

Vícerozměrná analýza dat



Jiří Jarkovský

Plán kurzu

- ☑ **Každých 14 dní 4 vyučovací hodiny**

- ☑ **Ukončení zkouškou**
 - **Písemná**
 - **Zaměřená na principy a aplikace analýz**

- ☑ **Cíl kurzu**
 - **Vysvětlit principy vícerozměrných analýz, jejich aplikaci v biologii a jejich interpretaci**
 - **Přehled základního software**
 - **Příklady na reálných datech**

Náplň kurzu I

- ☑ **Vícerozměná analýza dat – smysl a cíle**
 - Příklady užití vícerozměrných analýz
 - Výhody a nevýhody vícerozměrné analýzy dat
 - Parametrická a neparametrická vícerozměrná statistika
 - Statistické SW pro vícerozměrnou analýzu dat

- ☑ **Podobnost a vzdálenost objektů ve vícerozměrném prostoru**
 - **Metriky podobnosti a vzdálenosti a jejich úskalí**
 - ☐ Obecné metriky podobnosti a vzdálenosti
 - ☐ Metriky podobnosti pro biologická společenstva – problém double zero
 - **Asociační matice**
 - ☐ Struktura asociační matice
 - ☐ Práce s asociační maticí
 - ☐ Mantelův test

- ☑ **Vícerozměrné statistické testy a rozložení**
 - Vícerozměrné normální rozložení
 - Vícerozměrné charakteristiky - medoid
 - Hottelingovo T, Wishartovo rozdělení

- ☑ **Základy maticové algebry**
 - Typy matic a jejich využití při vícerozměrné analýze dat
 - Matematické operace s maticemi
 - Eigenvalues (vlastní čísla) a eigenvectory (vlastní vektory) matic

Náplň kurzu II

- ☑ **Shluková analýza**
 - **Kriteria posuzování výsledků shlukovacích metod**
 - ☐ **Minimální vnitroshluková varibilita**
 - ☐ **Maximální mezishluková variabilita**
 - ☐ **Silhouette width**
 - **Hierarchické aglomerativní shlukování**
 - ☐ **Shlukovací algoritmy**
 - **nearest neighbour (single linkage)**
 - **farthest neighbour (complete linkage)**
 - **UPGMA**
 - **WPGMA**
 - **UPGMC**
 - **WPGMC**
 - **Ward's method**
 - **Hierarchické divizivní shlukování**
 - ☐ **TWINSpan**
 - **Nehierarchické divizivní shlukování**
 - ☐ **K-means clustering**
 - ☐ **X-means clustering**
 - ☐ **Partitioning around medoids (PAM)**

Náplň kurzu III

- ☑ **Ordinační analýzy**
 - **Principy ordinačních analýz - redukce dimenzionality**
 - ☐ **Eigenvektor**
 - ☐ **Eigenvalue**
 - **Základní typy ordinační analýzy a jejich užití**
 - ☐ **PCA**
 - ☐ **CA**
 - ☐ **DCA**
 - ☐ **CCA**
 - ☐ **DCCA**
 - ☐ **RDA**
 - ☐ **MDS**
 - ☐ **PCoA**
 - ☐ **Kanonická korelace**
- ☑ **Analýza hlavních komponent**
 - **PCA na základě euklidovské vzdálenosti**
 - **PCA na základě korelací a kovariancí**
 - **Normalised PCA**
 - **Biplot a jeho interpretace**
- ☑ **Korespondenční analýza a její varianty**
 - **CA, DCA, CCA, DCCA**
- ☑ **MDS a PCoA – ordinační analýza na libovolné asociační matici**

Software pro vícerozměrnou analýzu

☑ „Klikací všeobecné SW“

- Statistica
- SPSS
- SAS

☑ Specializované SW

- PcORD
- CANOCO
- PAST
- WEKA
- ORANGE
- SW pro microarray analýzu
- Nejružnější utility na netu
-


☑ Univerzální SW

- R - ADE4 atd.

Vícerozměrná analýza dat

Základní statistické výpočty s vazbou na vícerozměrnou analýzu

Vztah klasické a vícerozměrné statistiky

- ✓ **Vícerozměrná analýza dat využívá přístupů klasické statistiky**
- ✓ **Zároveň je citlivá i na jejich problémy** 
- ✓ **Agregace dat přes sumární statistiku nebo kontingenční tabulky – korespondenční analýza**
- ✓ **Korelace – analýza hlavních komponent, faktorová analýza, diskriminační analýza**

Kontingenční tabulka

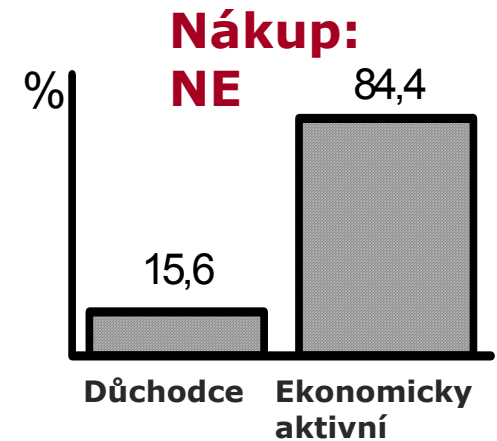
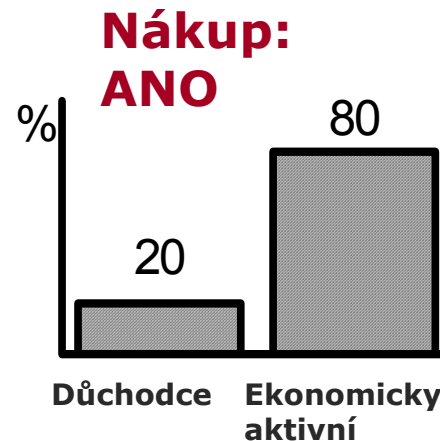
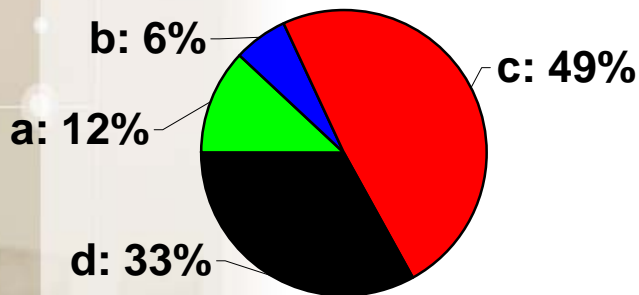
Nákup

Důchodový věk

	Ano	Ne	Σ
Ano	20	82	102
Ne	10	54	64
Σ	30	136	166

- ☑ Kontingenční tabulka je používána pro hodnocení vztahu kategoriálních proměnných

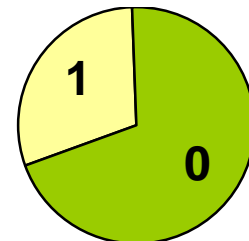
Kontingenční tabulka v obrázku



Kontingenční tabulky – princip analýzy

Binomické jevy (1/0)

$$\chi^2_{(1)} = \frac{\left[\begin{array}{c} \text{pozorovaná} \\ \text{četnost} \end{array} - \begin{array}{c} \text{očekávaná} \\ \text{četnost} \end{array} \right]^2}{\text{očekávaná četnost}} + \frac{\left[\begin{array}{c} \text{pozorovaná} \\ \text{četnost} \end{array} - \begin{array}{c} \text{očekávaná} \\ \text{četnost} \end{array} \right]^2}{\text{očekávaná četnost}}$$



I. jev 1

II. jev 2

Příklad



10 000 lidí hází mincí



rub: 4 000 případů (R)
líc: 6 000 případů (L)



