K., Theodoridis, S. (2008): *Pattern Recognition*, 4th Edition, Academic Press, 984 p. (online přístup, knihovna UK)

Brunton, S. L., Kutz, J. N. (2019). *Data-driven science and engineering: Machine learning, dynamical systems, and control*. Cambridge University Press.

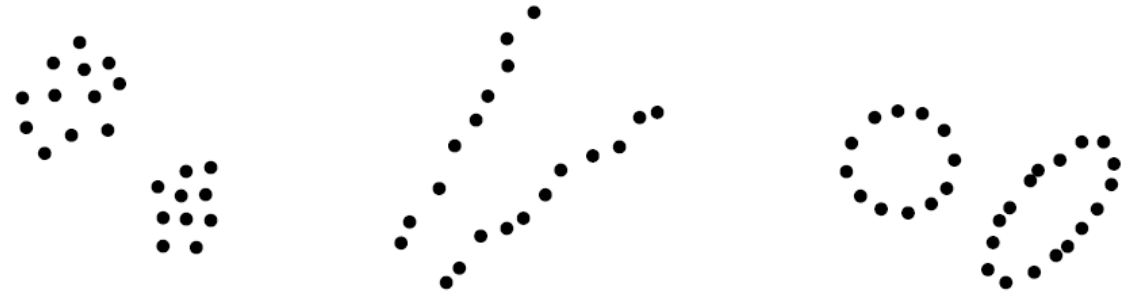
Shlukování

- Metoda identifikace podobných skupin dat v souboru dat
- Při daném počtu shluků, které je třeba vytvořit, se hledá řešení, při kterém jsou všechny prvky v každém shluku co nejpodobnější a všechny shluky se od sebe co nejvíce liší
- Hledání struktur v datech, aniž by tyto byly vázány na konkrétní výsledek

Shlukování

Množina pozorování $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

- \mathbf{x}_i ... vektor příznaků, $i = 1, \dots, L$
- N ... počet pozorování (např. bodů, pixelů, polygonů, textových řetězců, ...)



Shlukování = rozdělení X do m podmnožin C_1, \dots, C_m (shluků) tak, že platí

$$C_i \neq \emptyset, i = 1, \dots, m$$

$$\cup_{i=1}^m C_i = X$$

$$C_i \cap C_j = \emptyset, i \neq j, i, j = 1, \dots, m$$

fuzzy přístup definicí funkce příslušnosti u_j

$$u_j : X \rightarrow [0, 1], \quad j = 1, \dots, m$$

$$\sum_{j=1}^m u_j(\mathbf{x}_i) = 1, \quad i = 1, 2, \dots, N, \quad 0 < \sum_{i=1}^N u_j(\mathbf{x}_i) < N, \quad j = 1, 2, \dots, m$$

- vektory obsažené v shluku C_i jsou si „podobnější“ a jsou „méně podobné“ vektorům příznaků ostatních shluků.
- Kvantifikace „podobnosti“ závisí na tvaru shluku (kruhové, liniové, ...)

Dílčí kroky shlukové analýzy

- Jsou data vhodná ke shlukové analýze, mají-li tendenci vytvářet shluky v daném příznakovém prostoru – shluková analýza nebude fungovat na data zcela náhodného charakteru
- Výběr příznaků (minimalizace redundance/korelace mezi příznaky)
- Míry „blízkosti“/„podobnosti“/„rozdílnosti“
- Rozhodovací kritérium (např. nákladová funkce)
- Algoritmus shlukování
- Validace
- Interpretace výsledku

Míra podobnosti/rozdílnosti (similarity/dissimilarity measures)

A *similarity measure* (SM) s on X is defined as

$$s : X \times X \rightarrow \mathcal{R}$$

such that

$$\exists s_0 \in \mathcal{R} : -\infty < s(\mathbf{x}, \mathbf{y}) \leq s_0 < +\infty, \quad \forall \mathbf{x}, \mathbf{y} \in X$$

$$s(\mathbf{x}, \mathbf{x}) = s_0, \quad \forall \mathbf{x} \in X$$

and

$$s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in X$$

If in addition

$$s(\mathbf{x}, \mathbf{y}) = s_0 \quad \text{if and only if} \quad \mathbf{x} = \mathbf{y}$$

and

$$s(\mathbf{x}, \mathbf{y})s(\mathbf{y}, \mathbf{z}) \leq [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]s(\mathbf{x}, \mathbf{z}), \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X$$

s is called a *metric SM*.

