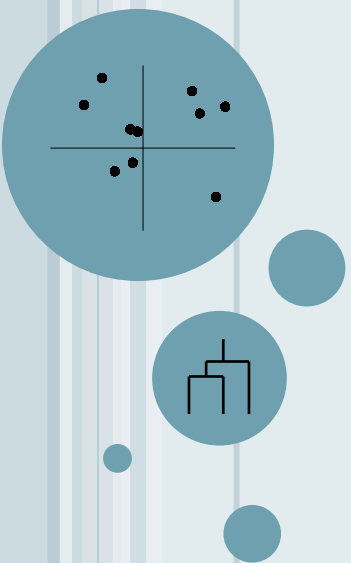


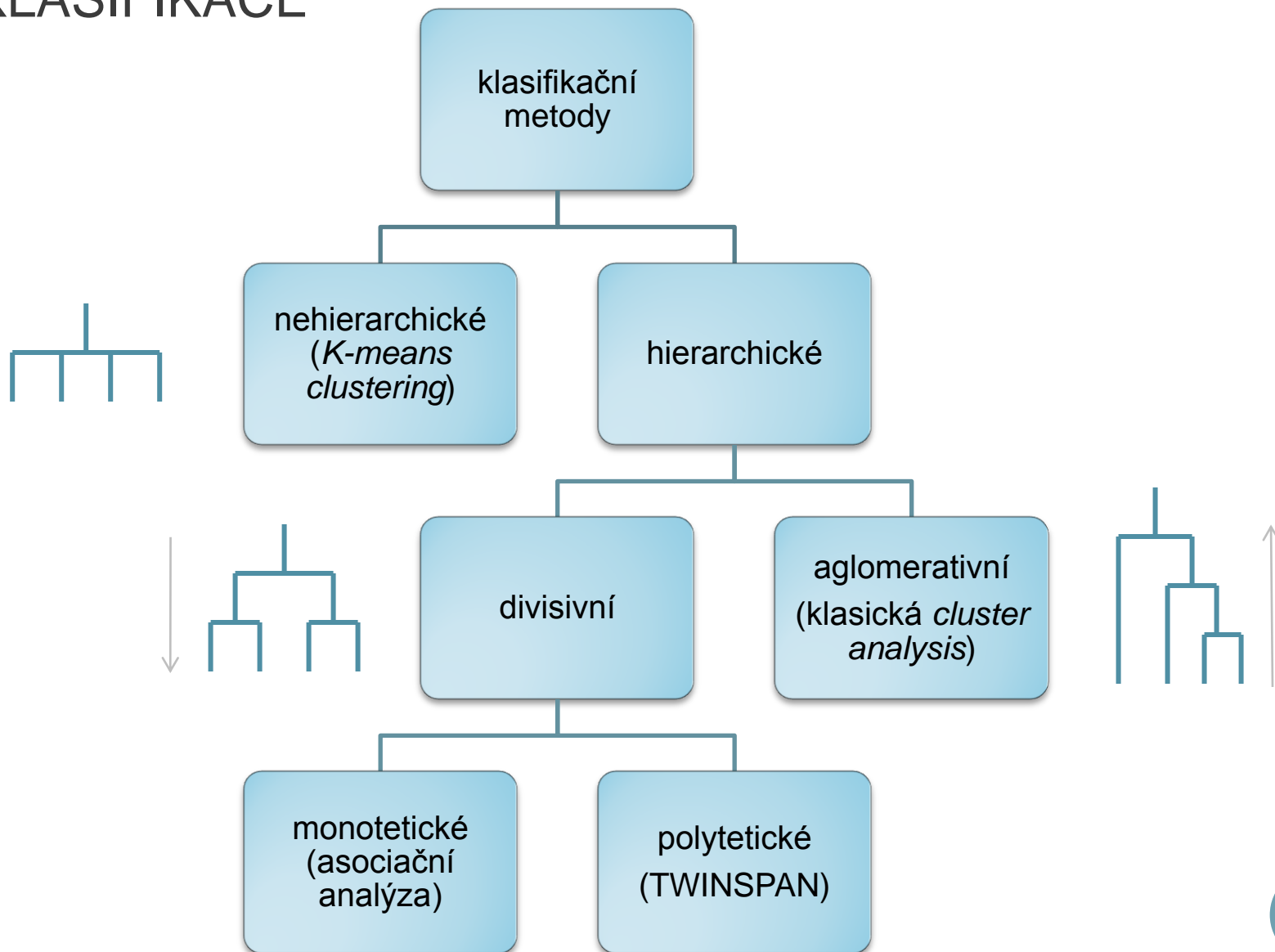
KLASIFIKACE



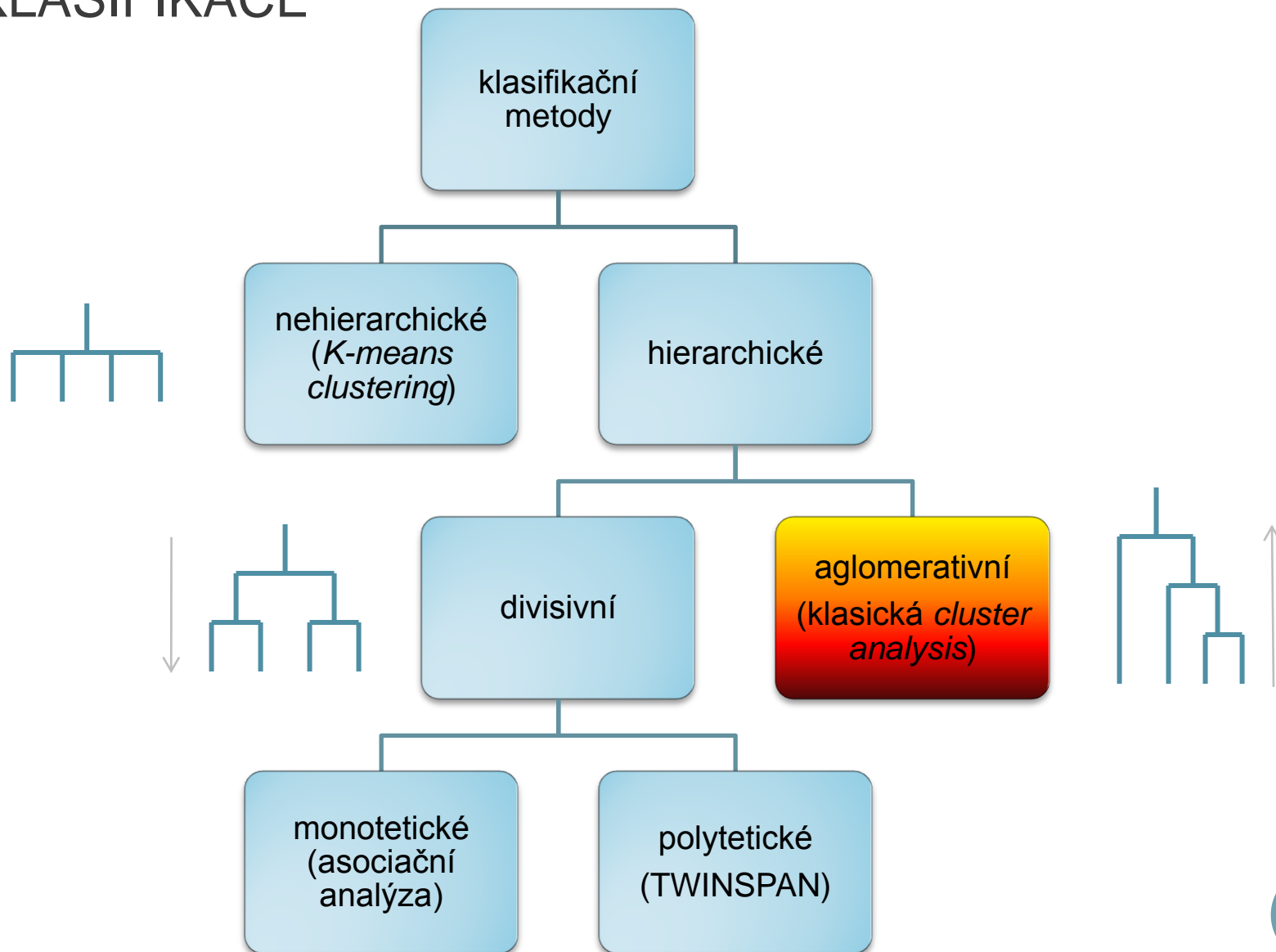
KLASIFIKACE

- cílem je seskupit podobné vzorky (případně druhy) do skupin, které jsou vnitřně homogenní, dobře popsatelné a odlišitelné od jiných skupin
- pokud analyzují vzorky – daná skupina obsahuje vzorky s podobným druhovým složením (např. podobná stanoviště)
- pokud analyzují druhy – daná skupina obsahuje druhy s podobným ekologickým chováním

KLASIFIKACE



KLASIFIKACE



KLASIFIKACE

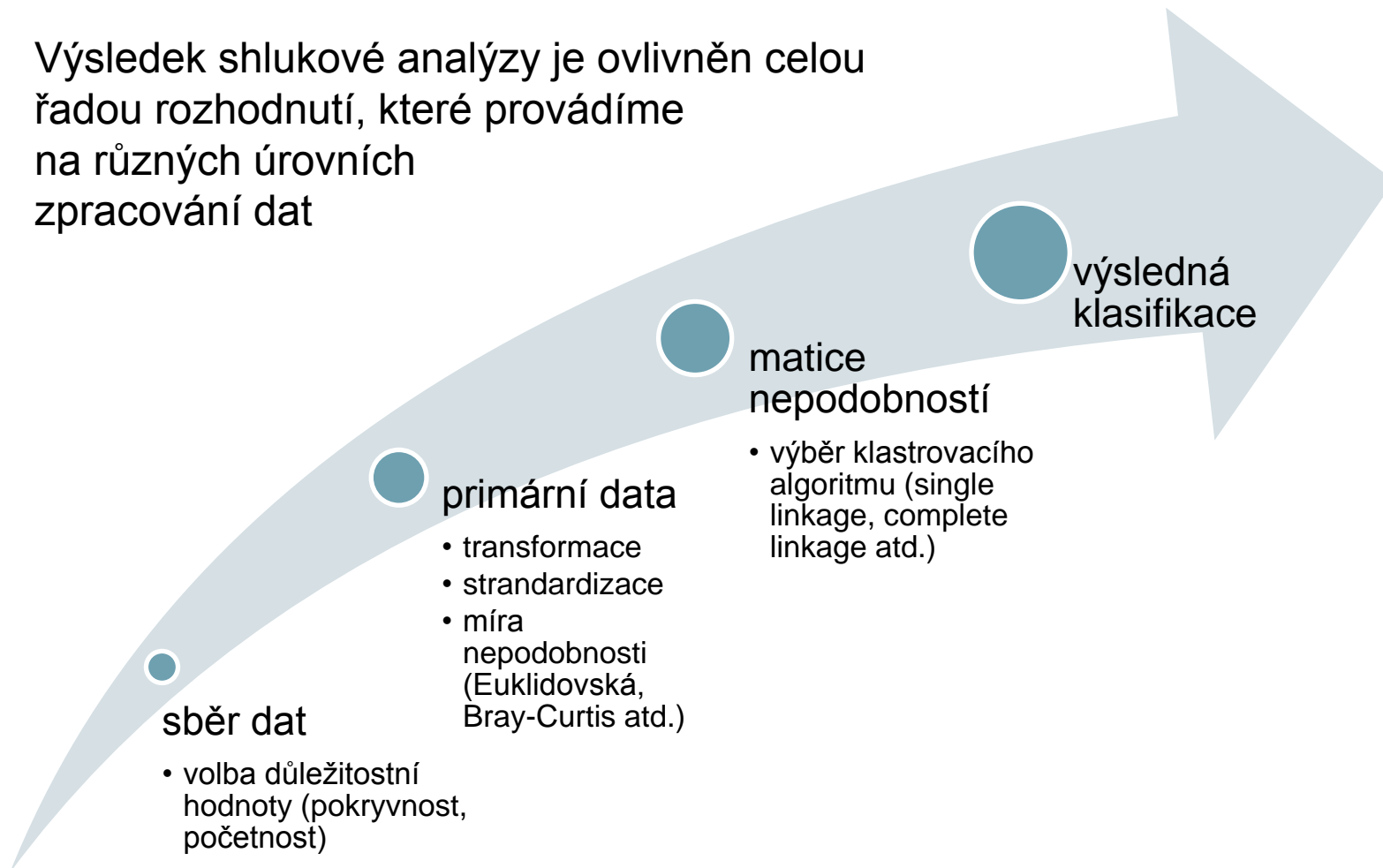
HIERARCHICKÁ A AGLOMERATIVNÍ

Cluster analysis – shluková analýza

- hierarchická metoda
 - shluky jsou hierarchicky uspořádány
- aglomerativní metoda
 - shluky jsou tvořeny ‘odspodu’, tzn. postupným shlukováním jednotlivých vzorků do větších skupin
- základní volby:
 - míra nepodobnosti mezi vzorky (*distance measure*)
 - shlukovací (klastrovací) algoritmus (*clustering algorithm*)
- pozor – nejde o OBJEKTIVNÍ metodu klasifikace (ta neexistuje), protože výsledná podoba klasifikace je ovlivněna řadou našich SUBJEKTIVNÍCH rozhodnutí