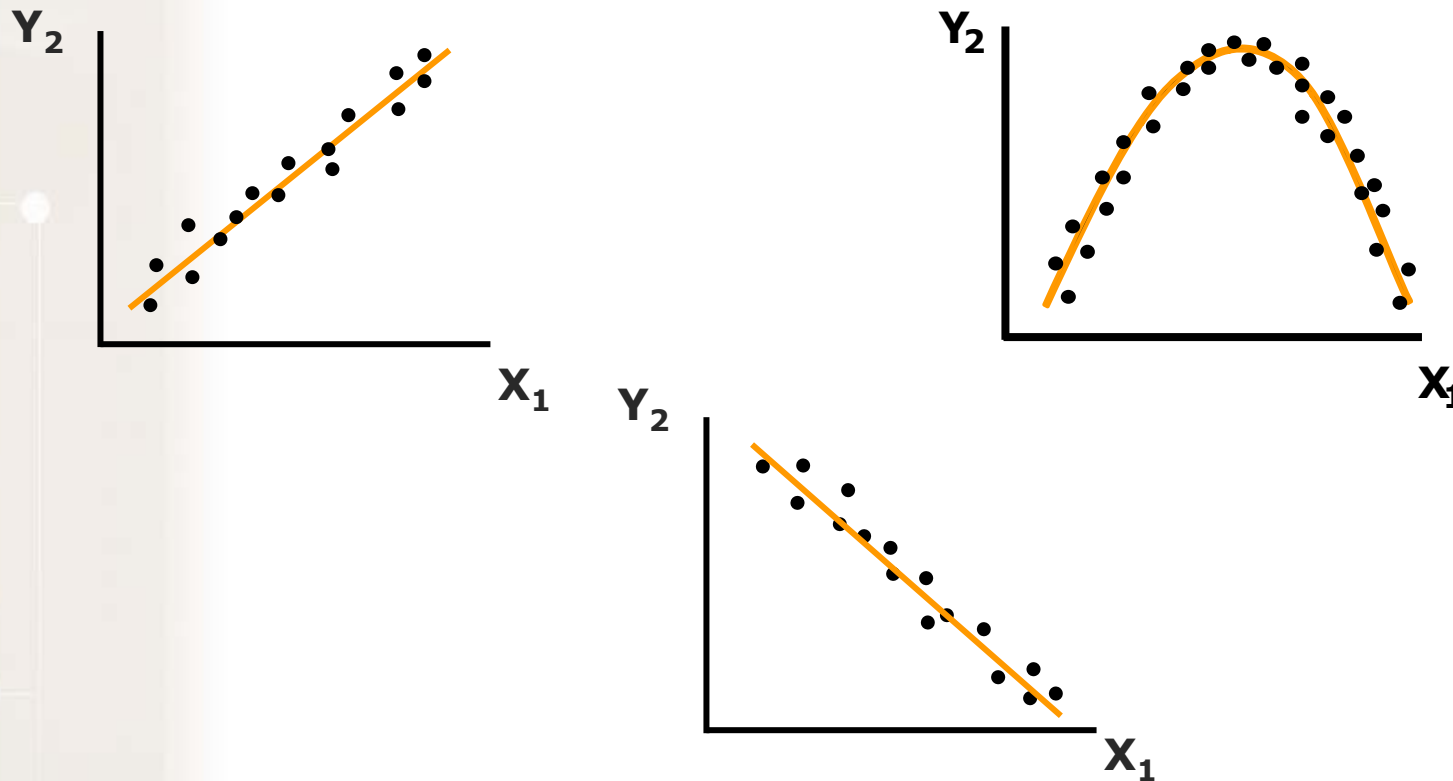
Lze výsledek považovat za statisticky významně odlišný (nebo neodlišný) od očekávaného poměru R : L = 1 : 1 ?



Stejným způsobem, tedy hodnocením odchylek od očekávaného vyrovnaného počtu případů hodnotí data i korespondenční analýza

Korelační analýza

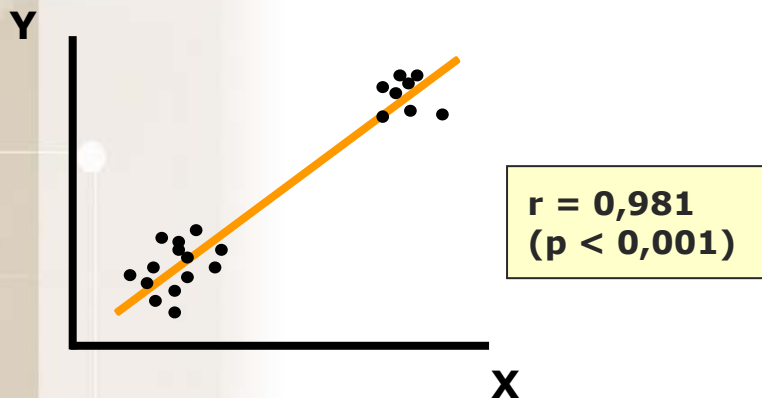
- Korelace - vztah (závislost) dvou znaků (parametrů)



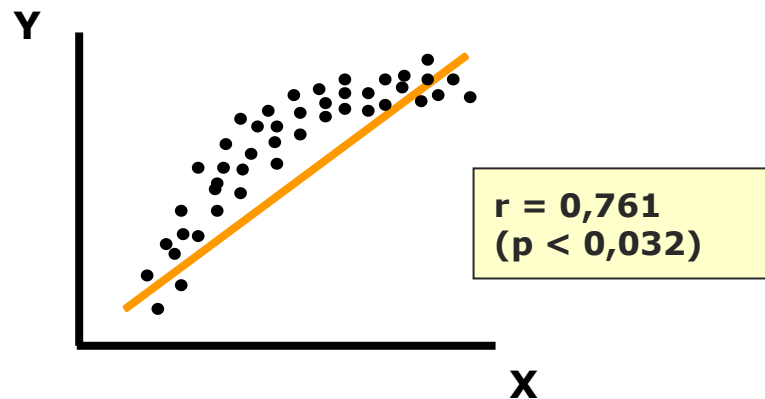
Korelace mezi parametry jsou základem faktorové analýzy a analýzy hlavních komponent, pokud vazby mezi parametry nejsou tyto metody postrádají smysl.

Rizika korelační analýzy

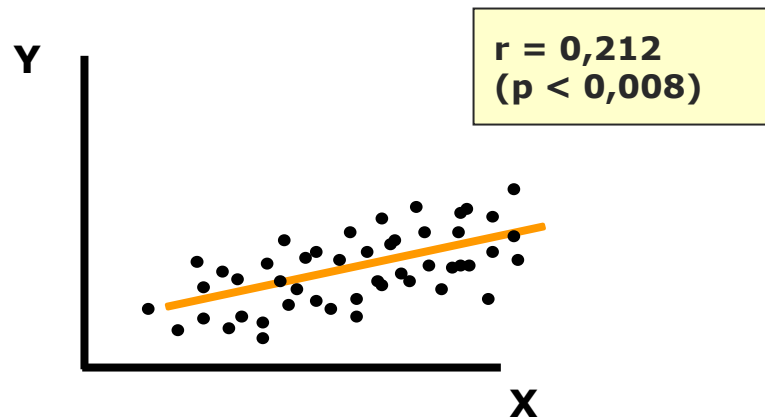
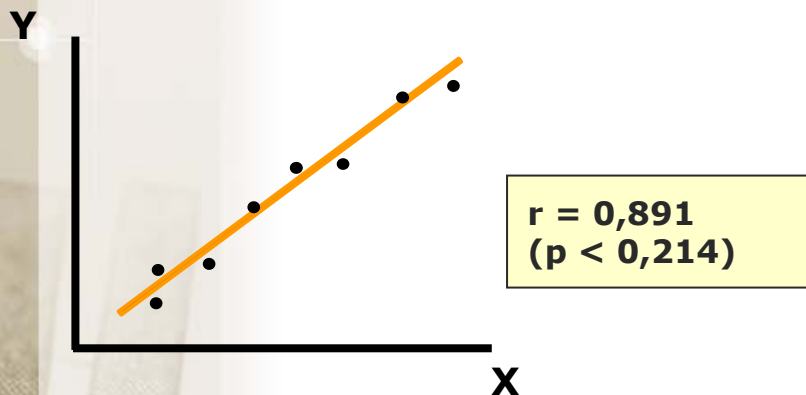
Problém rozložení hodnot