- Příklad DM: Eukleidovská vzdálenost
- DM – hledáme minimum, SM – hledáme maximum

A *dissimilarity measure* (DM) d on X is a function.

$$d : X \times X \rightarrow \mathcal{R}$$

where \mathcal{R} is the set of real numbers, such that

$$\exists d_0 \in \mathcal{R} : -\infty < d_0 \leq d(\mathbf{x}, \mathbf{y}) < +\infty, \quad \forall \mathbf{x}, \mathbf{y} \in X$$

$$d(\mathbf{x}, \mathbf{x}) = d_0, \quad \forall \mathbf{x} \in X$$

and

$$d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in X$$

If in addition

$$d(\mathbf{x}, \mathbf{y}) = d_0 \quad \text{if and only if} \quad \mathbf{x} = \mathbf{y}$$

and

$$d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}), \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X$$

d is called a *metric DM*.

Koutroumbas, K., Theodoridis, S. (2008)

Míry rozdílnosti/podobnosti mezi dvěma body

DM metriky

- Lp norma
$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^l w_i |x_i - y_i|^p \right)^{1/p}$$

- L_1 $d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^l w_i |x_i - y_i|$
- $L_2, w_i=1$ $d_2(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i - y_i)^2}$
- L_2 , vážená $d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T B (\mathbf{x} - \mathbf{y})}$

- L_∞ $d_\infty(\mathbf{x}, \mathbf{y}) = \max_{1 \leq i \leq l} w_i |x_i - y_i|$

SM metriky

- Skalární součin a kosinus úhlu příznakových vektorů

$$s_{\text{inner}}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^l x_i y_i \quad s_{\text{cosine}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

- Pearsonův korelační koeficient
$$r_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}_d^T \mathbf{y}_d}{\|\mathbf{x}_d\| \|\mathbf{y}_d\|}$$

V případě konečného počtu celočíselných hodnot vektorů \mathbf{x}, \mathbf{y} lze vytvořit kontingenční tabulku a z ní spočítat jako DM tzv. Hammingovu vzdálenost (suma nediagonálních prvků)

$$d_H(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{k-1} \sum_{j=0, j \neq i}^{k-1} a_{ij}$$

Míry rozdílnosti/podobnosti mezi bodem a množinou bodů

- Používá se v případě, že je třeba ohodnotit míru shody mezi vektorem \mathbf{x} a shlukem C

A) Hodnotí se vztah \mathbf{x} vůči všem prvkům C

The *max proximity function*:

$$\wp_{\max}^{ps}(\mathbf{x}, C) = \max_{\mathbf{y} \in C} \wp(\mathbf{x}, \mathbf{y})$$

The *min proximity function*:

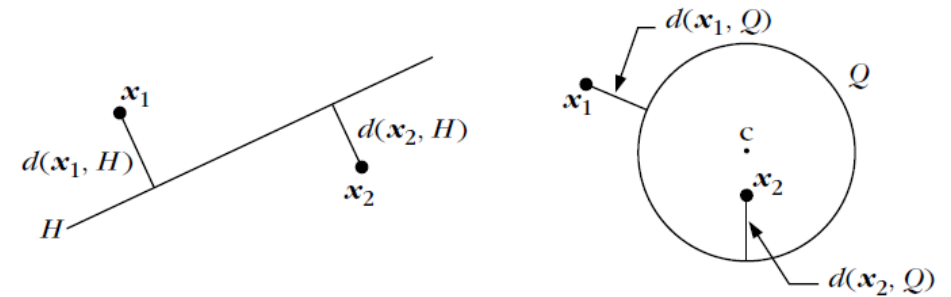
$$\wp_{\min}^{ps}(\mathbf{x}, C) = \min_{\mathbf{y} \in C} \wp(\mathbf{x}, \mathbf{y})$$

The *average proximity function*:

$$\wp_{\text{avg}}^{ps}(\mathbf{x}, C) = \frac{1}{n_C} \sum_{\mathbf{y} \in C} \wp(\mathbf{x}, \mathbf{y})$$

B) Hodnotí se vztah \mathbf{x} vůči reprezentativnímu výběru z C

- průměr, medián
- vzdálenost k nadrovině, „hypersféře“, ...