SHLUKOVÁ ANALÝZA (*CLUSTER ANALYSIS*)

Výsledek shlukové analýzy je ovlivněn celou řadou rozhodnutí, které provádíme na různých úrovních zpracování dat



sběr dat

- volba důležitosti hodnoty (pokryvnost, početnost)

primární data

- transformace
- standardizace
- míra nepodobnosti (Euklidovská, Bray-Curtis atd.)

matice nepodobností

- výběr klastrovacího algoritmu (single linkage, complete linkage atd.)

výsledná klasifikace

SHLUKOVÁ ANALÝZA (*CLUSTER ANALYSIS*)

SHLUKOVACÍ ALGORITMY

Single linkage, nearest neighbour

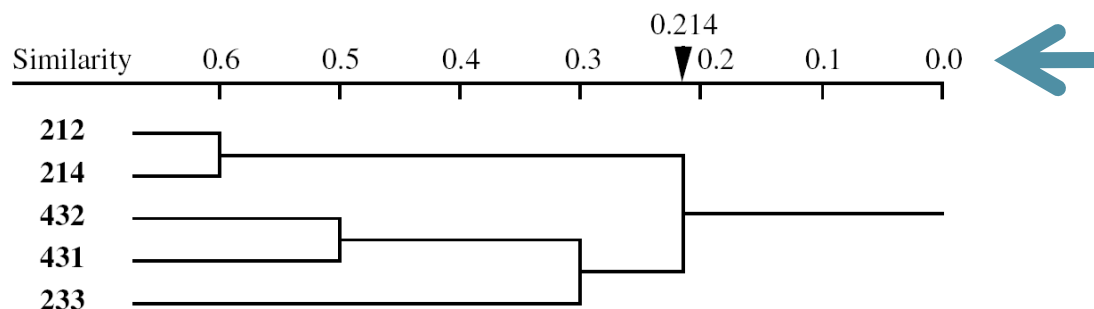
Ponds	Ponds				
	212	214	233	431	432
212	—				
214	0.600	—			
233	0.000	0.071	—		
431	0.000	0.063	0.300	—	
432	0.000	0.214	0.200	0.500	—

matice podobností

páry vzorků seřazené podle podobnosti

S_{20}	Pairs formed
0.600	212-214
0.500	431-432
0.300	233-431
0.214	214-432
0.200	233-432
0.071	214-233
0.063	214-431
0.000	212-233
0.000	212-431
0.000	212-432

Legendre & Legendre 1998



výsledný dendrogram

SHLUKOVÁ ANALÝZA (*CLUSTER ANALYSIS*)

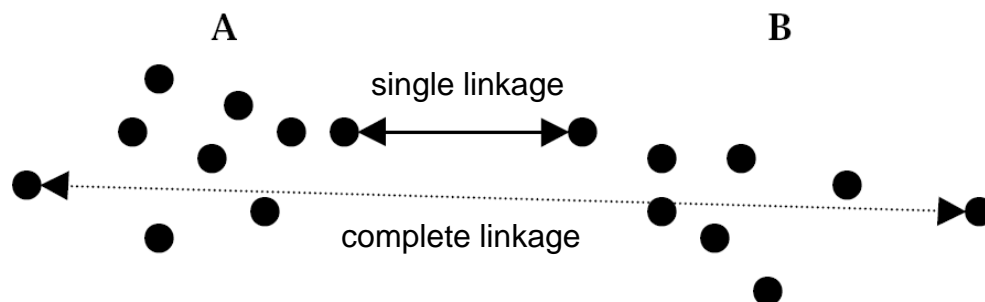
SHLUKOVACÍ ALGORITMY

Single linkage (nearest neighbour, metoda jednospojná)

- vzorky se pojí ke shluku, ve kterém je jim nejpodobnější vzorek
- *přidám se ke skupině, ve které je ten, kdo je mi nejvíc sympatický*

Complete linkage (farthest neighbour, metoda všespojná)

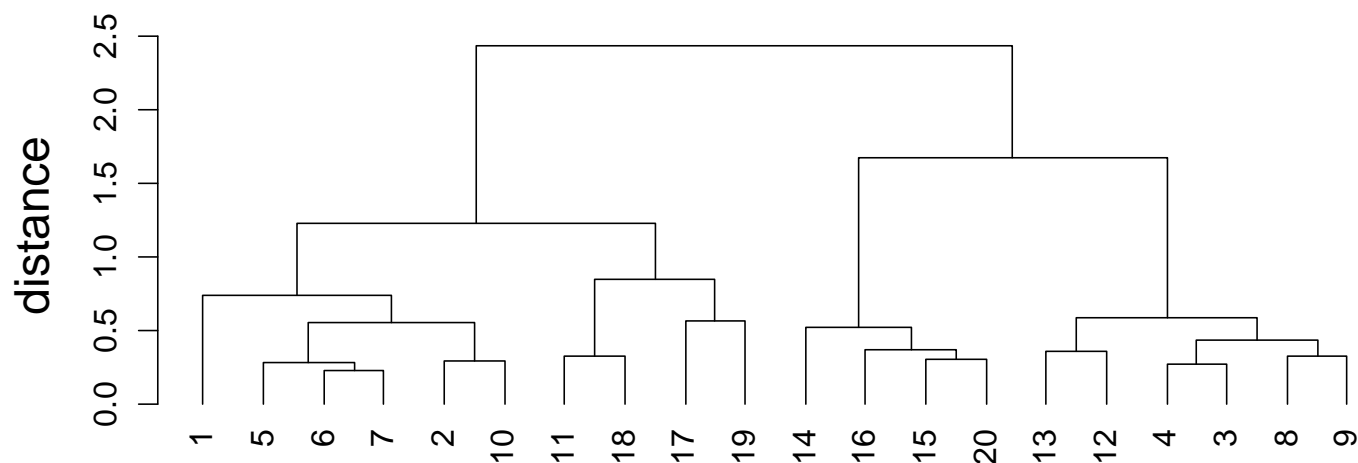
- vzorky se připojí ke shluku až v okamžiku, kdy shluk obsahuje všechny podobné vzorky
- *přidám se ke skupině ve které je ten, kdo je mi nejmíň nesympatický*