Problém typu modelu



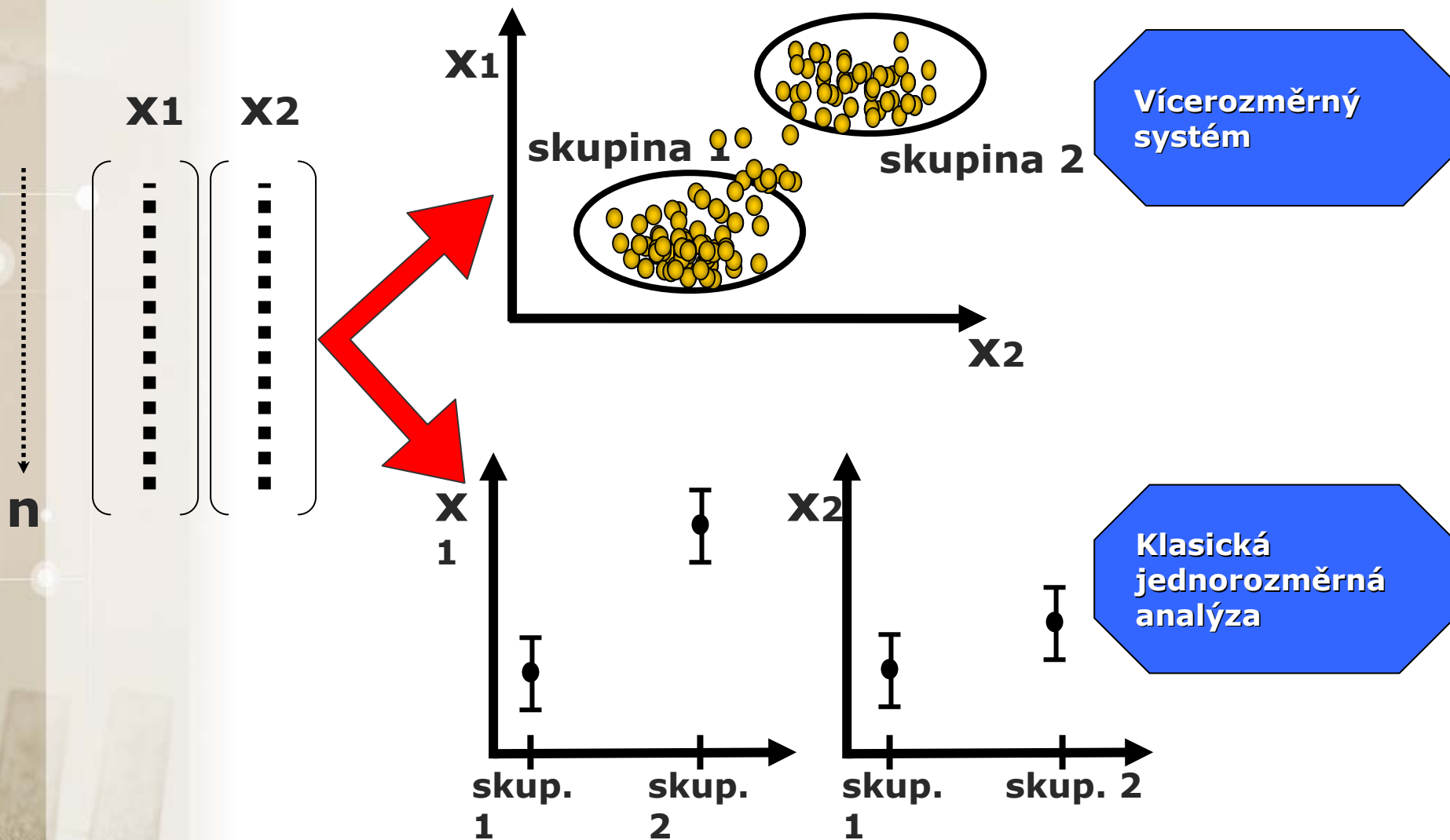
Problém velikosti vzorku



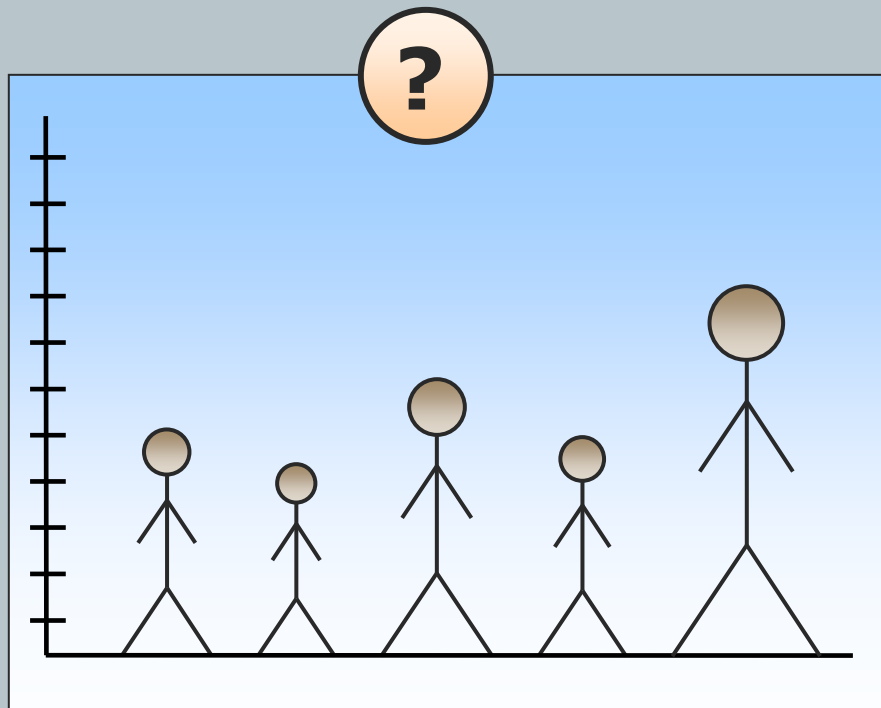
Vícerozměrná analýza dat

Význam vícerozměrného hodnocení dat

Vícerozměrné vnímání skutečnosti – nová kvalita analýzy dat



Běžná sumarizace dat „likviduje“ individualitu jedince



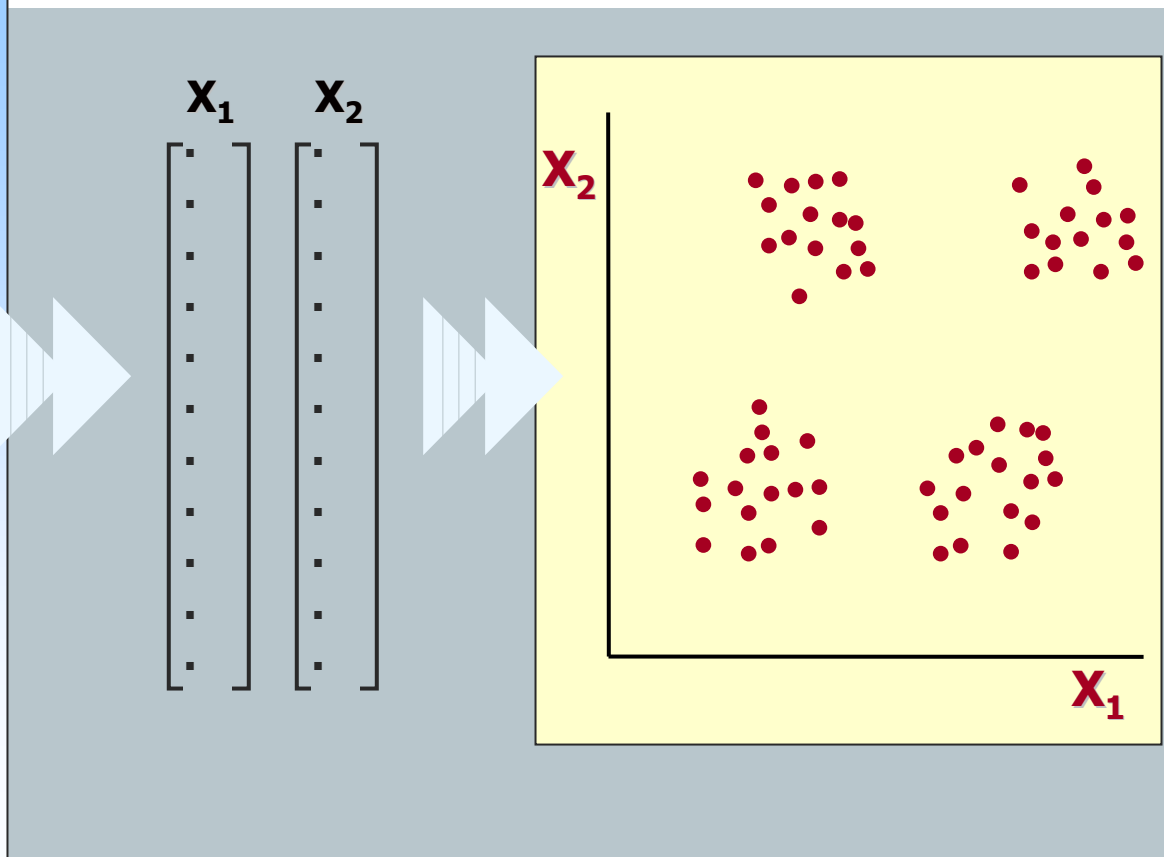
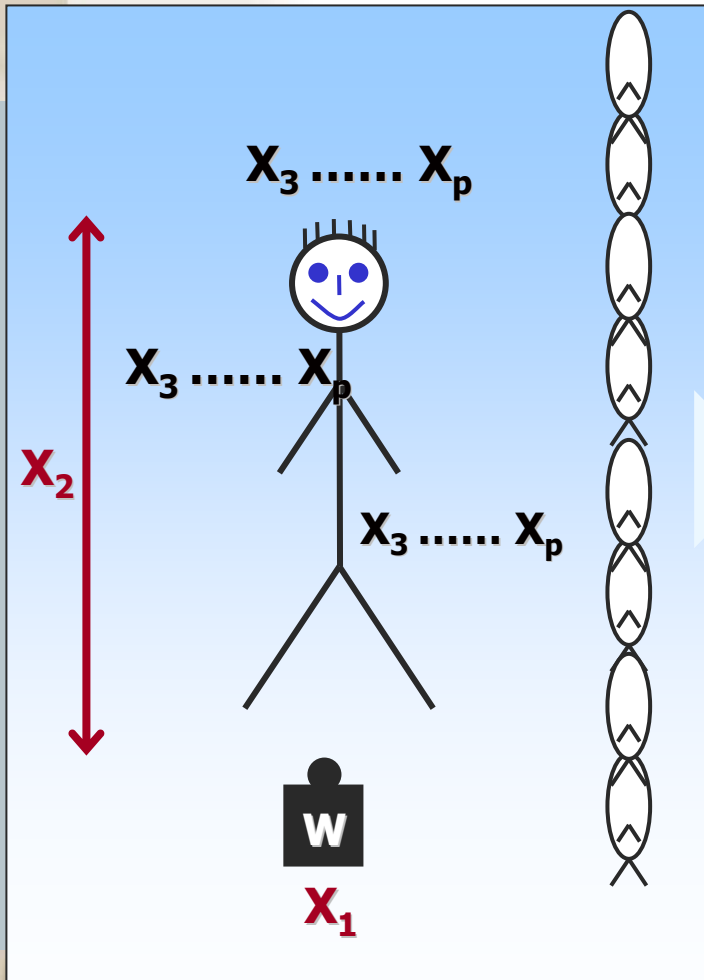
Průměr \pm SE