Koutroumbas, K., Theodoridis, S. (2008)

Kolika způsoby můžeme m shluků vytvořit?

Množina pozorování $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

- \mathbf{x}_i ... vektor příznaků, $i = 1, \dots, N$
- N ... počet pozorování (např. pixelů v obraze)

$S(N, m)$... možnosti shlukování N vektorů do m shluků

$$S(N, 1) = 1$$

$$S(N, N) = 1$$

$$S(N, m) = 0, \text{ for } m > N$$

jinak

$$S(N, m) = \frac{1}{m!} \sum_{i=0}^m (-1)^{m-i} \binom{m}{i} i^N$$

Příklad:

$$S(15, 3) = 2375101$$

$$S(20, 4) = 45232115901$$

$$S(25, 8) = 690223721118368580$$

$$S(100, 5) \simeq 10^{68}$$

Přístupy ke shlukové analýze

- Sekvenční přístup
- Hierarchický přístup (spojování či rozdělování datasetu)
- Optimalizace nákladové funkce
- Další přístupy:
 - Valley-seeking clustering algorithms
 - Density-based algorithms
 - Kernel-based methods

(dle Koutroumbas, K., Theodoridis, S., 2008)

Základní schéma sekvenčního algoritmu shlukování

- Počet tříd není znám, příznakové vektory vstupují do algoritmu pouze jednou
- Parametry definované uživatelem: prahová hodnota DM - Θ , maximální počet shluků - q
- Volba DM, tj. $d(\mathbf{x}, C)$ vede k různým algoritmům
- Co ovlivní výsledek?
 - Pořadí vstupujících vektorů
 - Volba Θ (čím větší, tím menší počet shluků)
- Je-li C reprezentován bodem, algoritmus zvýhodňuje kompaktní (kruhové) shluky
- Protože se očekává, že konečný počet shluků m bude mnohem menší než N , je časová složitost BSAS $O(N)$
- Možné úpravy:
 - Příznakové vektory vstupují do algoritmu 2x – prvně pro vytvoření shluků, podruhé pro konečné přiřazení ke shlukům
 - Volba více prahových hodnot Θ
 - Spojování shluků
 - Změna pořadí vstupních vektorů

Basic Sequential Algorithmic Scheme (BSAS)

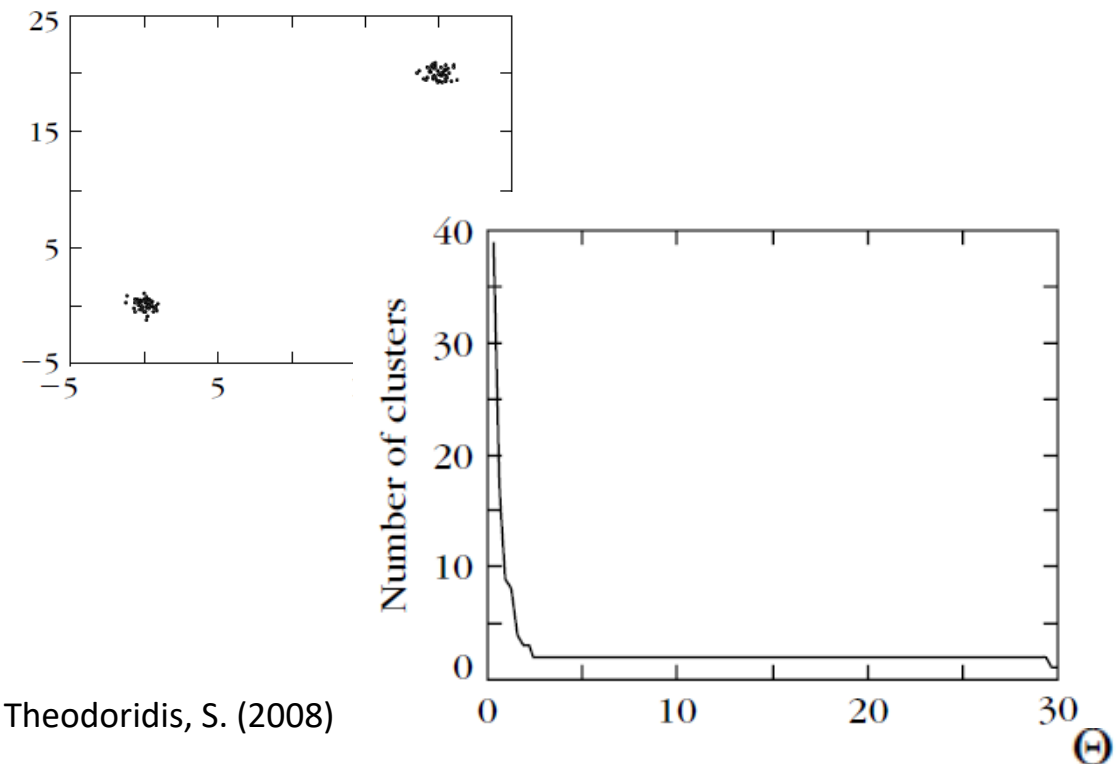
- $m = 1$
- $C_m = \{\mathbf{x}_1\}$
- For $i = 2$ to N
 - Find C_k : $d(\mathbf{x}_i, C_k) = \min_{1 \leq j \leq m} d(\mathbf{x}_i, C_j)$.
 - If $(d(\mathbf{x}_i, C_k) > \Theta)$ AND $(m < q)$ then
 - $m = m + 1$
 - $C_m = \{\mathbf{x}_i\}$
 - Else
 - $C_k = C_k \cup \{\mathbf{x}_i\}$
 - Where necessary, update representatives²
 - End {if}
- End {For}