SHLUKOVÁ ANALÝZA (*CLUSTER ANALYSIS*)

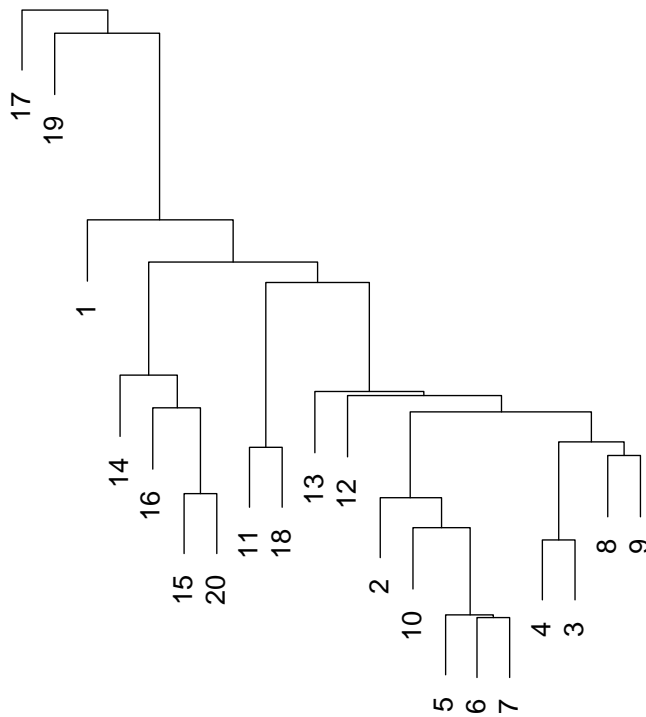
DENDROGRAM

- záleží na tom, které vzorky jsou spojeny na které úrovni
- nezáleží na tom, který vzorek (skupina) je vpravo a který vlevo