BĚŽNÁ STATISTICKÁ
SUMARIZACE

- ✓ *Zpřehlednění dat*
- ✓ *Neodliší původní měření*

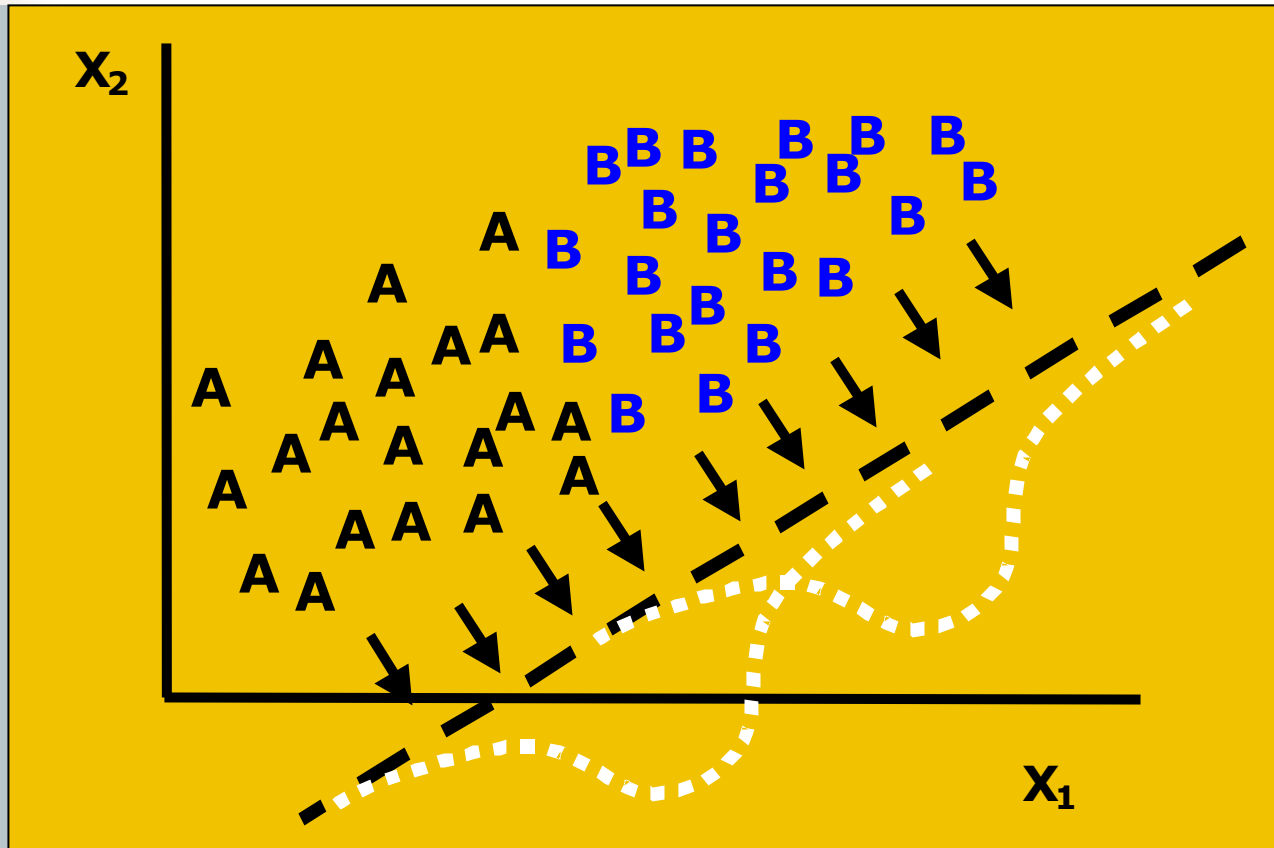
Vícerozměrné hodnocení

... s ohledem na individualitu !



Vícerozměrné hodnocení – nová kvalita

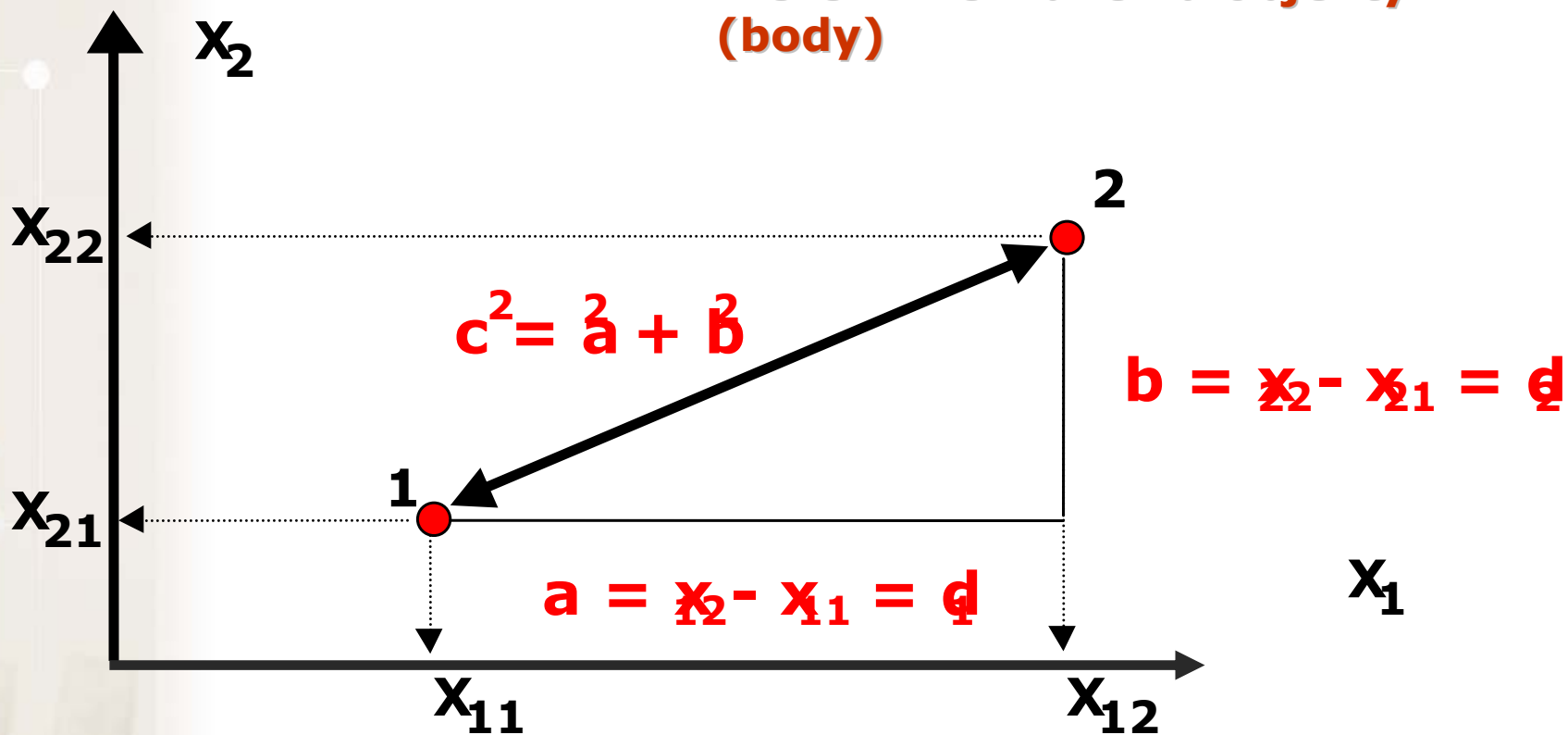
Pouze kombinované parametry mají odpovídající informační sílu



příklad: $X_1 =$

Vícerozměrné hodnocení vychází z jednoduchých principů

příklad: vícerozměrná vzdálenost měření mezi dvěma objekty (body)