Koutroumbas, K., Theodoridis, S. (2008)

Optimální počet shluků

- S využitím sekvenčního shlukování
- Hodnoty a , b jsou minimální a maximální hodnoty DM mezi všemi páry vektorů v X , tj. $a = \min_{i,j=1,\dots,N} d(\mathbf{x}_i, \mathbf{x}_j)$, $b = \max_{i,j=1,\dots,N} d(\mathbf{x}_i, \mathbf{x}_j)$.
Volba c je přímo ovlivněna typem $d(\mathbf{x}, C)$.
- Čím větší s , tím větší je statistický vzorek a tím vyšší přesnost výsledků

- For $\Theta = a$ to b step c
 - Run s times the algorithm BSAS(Θ), each time presenting the data in a different order.
 - Estimate the number of clusters, m_Θ , as the most frequent number resulting from the s runs of BSAS(Θ).
- Next Θ



Koutroumbas, K., Theodoridis, S. (2008)

Hierarchické shlukování

Algoritmy zahrnují N kroků (počet vstupních vektorů). V každém kroku t se získá nové shlukování na základě shlukování vytvořeného v předchozím kroku $t-1$.

Dvě hlavní kategorie

- Aglomerativní hierarchické algoritmy
 - Vychází z N shluků (= každý vstupní vektor reprezentuje jeden shluk) a ty postupně spojuje, výsledný shluk obsahuje celý dataset X
- Dělicí hierarchické algoritmy
 - Výchozí shluk obsahuje celý dataset X a postupně je dělen až do úrovně dílčích vektorů

Obecný postup aglomerativního hierarchického shlukování

- Pokud se dva vektory spojí do jednoho shluku na úrovni t , zůstanou ve stejném shluku pro všechny následující shlukování
→ neexistuje způsob, jak napravit „špatné“ shlukování, ke kterému mohlo dojít na nižší úrovni hierarchie
- Celkový počet dvojic, které je třeba prozkoumat během celého procesu shlukování, je

$$\frac{(N-1)N(N+1)}{6}$$

tj. výpočetní náročnost algoritmu obecně odpovídá N^3 , ale existují varianty $O(N^2 \log N)$ či $O(N^2)$

- $g(C_i, C_j)$ je řešena na základě matice rozdílnosti (disimilarity matrix) či grafových algoritmů

■ Initialization:

- Choose $\mathfrak{R}_0 = \{C_i = \{\mathbf{x}_i\}, i = 1, \dots, N\}$ as the initial clustering.
- $t = 0$.