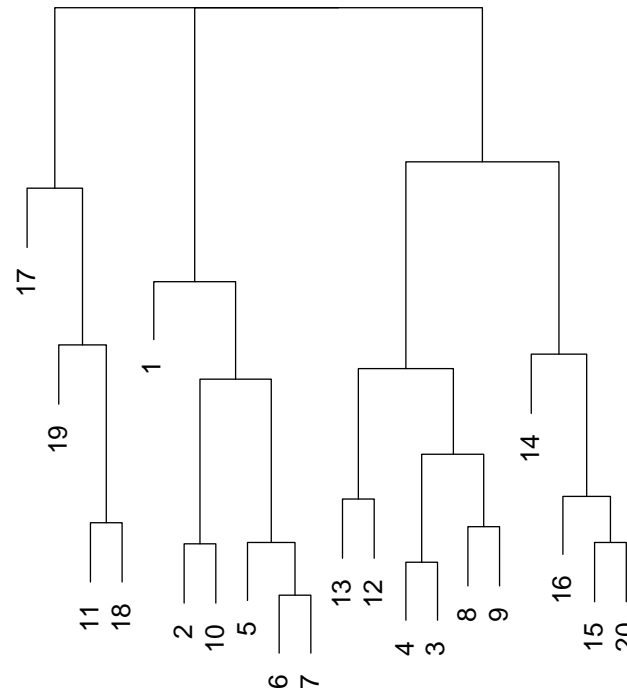


SINGLE LINKAGE VS COMPLETE LINKAGE

Bray-Curtis distance / Single linkage



Bray-Curtis distance / Complete linkage

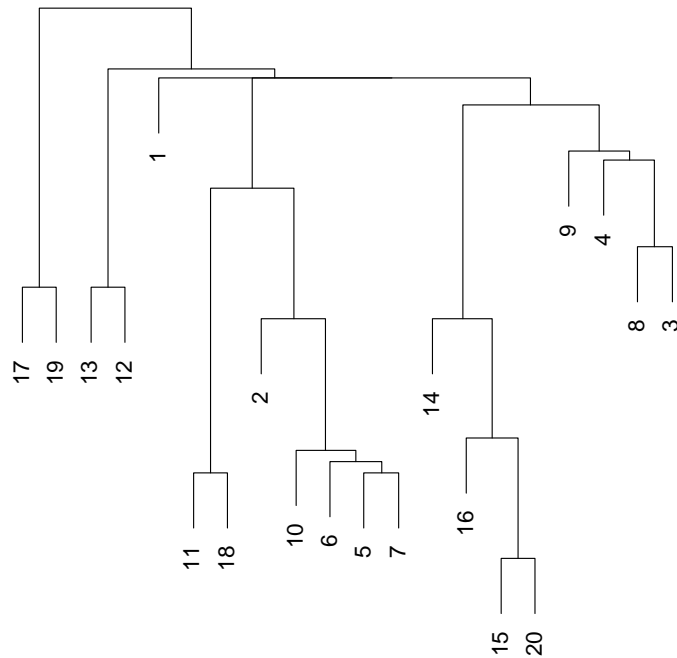


single linkage se výrazně řetězí

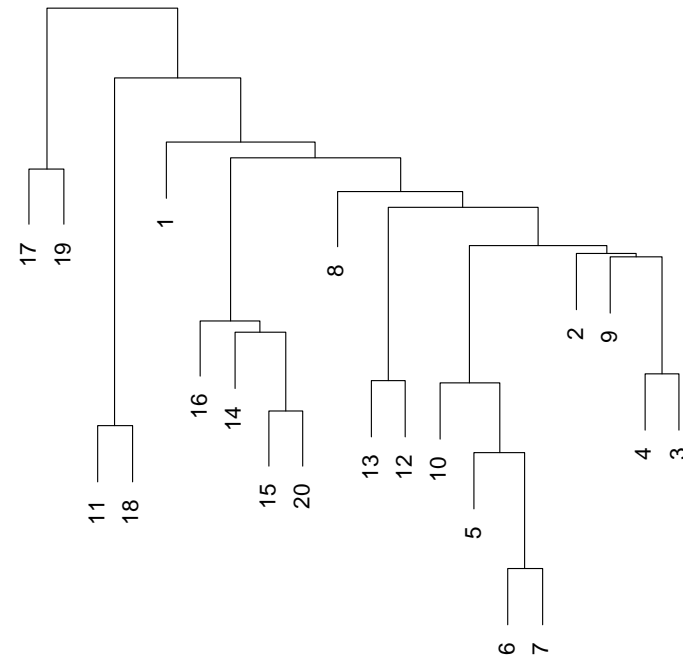
SINGLE LINKAGE

VLIV TRANSFORMACE DRUHOVÝCH DAT

Single linkage / Euclidean distance / no transformation



Single linkage / Euclidean distance / LOG transformation



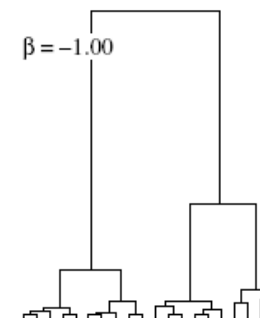
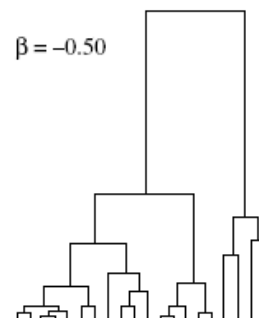
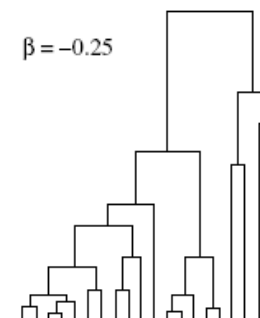
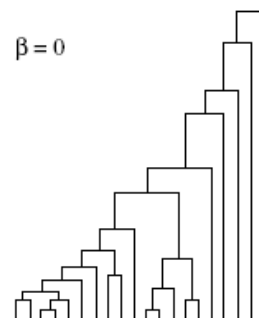
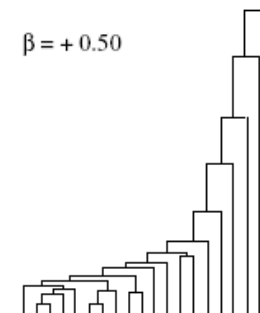
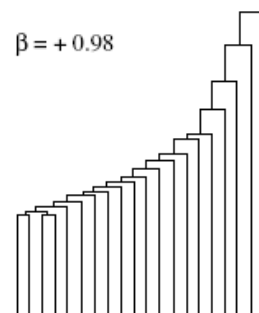
transformace dat (např. logaritmická) může výrazně ovlivnit výsledný dendrogram – v případě euklidovských vzdáleností a *single linkage* metody obzvlášť

SHLUKOVÁ ANALÝZA (*CLUSTER ANALYSIS*)

SHLUKOVACÍ ALGORITMY

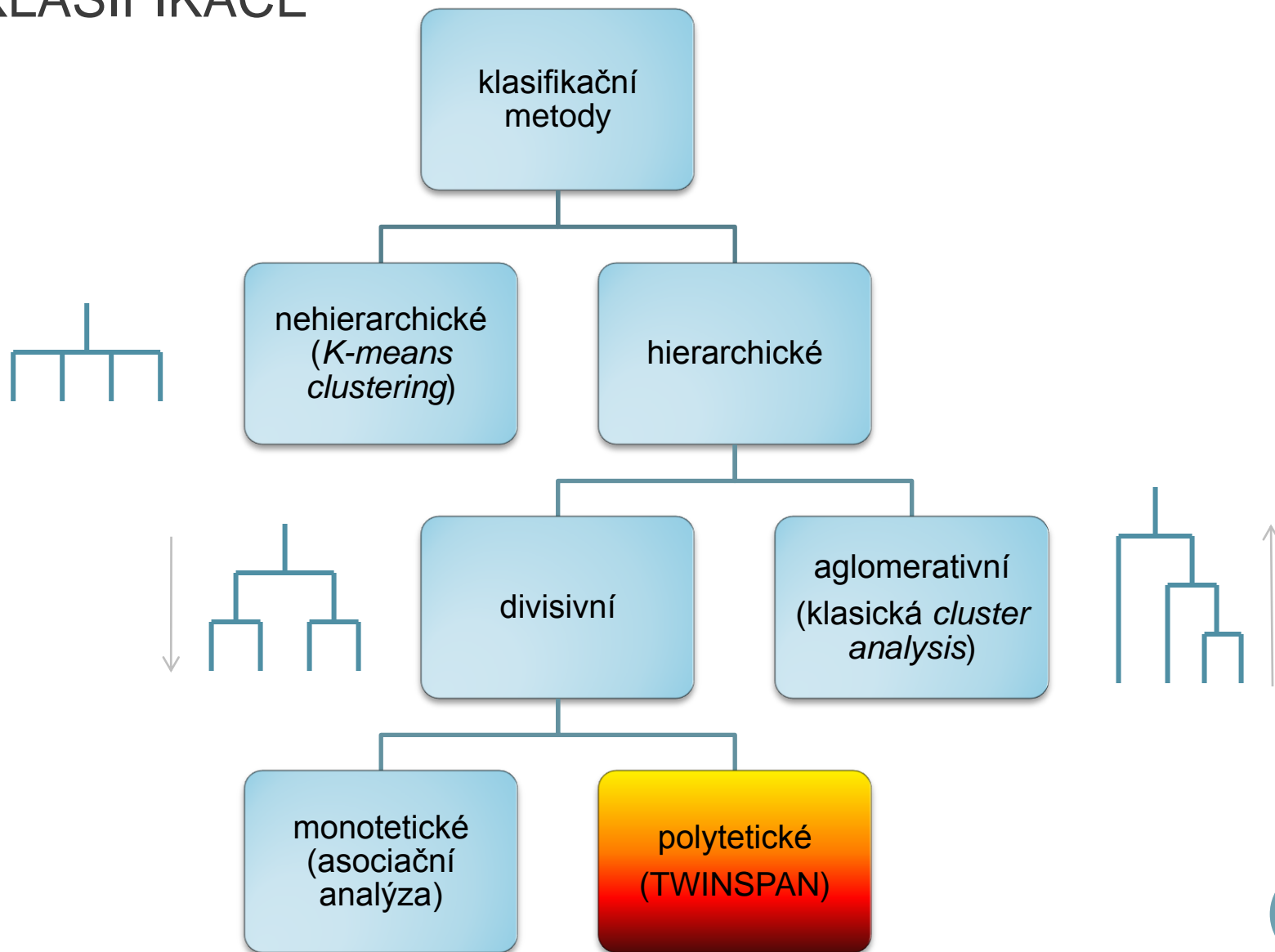
Flexible clustering (beta flexible)