Vícerozměrné modelování je strategickou disciplínou

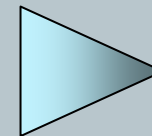


$X_1 \dots X_n$

technické parametry
automobilu

$X_{n+1} \dots X_p$

řidičovy schopnosti
a jeho stav



$X_{p+1} \dots X_2$

rychlost, povrch,
situace

X_1

⋮

X_2

⋮

X_3

⋮

X_4

⋮

X_5

⋮

⋮

X_p

⋮

Vícerozměrná analýza dat

Základní principy vícerozměrného hodnocení dat

Pojmy vícerozměrných analýz

- ✓ **Vícerozměrné metody:** Název vícerozměrné vychází z typu vstupních dat, tato data jsou tvořena jednotlivými objekty (i.e. klienti) a každý z nich je charakterizován svými parametry (věk, příjem atd.) a každý z těchto parametrů můžeme považovat za jeden rozměr objektu.
- ✓ **Maticová algebra:** Základem práce s daty a výpočtů vícerozměrných metod je maticová algebra, matice tvoří jak vstupní, tak výstupní data a probíhají na nich výpočty.
- ✓ **$N \times P$ matice:** N objektů s p parametry pak vytváří tzv. $N \times P$ matici, která je prvním typem vstupu dat do vícerozměrných analýz.
- ✓ **Asociační matice:** Na základě těchto matic jsou počítány matice asociační na nichž pak probíhají další výpočty, jde o čtvercové matice obsahující informace o podobnosti nebo rozdílnosti (tzv. metriky) buď objektů (Q mode analýza) nebo parametrů (R mode analýza). Měřítko podobnosti se liší podle použité metody a typu dat, některé metody umožňují použití uživatelských metrik.

Vstupní matice vícerozměrných analýz

NxP MATICE

	parametr 1	parametr 2	parametr 3
objekt 1			
objekt 2			
objekt 3			
objekt 4			
objekt 5			
objekt 6			

Výpočet metriky
podobnosti/
vzdáleností



ASOCIAČNÍ MATICE

	objekt 1	objekt 2	objekt 3	objekt 4	objekt 5	objekt 6
objekt 1						
objekt 2						
objekt 3						
objekt 4						
objekt 5						
objekt 6						

Hodnoty parametrů pro jednotlivé objekty

Korelace, kovariance, vzdálenost, podobnost

Základní typy vícerozměrných analýz

SHLUKOVÁ ANALÝZA

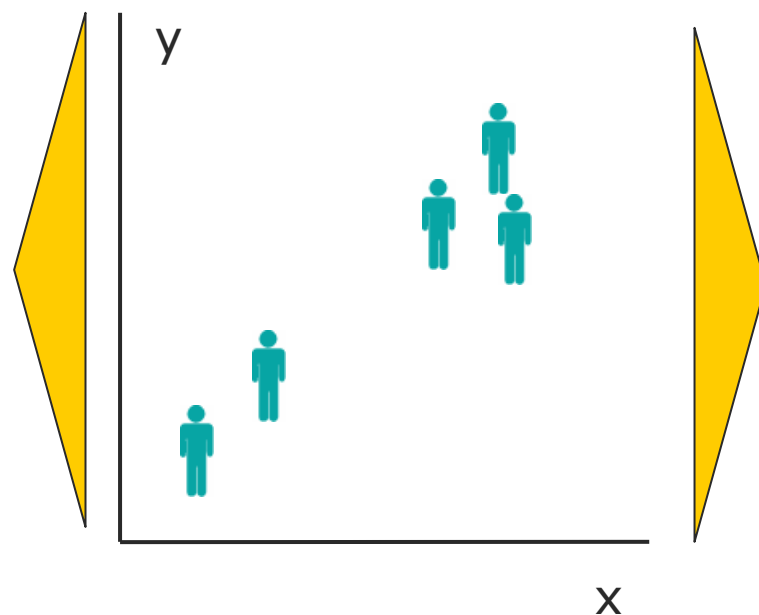
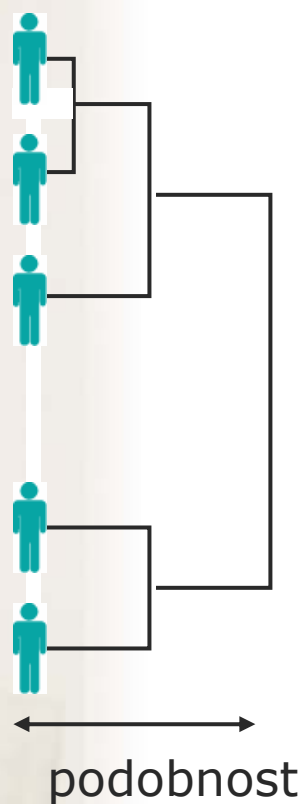
- ☑ vytváření shluků objektů na základě jejich podobnosti
- ☑ identifikace typů objektů

ORDINAČNÍ METODY

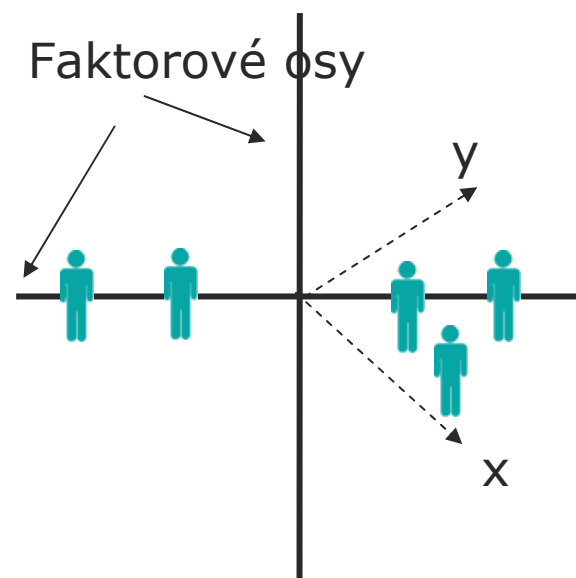
- ☑ zjednodušení vícerozměrného problému do menšího počtu rozměrů
- ☑ principem je tvorba nových rozměrů, které lépe vyčerpávají variabilitu dat

Typy vícerozměrných analýz

SHLUKOVÁ ANALÝZA



ORDINAČNÍ METODY

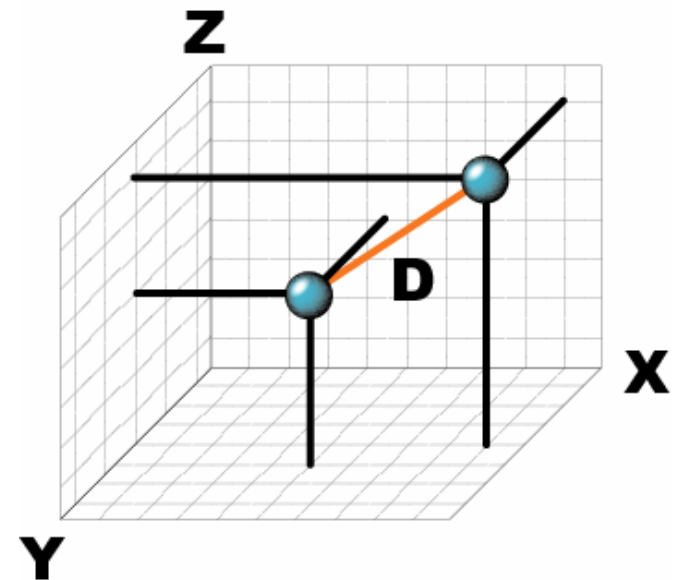


Vícerozměrná analýza dat

Asociační matice
Vícerozměrná vzdálenost a podobnost

Seznam taxonů – vícerozměrný popis společenstva

- ✓ Na seznam taxonů lze pohlížet také jako seznam rozměrů společenstva
- ✓ Záznam o nalezených taxonech tak vlastně tvoří vícerozměrný popis daného společenstva
- ✓ Společenstva můžeme srovnávat podle jejich vzájemné pozice v n-rozměrném prostoru
- ✓ Pro srovnání společenstev lze teoreticky využít libovolnou metriku vícerozměrné podobnosti nebo vzdálenosti