■ Repeat:

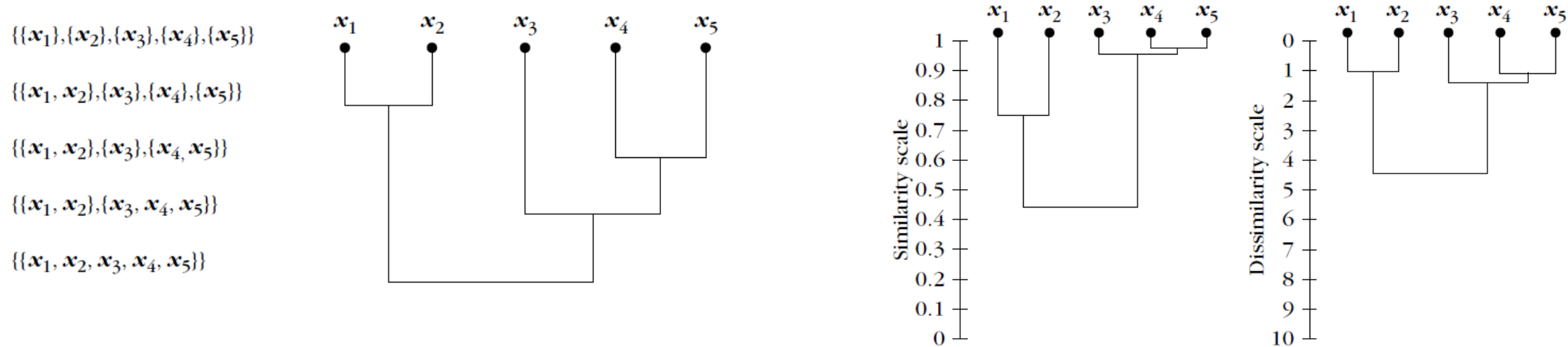
- $t = t + 1$
- Among all possible pairs of clusters (C_r, C_s) in \mathfrak{R}_{t-1} find the one, say (C_i, C_j) , such that

$$g(C_i, C_j) = \begin{cases} \min_{r,s} g(C_r, C_s), & \text{if } g \text{ is a dissimilarity function} \\ \max_{r,s} g(C_r, C_s), & \text{if } g \text{ is a similarity function} \end{cases} \quad (13.1)$$

- Define $C_q = C_i \cup C_j$ and produce the new clustering $\mathfrak{R}_t = (\mathfrak{R}_{t-1} - \{C_i, C_j\}) \cup \{C_q\}$.

■ Until all vectors lie in a single cluster.

Aglomerativní hierarchické shlukování - dendrogram



Která úroveň je dostačující?

Příklad kritéria:

Definujeme DM uvnitř shluku C např. $b_1(C) = \max\{d(\mathbf{x}, \mathbf{y}), \mathbf{x}, \mathbf{y} \in C\}$ $b_2(C) = \text{med}\{d(\mathbf{x}, \mathbf{y}), \mathbf{x}, \mathbf{y} \in C\}$

a prahovou hodnotu θ

Algoritmus bude ukončen na úrovni t , platí-li $\exists C_j \in \mathfrak{R}_{t+1} : b(C_j) > \theta$

Shlukování s optimalizační funkcí

- Nejvyžívanější postup shlukové analýzy
- Založen na optimalizaci nákladové funkce J (funkce vektorů datové sady X) parametrizované neznámým vektorem Θ
- Většina algoritmů vyžaduje znalost počtu shluků m
- Využíváno i pro geometrické úlohy (např. robotika a rozeznávání objektů splňujících dané geometrické parametry – rovina, sféra atd.)
- Různé přístupy:
 - Mixture decomposition – nákladová funkce je konstruována na základě náhodných vektorů a přiřazení do shluků se řídí rozdělením pravděpodobnosti (paralela s Bayesovskou klasifikací); podmíněné pravděpodobnosti jsou výsledkem procesu optimalizace
 - Fuzzy – pracuje s DM/SM jako funkcí „blízkosti“ hodnoceného vektoru a shluku a zavádí funkce příslušnosti
 - Hard - speciální případ přístupu fuzzy shlukování, kde každý vektor patří výhradně do jednoho shluku

Obecné schéma tvrdého klasifikátoru s optimalizační funkcí

- Choose $\theta_j(0)$ as initial estimates for $\theta_j, j = 1, \dots, m$.