- nastavení parametru β ovlivňuje řetězení dendrogramu
- nejvíc se řetězí pro $\beta \sim 1$, nejméně pro $\beta = -1$
- optimální reprezentace vzdáleností mezi vzorky je při $\beta = -0,25$



Legendre & Legendre 1998

KLASIFIKACE



KLASIFIKACE

HIERARCHICKÁ A **DIVISIVNÍ**

TWINSpan (*Two Way INdicator Species ANalysis*)

- divisivní metoda
 - začíná dělením celého souboru vzorků a postupuje směrem dolů
- polytetická metoda
 - každé dělení závisí na **několika (indikačních) druzích** (x monotetická metoda – dělení ovlivňuje jediný druh)
- metoda velmi oblíbená mezi vegetačními ekology
 - **ale** – algoritmus je poměrně složitý, s řadou arbitrárních kroků, a proto má také řadu zarytých odpůrců
- vzorky jsou uspořádány podle první osy korespondenční analýzy (CA, DCA) a podle ní jsou rozděleny do dvou shluků (vzorky s pozitivním skóre a negativním skóre)
- metoda ošetří vzorky, které leží blízko středu osy, a které tak mají velkou pravděpodobnost, že budou špatně klasifikovány
- stejný postup je následně aplikován i na jednotlivé shluky

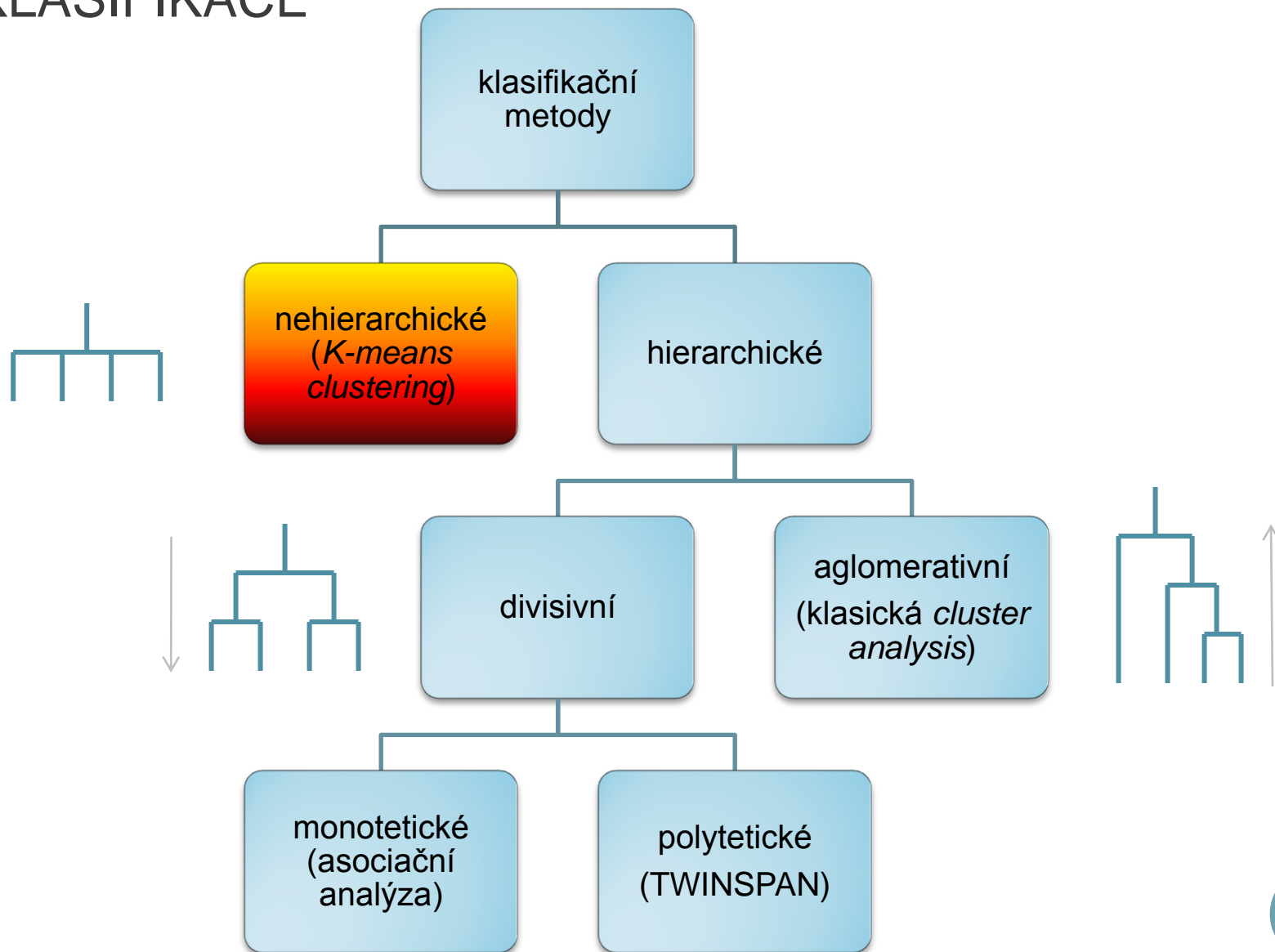
KLASIFIKACE

HIERARCHICKÁ A **DIVISIVNÍ**

TWINSPAN (*Two Way INdicator Species ANalysis*)

- *pseudospecies*
 - metoda primárně funguje pro kvalitativní data
 - kvantitativní informace se dodává rozdělením druhů na *pseudospecies* podle abundance (*cut levels*)
- výsledkem je (mimo jiné) tabulka podobná fytocenologické
 - snímky z určitých klastrů a druhy s vysokou fidelitou k dané skupině jsou seskupeny dohromady
- metoda vhodná v případě, že jsou data strukturovaná podle jednoho výrazného gradientu
- vhodné na hledání (několika málo) ekologicky interpretovatelných skupin v datech
- PC-ORD, JUICE