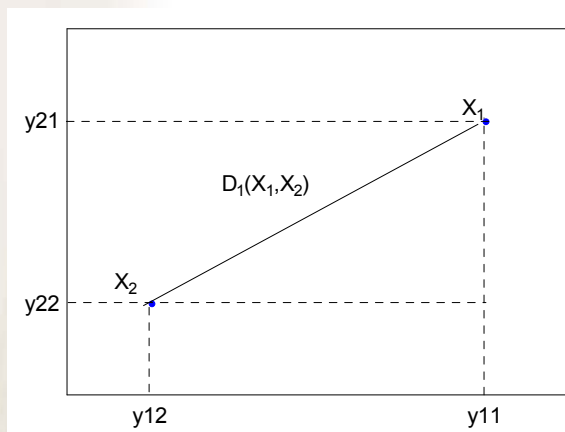
Euklidovská vzdálenost

- ☑ Jde o základní metrické měřítko vzdálenosti a počítá vzdálenost objektů obdobně jako Pythagorova věta počítá přeponu pravoúhlého trojúhelníku. Metoda je citlivá na rozdílný rozsah hodnot vstupujících proměnných (vhodným řešením může být standardizace) a double zero problém. Nemá horní hranici hodnot.

$$D_1(x_1, x_2) = \sqrt{\sum_{j=1}^p (y_{1j} - y_{2j})^2}$$

- ☑ Jako další měřítko se používá také čtverec této vzdálenosti. Jeho nevýhodou jsou semimetrické vlastnosti.

$$D_1^2(x_1, x_2) = \sum_{j=1}^p (y_{1j} - y_{2j})^2$$



Double zero problém !!!

- ✓ **V případě binárních metrik (druh se vyskytuje/nevyskytuje) není možné uvažovat stejnou váhu pro souhlas přítomnosti (11) a nepřítomnosti (00) taxonů (symetrický koeficient)**
- ✓ **Problémem využití všech typů metrik pro data abundancí spočívá v odlišném významu přítomnosti a nepřítomnosti taxonů**
- ✓ **Pokud se taxon nachází v obou srovnávaných společenstvech – znamená to že společenstva si budou v tomto ohledu podobná, protože mají podmínky umožňující přítomnost taxonu**
- ✓ **Pokud se taxon nenachází ani v jednom ze dvou srovnávaných společenstev – příčina může být nejrůznější – double zero problem**
- ✓ **Pro odstranění tohoto problému je použito asymetrické hodnocení souhlasné přítomnosti (11) a nepřítomnosti (00) taxonů (asymetrické koeficienty)**

Koeficienty podobnosti (indexy podobnosti)

- ☑ V ekologii se využívá řada indexů podobnosti založených buď na přítomnosti/nepřítomnosti taxonů nebo na abundancích

Binární koeficienty podobnosti

	Společenstvo 1	
	1	0
Společenstvo 2	1	0
	a	b
	c	d

a, b, c, d = počet případů, kdy souhlasí binární charakteristika společenstev 1 a 2
 $a+b+c+d=p$

Symetrické binární koeficienty - není rozdíl mezi případem 1-1 a 0-0
Asymetrické binární koeficienty - rozdíl mezi případem 1-1 a 0-0

Více informací a další měření vzdáleností a podobností najdete v knize
**LEGENDRE, P. & LEGENDRE, L. (1998). *Numerical ecology*.
Elsevier Science BV, Amsterdam.**

Vícerozměrná analýza dat

Symetrické binární koeficienty

Simple matching coefficient (Sokal & Michener, 1958)

- ✓ Obvyklou metodou pro výpočet podobnosti mezi dvěma objekty je podíl počtu deskriptorů, které kódují objekt stejně, a celkového počtu deskriptorů. Při použití tohoto koeficientu předpokládáme, že není rozdíl mezi nastáním 0 a 1 u deskriptorů.

$$S_1(x_1, x_2) = \frac{a + d}{p}$$

Rogers & Tanimoto koeficient (1960)

- ☑ **Dává větší váhu rozdílům než podobnostem.**

$$S_2(x_1, x_2) = \frac{a + d}{a + 2b + 2c + d}$$

Sokal & Sneath (1963)

- ☑ Další čtyři navržené koeficienty obsahují double-zero, ale jsou navrženy tak, aby se snížil vliv double-zero:

$$S_3(x_1, x_2) = \frac{2a + 2d}{2a + b + c + 2d}$$

- ☑ tento koeficient dává dvakrát větší váhu shodným deskriptorům než rozdílným;

$$S_4(x_1, x_2) = \frac{a + d}{b + c}$$

- ☑ porovnává shody a rozdíly prostým podílem v měřítku jdoucím od 0 do nekonečna;

$$S_5(x_1, x_2) = \frac{1}{4} \left[\frac{a}{a+b} + \frac{a}{a+c} + \frac{d}{b+d} + \frac{d}{c+d} \right]$$

- ☑ porovnává shodné deskriptory se součty okrajů tabulky;

$$S_6(x_1, x_2) = \frac{a}{\sqrt{(a+b)(a+c)}} \frac{d}{\sqrt{(b+d)(c+d)}}$$

- ☑ je vytvořen z geometrických průměrů členů vztahujících se k a a d , podle koeficientu S_5 .

Hammannův koeficient

$$S = \frac{a + d - b - c}{p}$$

Yuleho koeficient

$$S = \frac{ad - bc}{ad + bc}$$

Pearsonovo Φ (phi)

$$\phi = \frac{ad - bc}{\sqrt{(a + b)(c + d)(a + c)(b + d)}}$$

Vícerozměrná analýza dat

Asymetrické binární koeficienty

Jaccardův koeficient (1900, 1901, 1908)

- ☑ **Všechny členy mají stejnou váhu**

$$S_7(x_1, x_2) = \frac{a}{a + b + c}$$

Sørensenův koeficient (1948) (Coincidence index, Dice(1945))

- ☑ varianta předchozího koeficientu dává dvojnásobnou váhu dvojitým prezencím , protože se může zdát, že přítomnost druhů je více informativní než jejich absence, která může být způsobena různými faktory a nemusí nutně odrážet rozdílnost prostředí. Prezence druhu na obou lokalitách je silným ukazatelem jejich podobnosti. S_7 je monotónní k S_8 , proto podobnost pro dvě dvojice objektů vypočítaná podle S_7 bude podobná stejnému výpočtu S_8 . Oba koeficienty se liší pouze v měřítku. Tento index byl poprvé použit Dicem v R-mode studii asociací druhů. Jiná varianta tohoto koeficientu dává duplicitním prezencím trojnásobnou váhu.

$$S_8(x_1, x_2) = \frac{2a}{2a + b + c}$$

$$S_8(x_1, x_2) = \frac{3a}{3a + b + c}$$

Sokal & Sneath (1963)

- ☑ navržen jako doplněk Rogers & Tanimotova koeficientu (S_2), dává dvojnásobnou váhu rozdílům ve jmenovateli.

$$S_{10}(x_1, x_2) = \frac{a + d}{a + 2b + 2c}$$

Russel & Rao (1940)

- ☑ navržená míra umožňuje porovnání počtu duplicitních prezencí (v čitateli) proti celkovému počtu druhů, nalezených na všech lokalitách, zahrnujícím druhy, které chybějí (d) na obou uvažovaných lokalitách.

$$S_{11}(x_1, x_2) = \frac{a}{p}$$

Kulczyński (1928)

- ☑ koeficient porovnávající duplicitní prezence s diferencemi

$$S_{12}(x_1, x_2) = \frac{a}{b + c}$$

Binární verze asymetrického kvantitativního Kulczyński koeficientu (1928)

- ✓ Mezi svými koeficienty pro presence/absence data zmiňují Sokal & Sneath (1963) tuto verzi kvantitativního koeficientu S_{18} , kde jsou duplicitní prezence srovnávány se součty okrajů tabulky $(a+b)$ a $(a+c)$.

$$S_{13}(x_1, x_2) = \frac{1}{2} \left[\frac{a}{a+b} + \frac{a}{a+c} \right]$$

Ochiachi (1957)

- ☑ použil jako míru podobnosti geometrický průměr poměrů a k počtu druhů na každé lokalitě, tj. se součty okrajů tabulky $(a+b)$ a $(a+c)$, tento koeficient je obdobou S_6 , bez části, týkající se double-zero (d).

$$S_{14}(x_1, x_2) = \sqrt{\frac{a}{(a+b)} \frac{a}{(a+c)}} = \frac{a}{\sqrt{(a+b)(a+c)}}$$

Faith (1983)