- $t = 0$

- Repeat

- For $i = 1$ to N

- For $j = 1$ to m

- *Determination of the partition:*¹⁰

$$u_{ij}(t) = \begin{cases} 1, & \text{if } d(\mathbf{x}_i, \theta_j(t)) = \min_{k=1, \dots, m} d(\mathbf{x}_i, \theta_k(t)) \\ 0, & \text{otherwise,} \end{cases}$$

- End {For- j }

- End {For- i }

- $t = t + 1$

- For $j = 1$ to m

- *Parameter updating:* Solve

$$\sum_{i=1}^N u_{ij}(t-1) \frac{\partial d(\mathbf{x}_i, \theta_j)}{\partial \theta_j} = \mathbf{0}$$

with respect to θ_j and set $\theta_j(t)$ equal to the computed solution.

- End {For- j }.

- Until a termination criterion is met.

Optimalizační nákladová funkce

$$J(\theta, U) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} d(\mathbf{x}_i, \theta_j)$$

$$u_{ij} \in \{0, 1\}, \quad j = 1, \dots, m$$

$$\sum_{j=1}^m u_{ij} = 1$$

Ukončení výpočtu

$$\|\theta(t) - \theta(t-1)\| < \varepsilon$$

Koutroumbas, K., Theodoridis, S. (2008)

ISODATA, k-means

- Hledáme minimum nákladové funkce $J(\Theta, U)$

$$J(\Theta, U) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} \|\mathbf{x}_i - \boldsymbol{\theta}_j\|^2$$

- Algoritmus generuje co nejkompaktnější shluky
- Nízká výpočetní náročnost $O(Nmq)$
 - q ... počet nutných iterací
 - vhodný pro velké datové sady
- Mnoho variant algoritmu
- Algoritmus nemusí vést k nalezení globálního minima funkce $J(\Theta, U)$, různé inicializace algoritmu mohou vést k různým řešením
- Počet shluků je vstupním parametrem
- Otázka optimálního počtu shluků („NP-úplný problém“)
- Citlivý k šumu a odlehlým hodnotám
- Optimalizace středu shluků k-means++