KLASIFIKACE



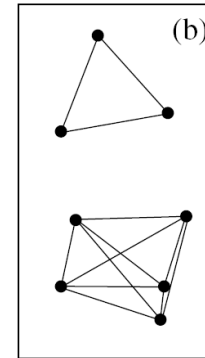
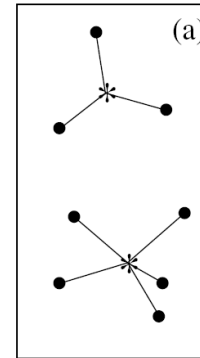
KLASIFIKACE

NEHIERARCHICKÁ

K-means clustering

(shlukování metodou K-průměrů)

- nehierarchická metoda – všechny shluky jsou si rovny
- minimalizuje sumy čtverců vzdáleností mezi vzorky uvnitř shluků
- na začátku uživatel zvolí počet shluků
- iterativní metoda, začne od náhodného přiřazení vzorků do shluků, postupně přehazuje vzorky mezi shluky a hledá optimální řešení
- výsledek do určité míry záleží na počátečním rozmístění shluků do vzorků a je proto dobré proces mnohokrát zopakovat (najít stabilní řešení)
- STATISTICA, SYN-TAX 2000



Legendre & Legendre 1998

KLASIFIKACE

OBECNÉ ROZDĚLENÍ

Subjektivní

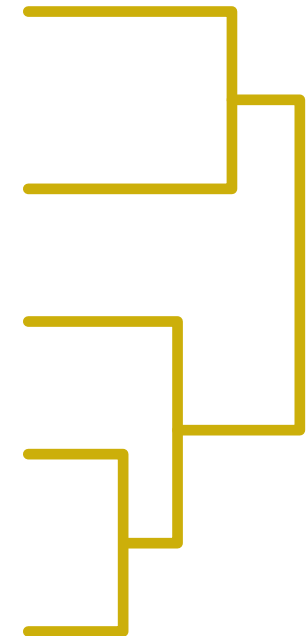
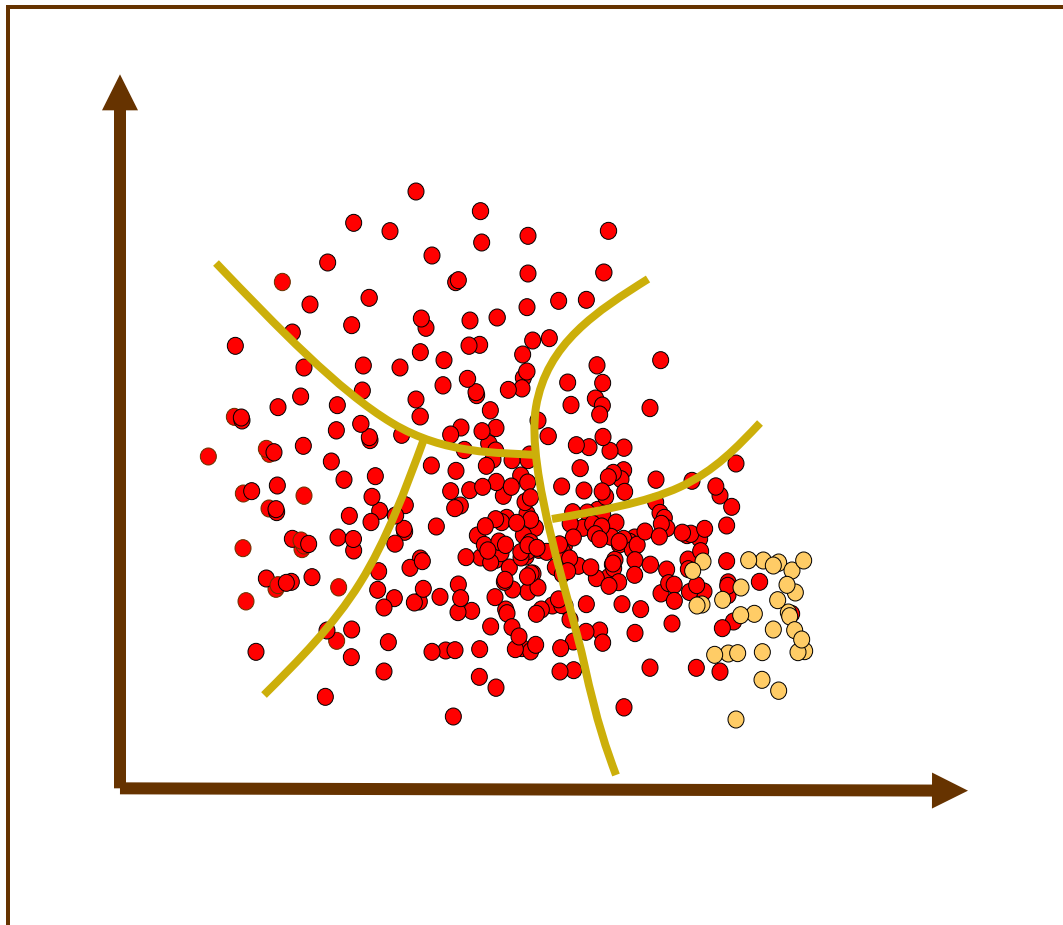
- založená na subjektivních kritériích badatele, špatně (pokud vůbec) reprodukovatelná

Formalizovaná (ne objektivní!)

- omezený výběr jasně specifikovaných (formálních) kritérií, dobře reprodukovatelná
- neřízená (*unsupervised*)
 - numerické metody klasifikace (*cluster analysis*, **TWINSpan**)
 - výslednou klasifikaci můžeme ovlivnit pouze výběrem metody (kombinace klasifikačního algoritmu a míry podobnosti), případně požadovaného počtu shluků
- řízená (*supervised*)
 - ANN – *artificial neural networks*, klasifikační stromy, náhodné lesy (*random forests*), **COCKTAIL**
 - klasifikační systém musíme nejdříve naučit, jak má vypadat výsledná klasifikace (*training*), a systém ji pak reprodukuje na dalších vzorcích