- ☑ **V tomto koeficientu je neshoda (přítomnost na jedné a absence na druhé lokalitě) vážena proti duplicitní prezenci. Hodnota S_{26} klesá s růstem double-zero**

$$S_{26}(x_1, x_2) = \frac{a + d / 2}{p}$$

Vícerozměrná analýza dat

Kvantitativní koeficienty

„Klasické“ indexy podobnosti

- ☑ **Sørensenův kvantitativní koeficient**, kde aN a bN jsou celkové počty jedinců v společenstvech A a B, jN je pak suma abundancí pokud se druh nachází v obou společenstvech, je počítána vždy z nižší abundance daného druhu ve společenstvu

$$C_N = \frac{2jN}{(aN + bN)}$$

- ☑ **Morisita-Horn index**, kde aN je celkový počet jedinců ve společenstvu A a an_i počet jedinců druhu i ve společenstvu A (obdobně platí pro společenstvo B)

$$C_{mH} = \frac{2 \sum (an_i \cdot bn_i)}{(da + db) \cdot aN \cdot bN} \quad da = \frac{\sum an_i^2}{aN^2}$$

Jednoduchý srovnávací koeficient (Sokal & Michener, 1958)

- ✓ modifikovaný simple matching coefficient může být použit pro multistavové deskriptory - číselník obsahuje počet deskriptorů, pro které jsou dva objekty ve stejném stavu – např. je-li dvojice objektů popsána následujícími deseti multistavovými deskriptory: hodnota S_1 , vypočítaná pro 10 multistavových deskriptorů bude $S_1(x_1, x_2) = 4 \text{ agreements} / 10 \text{ descriptors} = 0.4$
- ✓ Podobným způsobem je možné rozšířit všechny binární koeficienty pro multistavové deskriptory.

$$S_1(x_1, x_2) = \frac{\text{agreements}}{p}$$

	Deskriptors										Σ
Object x_1	9	3	7	3	4	9	5	4	0	6	
Object x_2	2	3	2	1	2	9	3	2	0	6	
Agreements	0	+	+	+	+	+	+	+	+	+	4

Gowerův obecný koeficient podobnosti (1971)

I.

- ☑ **Gover navrhl obecný koeficient podobnosti, který může kombinovat různé typy deskriptorů. Podobnost mezi dvěma objekty je vypočítána jako průměr podobností, vypočítaných pro všechny deskriptory. Pro každý deskriptor j je hodnota parciální podobnosti s_{12j} mezi objekty x_1 a x_2 vypočítána následovně:**

$$S_{15}(x_1, x_2) = \frac{1}{p} \sum_{j=1}^p s_{12j}$$

- ✓ **Pro binární deskriptory $s_j=1$ (shoda) nebo 0 (neshoda). Gower navrhl dvě formy tohoto koeficientu. Následující forma je symetrická, dává $s_j=1$ double-zero. Druhá forma, Gowerův asymetrický koeficient S_{19} dává pro double-zero $s_j=0$**
- ✓ **Kvalitativní a semikvantitativní deskriptory jsou upraveny podle jednoduchého zaměňovacího pravidla, $s_j=1$ při souhlasu a $s_j = 0$ při nesouhlasu deskriptorů. Double zero jsou ošetřeny stejně jako v předchozím odstavci.**
- ✓ **Kvantitativní deskriptory (reálná čísla) jsou zpracovány následovně: pro každý deskriptor se nejprve vypočte rozdíl mezi stavy obou objektů který je poté vydělen největším rozdílem (R_j), nalezeným pro daný deskriptor mezi všemi objekty ve studii (nebo v referenční populaci – doporučuje se vypočítat největší diferenci R_j každého deskriptoru j pro celou populaci, aby byla zajištěna konzistence výsledků pro všechny parciální studie).**

Gowerův obecný koeficient podobnosti (1971)

II.

- ☑ **normalizovaná vzdálenost může být odečtena od 1 aby byla transformována na podobnost:**

$$s_{12j} = 1 - \left[\frac{|y_{1j} - y_{2j}|}{R_j} \right]$$

- ☑ **Gowerův koeficient může být nastaven tak, aby zahrnoval přídatný flexibilní prvek: žádné porovnání není vypočítáno u deskriptorů, u nichž chybí informace buď u jednoho, nebo u druhého objektu. Toto zajišťuje člen w_j , nazývaný Kroneckerovo delta, popisující přítomnost/nepřítomnost informace v obou objektech: je-li informace o deskriptoru y_j přítomna u obou objektů ($w_j=1$), jinak ($w_j=0$), tento koeficient nabývá hodnot podobnosti mezi 0 a 1 (největší podobnost objektů). Další možností je vážení různých deskriptorů prostým přiřazením čísla v rozsahu 0-1 w_j .**

$$S_{15}(x_1, x_2) = \frac{\sum_{j=1}^p w_{12j} s_{12j}}{\sum_{j=1}^p w_{12j}}$$

Vícerozměrná analýza dat

Různé vícerozměrné metriky vzdáleností

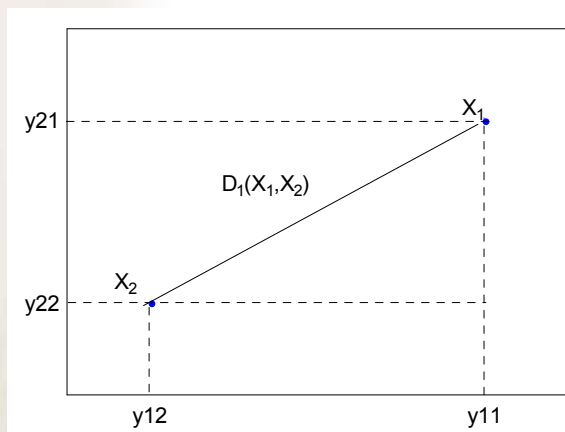
Euklidovská vzdálenost

- ☑ Jde o základní metrické měřítko vzdálenosti a počítá vzdálenost objektů obdobně jako Pythagorova věta počítá přeponu pravoúhlého trojúhelníku. Metoda je citlivá na rozdílný rozsah hodnot vstupujících proměnných (vhodným řešením může být standardizace) a double zero problém. Nemá horní hranici hodnot.

$$D_1(x_1, x_2) = \sqrt{\sum_{j=1}^p (y_{1j} - y_{2j})^2}$$

- ☑ Jako další měřítko se používá také čtverec této vzdálenosti. Jeho nevýhodou jsou semimetrické vlastnosti.

$$D_1^2(x_1, x_2) = \sum_{j=1}^p (y_{1j} - y_{2j})^2$$



Průměrná vzdálenost

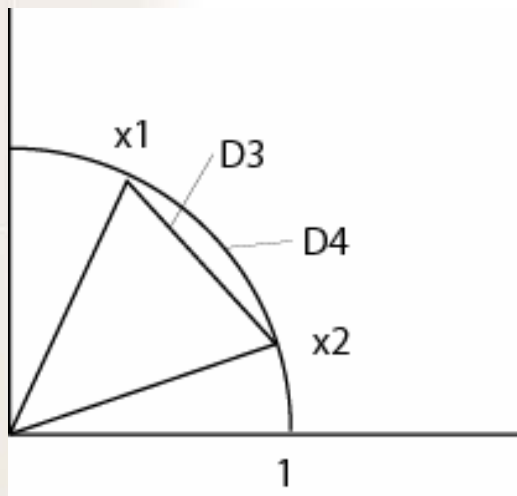
- ☑ Euklidovská vzdálenost je přepočítána na počet parametrů (druhů v případě vzdálenosti společenstev odběrů).

$$D_2^2(x_1, x_2) = \frac{1}{p} \sum_{j=1}^p (y_{1j} - y_{2j})^2$$

$$D_2(x_1, x_2) = \sqrt{D_2^2}$$

Chord distance (Orlóci, 1967)

- ✓ **Odstraňuje double zero problém a vliv rozdílného počtu jedinců druhů ve vzorcích při výpočtu Euklidovské vzdálenosti. Její maximální hodnota je druhá odmocnina ze dvou a minimum 0. Při výpočtu počítá pouze s poměry druhů v rámci jednotlivých vzorků. Jde vlastně o Euklidovskou vzdálenost počítanou pro vektory vzorků standardizované na délku 1, nebo je možný přímý výpočet už zahrnující standardizaci. Vnitřní část výpočtu je vlastně cosinus úhlu svíraného vektory, zápis vzorce je možný i v této formě.**