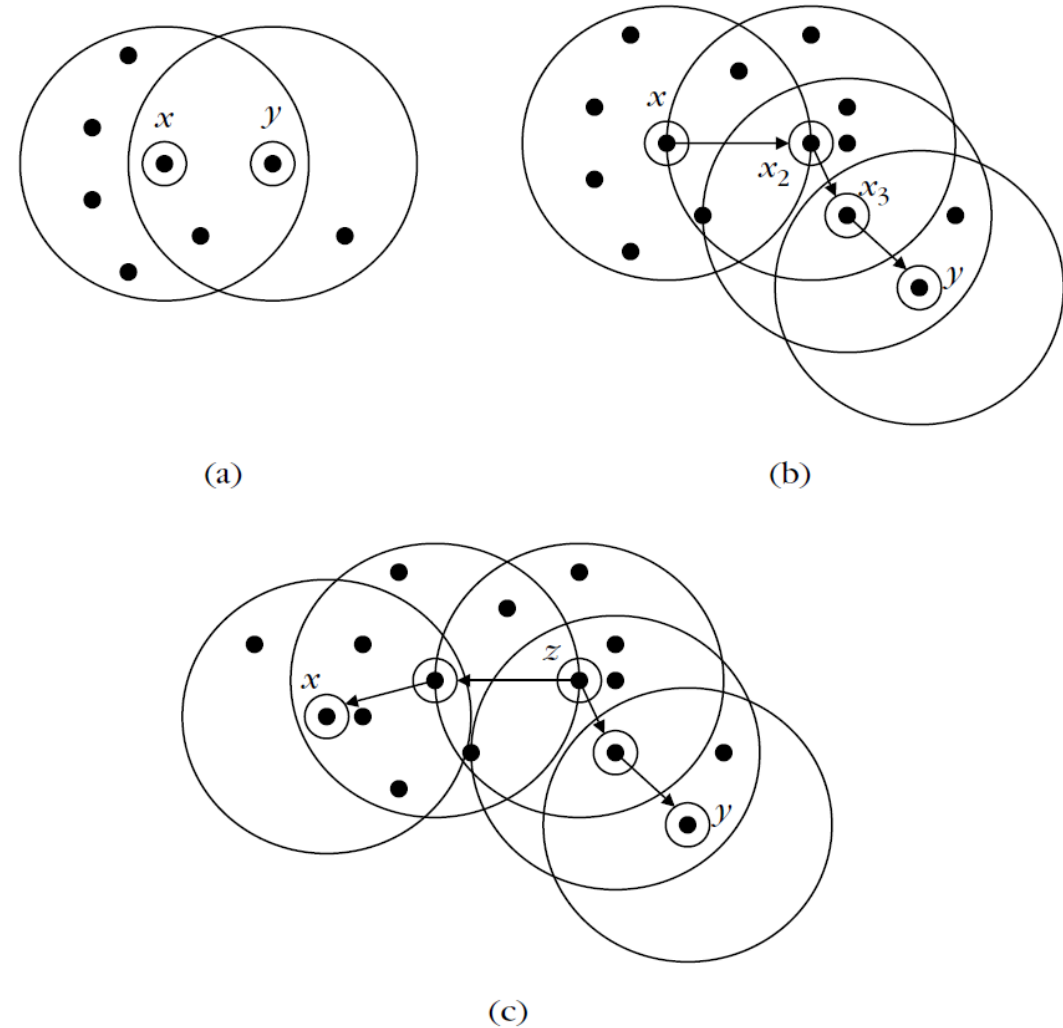
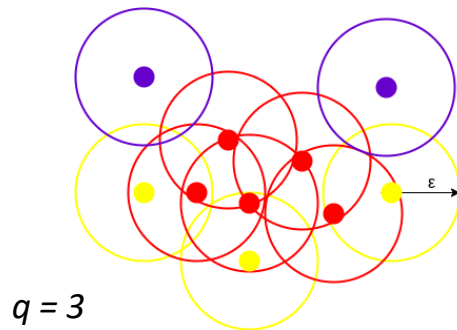
- Choose arbitrary initial estimates $\boldsymbol{\theta}_j(0)$ for the $\boldsymbol{\theta}_j$'s, $j = 1, \dots, m$.
- Repeat
 - For $i = 1$ to N
 - Determine the closest representative, say $\boldsymbol{\theta}_j$, for \mathbf{x}_i .
 - Set $b(i) = j$.
 - End {For}
 - For $j = 1$ to m
 - Parameter updating: Determine $\boldsymbol{\theta}_j$ as the mean of the vectors $\mathbf{x}_i \in X$ with $b(i) = j$.
 - End {For}.
- Until no change in $\boldsymbol{\theta}_j$'s occurs between two successive iterations.

Koutroumbas, K., Theodoridis, S. (2008)

Rozdíl ISODATA, k-means?

DBSCAN

- Density-Based Spatial Clustering of Applications with Noise
- „hustota“ bodu \mathbf{x} = počet bodů X , které se nacházejí v definovaném okolí bodu \mathbf{x}
- Obvykle se uvažuje sférické okolí $V_\varepsilon(\mathbf{x})$ o uživatelem definovaném poloměru ε
- Bod \mathbf{x} označíme jako **vnitřní (jádrový, „core“ bod)** shluku, pokud počet bodů v jeho okolí $N_\varepsilon(\mathbf{x}) > q$ (parametr q volí uživatel)
- **Hraniční bod** – dosažitelný z vnitřního bodu
- Šum



Koutroumbas, K., Theodoridis, S. (2008)