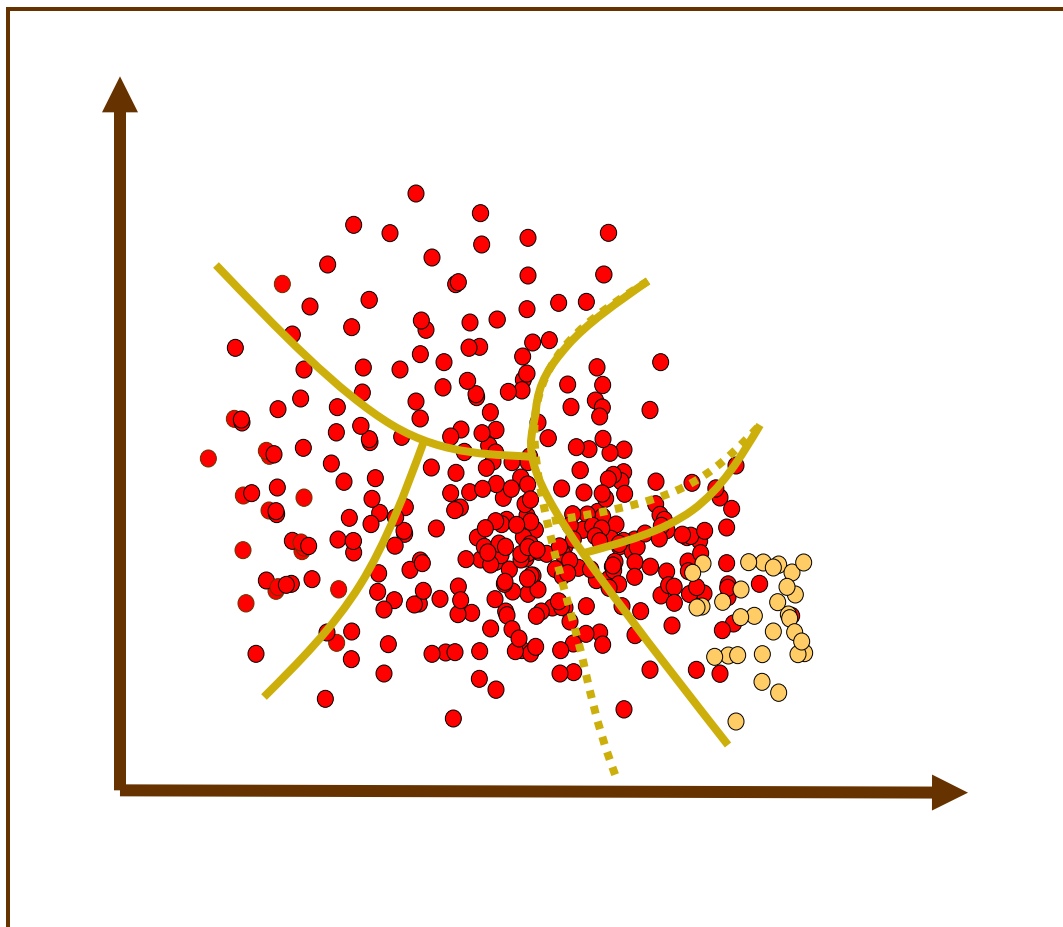
NEVÝHODY NEŘÍZENÉ KLASIFIKACE

(PŘÍŘAZENÍ VZORKŮ DO SKUPIN SE MĚNÍ PODLE KONTEXTU)

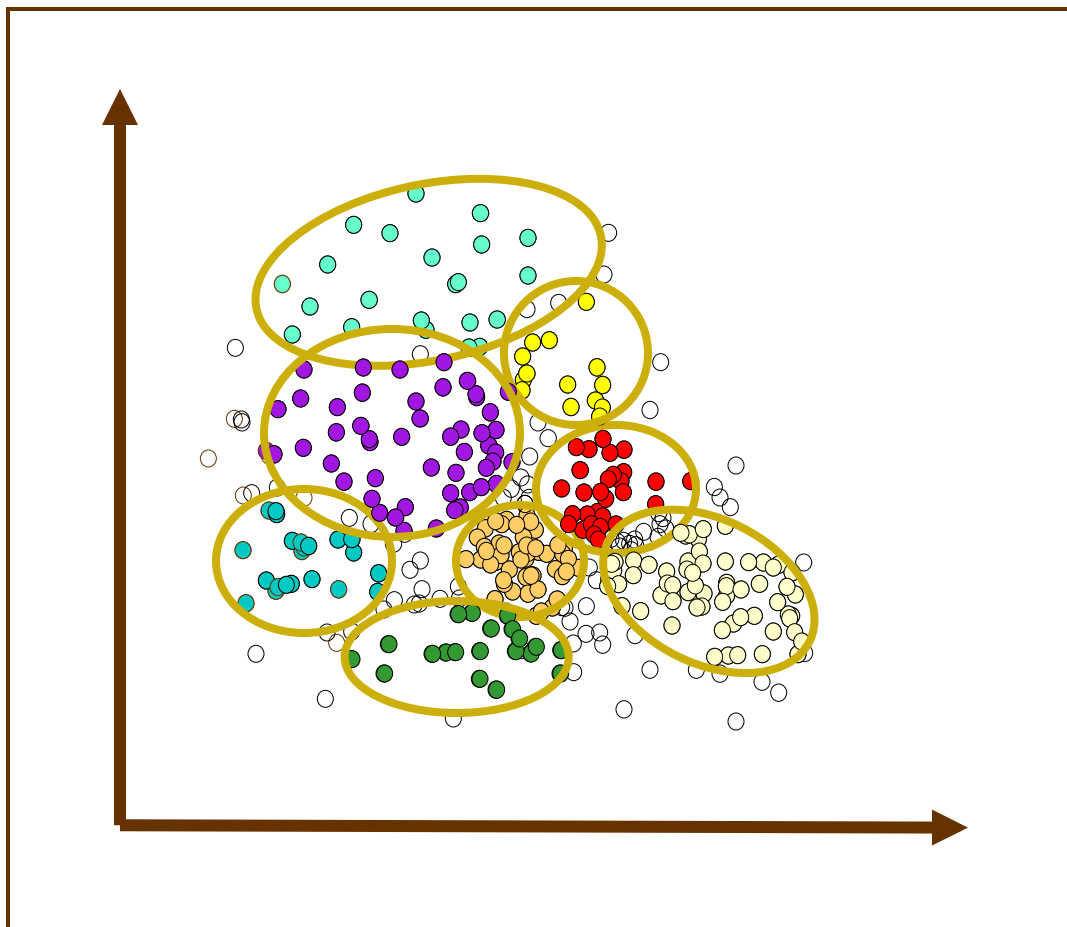


NEVÝHODY NEŘÍZENÉ KLASIFIKACE

(PŘÍŘAZENÍ VZORKŮ DO SKUPIN SE MĚNÍ PODLE KONTEXTU)



ŘÍZENÁ KLASIFIKACE



ŘÍZENÁ KLASIFIKACE