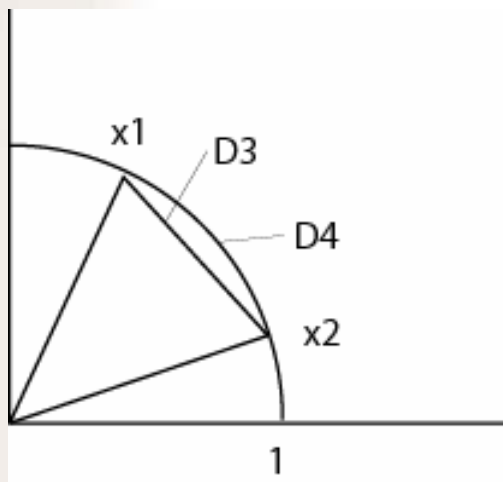


$$D_3(x_1, x_2) = \sqrt{2 \left(1 - \frac{\sum_{j=1}^p y_{1j} y_{2j}}{\sqrt{\sum_{j=1}^p y_{1j}^2 \sum_{j=1}^p y_{2j}^2}} \right)}$$

$$D_3 = \sqrt{2(1 - \cos \theta)}$$

Geodetická metrika

- ✓ **Počítá délku výseče jednotkové kružnice mezi normalizovanými vektory (viz. Chord distance).**



$$D_4(x_1, x_2) = \arccos \left[1 - \frac{D_3^2(x_1, x_2)}{2} \right]$$

Mahalanobisova vzdálenost (Mahalanobis 1936)

- ☑ Jde o obecné měřítko vzdálenosti beroucí v úvahu korelaci mezi parametry a je nezávislá na rozsahu hodnot parametrů. Počítá vzdálenost mezi objekty v systému souřadnic jehož osy nemusí být na sebe kolmé. V praxi se používá pro zjištění vzdálenosti mezi skupinami objektů. Jsou dány dvě skupiny objektů w_1 a w_2 o n_1 a n_2 počtu objektů a popsané p parametry:

$$D_5^2(w_1, w_2) = \overline{d}_{12} V^{-1} \overline{d}_{12}$$

- ☑ Kde \overline{d}_{12} je vektor o délce p rozdílů mezi průměry p parametrů v obou skupinách. V je vážená disperzní matice (matice kovariancí parametrů) uvnitř skupin objektů.

$$V = \frac{1}{n_1 + n_2 - 2} [(n_1 - 1)S_1 + (n_2 - 1)S_2]$$

- ☑ kde S_1 a S_2 jsou disperzní matice jednotlivých skupin. Vektor měří rozdíl mezi p - rozměrnými průměry skupin a V vkládá do rovnice kovarianci mezi parametry.

Minkowskeho metrika

- ☑ Je obecnou formou výpočtu vzdálenosti – podle zadaného koeficientu může odpovídat např. Euklidovské nebo Manhattanské metrice. Se stoupající koeficientem umocňování stoupá významnost větších rozdílů. Existuje ještě obecnější forma, kdy koeficient umocňování a odmocňování je zadáván zvlášť.

$$D_r(x_1, x_2) = \left[\sum_{j=1}^p |y_{1j} - y_{2j}|^r \right]^{1/r}$$

Manhattanská vzdálenost

- ☑ **Jde vlastně o součet rozdílů jednotlivých parametrů popisujících objekty**

$$D_7(x_1, x_2) = \sum_{j=1}^p |y_{1j} - y_{2j}|$$

Mean character difference (Czekanowski 1909)

- ☑ **Manhattanská vzdálenost přepočítaná na počet parametrů.**

$$D_8(x_1, x_2) = \frac{1}{p} \sum_{j=1}^p |y_{1j} - y_{2j}|$$

Whittakerův asociační index (Whittaker 1952)

- ☑ Je dobře použitelný pro data abundancí, každý druh je nejprve transformován ve svůj podíl ve společenstvu, následující výpočet je opět obdobou Manhattané vzdálenosti.

$$D_9(x_1, x_2) = \frac{1}{2} \sum_{j=1}^p \left| \frac{y_{1j}}{\sum_{j=1}^p y_{ij}} - \frac{y_{2j}}{\sum_{j=1}^p y_{2j}} \right|$$

- ☑ Jeho hodnota je 0 v případě identických proporcí druhů. Stejný výsledek lze získat i jako součet nejmenších podílů v rámci obou vzorků.

$$D_9(x_1, x_2) = \left[1 - \min \left(\frac{y_j}{\sum_{j=1}^p y_j} \right) \right]$$

Canberra metric (Lance & Williams 1966)

- ☑ Varianta Manhattané vzdálenosti (před výpočtem musí být odstraněny double zero a není jimy tedy ovlivněna). Stejný rozdíl mezi početnými druhy ovlivňuje vzdálenost méně než mezi druhy vzácnějšími.

$$D_{10}(x_1, x_2) = \sum_{j=1}^p \left[\frac{|y_{1j} - y_{2j}|}{(y_{1j} + y_{2j})} \right]$$

- ☑ Stephenson et al. (1972) a Moreau & Legendre (1979) použili tuto metriku jako součást koeficientu podobnosti

$$S(x_1, x_2) = 1 - \frac{1}{p} D_{10}$$

Koeficient divergence

- ☑ **Obdobná metrika jako D10 ale založená na Euklidovské vzdálenosti a vztažená na počet parametrů.**

$$D_{11}(x_1, x_2) = \sqrt{\frac{1}{p} \sum_{j=1}^p \left(\frac{y_{1j} - y_{2j}}{y_{1j} + y_{2j}} \right)^2}$$

Coefficient of racial likeness (Pearson 1926)

- ☑ Umožňuje srovnávat skupiny objektů podobně jako Mahalanobisova vzdálenost, ale na rozdíl od ní neeliminuje vliv korelace parametrů. Dvě skupiny objektů w_1 a w_2 jsou charakterizovány \bar{y}_j (průměr parametrů ve skupinách) a s_{jj}^2 (rozptyl parametrů ve skupinách).

$$D_{12}(w_1, w_2) = \sqrt{\frac{1}{p} \sum_{j=1}^p \frac{(\bar{y}_{1j} - \bar{y}_{2j})^2}{\left(\frac{s_{1j}^2}{n_1}\right) + \left(\frac{s_{2j}^2}{n_2}\right)}} - \frac{2}{p}$$

χ^2 metrika (Roux & Reyssac 1975)

- ✓ První ze skupiny metrik založených na χ^2 pro výpočet vzdáleností odběrů založených na abundancích druhů nebo jiných frekvenčních datech (nejsou přípustné žádné záporné hodnoty). Data původní matice abundancí/frekvencí Y jsou nejprve přepočítána do matice poměrných frekvencí (součty frekvencí v řádcích (odběry) jsou rovny 1). Jako dodatečné charakteristiky uplatňované při výpočtu jsou spočteny součty řádků y_{i+} a sloupců y_{+j} celé matice $n(i)$ odběrů \times $p(j)$ druhů.

$$Y = \begin{matrix} \begin{bmatrix} y_{ij} \\ \vdots \\ y_{+j} \end{bmatrix} & \begin{bmatrix} y_{i+} \\ \vdots \\ y_{++} \end{bmatrix} \end{matrix} \rightarrow \begin{bmatrix} y_{ij}/y_{i+} \\ \vdots \\ y_{ij}/y_{+j} \end{bmatrix}$$

$$D(x_1, x_2) = \sqrt{\sum_{j=1}^p \left(\frac{y_{1j}}{y_{1+}} - \frac{y_{2j}}{y_{2+}} \right)^2}$$

- ✓ Výpočet odstraňuje problém double zero. Nejjednodušším výpočtem je obdoba Euklidovské vzdálenosti
- ✓ která je dále vážena součty jednotlivých druhů

$$D_{15}(x_1, x_2) = \sqrt{\sum_{j=1}^p \frac{1}{y_{+j}} \left(\frac{y_{1j}}{y_{1+}} - \frac{y_{2j}}{y_{2+}} \right)^2}$$

χ^2 vzdálenost (Lébart & Fénelon 1971)

- ✓ Výpočet je podobný χ^2 metrice, ale vážení je prováděno relativní četností řádku v matici místo jeho absolutního součtu, při výpočtu se užívá parametr y_{++} (celkový součet matice). Je využívána také při výpočtu vztahů řádků a sloupců kontingenční tabulky.

$$D_{16}(x_1, x_2) = \sqrt{\sum_{j=1}^p \frac{1}{y_{+j} / y_{++}} \left(\frac{y_{1j}}{y_{1+}} - \frac{y_{2j}}{y_{2+}} \right)^2} = \sqrt{y_{++}} \sqrt{\sum_{j=1}^p \frac{1}{y_{+j}} \left(\frac{y_{1j}}{y_{1+}} - \frac{y_{2j}}{y_{2+}} \right)^2}$$

Hellingerova vzdálenost (Rao 1995)

☑ Koeficient související s D15 a D16.

$$D_{17}(x_1, x_2) = \sqrt{\sum_{j=1}^p \left[\sqrt{\frac{y_{1j}}{y_{1+}}} - \sqrt{\frac{y_{2j}}{y_{2+}}} \right]^2}$$