FIGURE 15.22

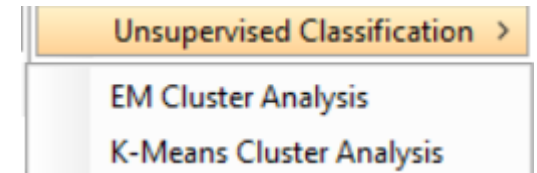
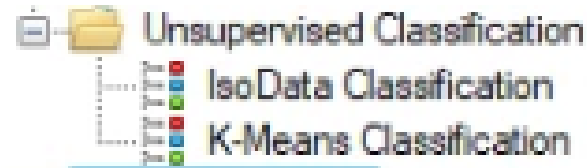
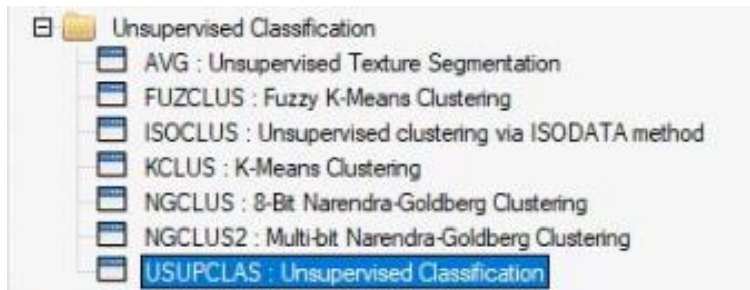
Assuming that $q = 5$, (a) \mathbf{y} is directly density reachable from \mathbf{x} , but not vice versa, (b) \mathbf{y} is density reachable from \mathbf{x} , but not vice versa, and (c) \mathbf{x} and \mathbf{y} are density connected (in addition, \mathbf{y} is density reachable from \mathbf{x} , but not vice versa).

DBSCAN

- Set $X_{un} = X$
- Set $m = 0$
- While $X_{un} \neq \emptyset$ do
 - Arbitrarily select a $\mathbf{x} \in X_{un}$.
 - If \mathbf{x} is a noncore point then
 - Mark \mathbf{x} as noise point.
 - $X_{un} = X_{un} - \{\mathbf{x}\}$
 - If \mathbf{x} is a core point then
 - $m = m + 1$
 - Determine all density-reachable points in X from \mathbf{x} .
 - Assign \mathbf{x} and the previous points to the cluster C_m . The border points that may have been marked as noise are also assigned to C_m .
 - $X_{un} = X_{un} - C_m$
 - End { if }
- End { while }
- The results of the algorithm are greatly influenced by the choice of ε and q . Different values of the parameters may lead to totally different results. One should select these parameters so that the algorithm is able to detect the least “dense” cluster. In practice, one has to experiment with several values for ε and q in order to identify their “best” combination for the data set at hand.
- The DBSCAN is not appropriate for cases where the clusters in X exhibit significant differences in density, and it is not well suited for high-dimensional data ([Xu 05]).

Koutroumbas, K., Theodoridis, S. (2008)

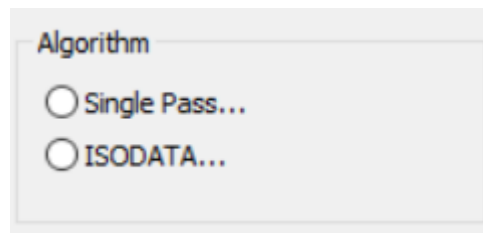
Co nabízejí softwary pro zpracování obrazu DPZ?



Iso Cluster Unsupervised Classification



Multispec



Shlukové analýzy v Matlabu



MATLAB® supports many popular cluster analysis algorithms:

- **Hierarchical clustering** builds a multilevel hierarchy of clusters by creating a cluster tree.
- **k-Means clustering** partitions data into k distinct clusters based on distance to the centroid of a cluster.
- **Gaussian mixture models** form clusters as a mixture of multivariate normal density components.
- **Spatial clustering** (such as the popular density-based DBSCAN) groups points that are close to each other in areas of high density, keeping track of outliers in low-density regions. Can handle arbitrary non-convex shapes.
- **Self-organizing maps** use neural networks that learn the topology and distribution of the data.
- **Spectral clustering** transforms input data into a graph-based representation where the clusters are better separated than in the original feature space. The number of clusters can be estimated by studying eigenvalues of the graph.
- **Hidden Markov models** can be used to discover patterns in sequences, such as genes and proteins in bioinformatics.

Cvičení - Matlab

1. Vygenerujte alespoň 3 množiny bodů/klastry (např. funkce *randn*) ve 2D.
2. Vytvořte vlastní implementaci algoritmu shlukování k-means.
3. Porovnejte svůj výsledek s výsledkem shlukování nad totožnou množinou s využitím funkce *kmeans*. Případné rozdíly komentujte.

Bonusová úloha: rozšiřte své řešení shlukování k-means do n-dimensionálního prostoru.

Děkuji za pozornost

marketa.potuckova@natur.cuni.cz