

# Výběr proměnných v kompozičních datech

S. Donevska   E. Fišerová   P. Filzmoser   K. Hron

ROBUST 2014

**Kompoziční data (CoDa)** = kvantitativní popis části nějakého celku, tudíž data nesoucí pouze **relativní informaci**.

- **Simplex s Aitchisonovou geometrií** = výběrový prostor CoDa,

$$S^D = \{\mathbf{x} = (x_1, \dots, x_D)', x_i > 0, \sum_{i=1}^D x_i = \kappa\}.$$

- Aitchisonova geometrie není úplně vhodná pro provádění standardních statistických metod.
  - ⇒ Potřeba najít vhodnou reprezentaci CoDa v reálném prostoru.
- Navrženy transformace logaritmu podílů (logratio):  
aditivní logratio (alr) transformace, centrovaná logratio (clr) transformace a i izometrická logratio (ilr) transformace.
- CoDa se reprezentují pomocí podílu složek.

Clr transformace je izometrické zobrazení mezi  $\mathcal{S}^D$  a nadrovinou v  $\mathbb{R}^D$ ,

$$\mathbf{y} = \text{clr}(\mathbf{x}) = (y_1, y_2, \dots, y_D)' = \left( \ln \frac{x_1}{\sqrt[D]{\prod_{i=1}^D x_i}}, \dots, \ln \frac{x_D}{\sqrt[D]{\prod_{i=1}^D x_i}} \right)'. \quad (1)$$

- Nevýhody clr proměnných:

- nejsou souřadnice vzhledem k bázi na simplexu,
- vedou k singulární varianční matici ( $y_1 + \dots + y_D = 0$ ),
- nejsou subkompozičně soudržné.

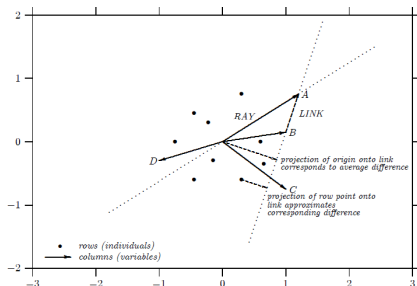
- Výhody clr proměnných:

- Převádí operace perturbace a mocninné transformace CoDa do obyčejného součtu a násobení skalárem clr koeficientů,
- Euklidovská vzdálenost mezi vektory clr proměnných = Aitchisonova vzdálenost jejich příslušných kompozic. Platí i pro skalární součin a normu.

# Interpretace CoDa biplotu

**Kompoziční biplot** je získaný jako standardní biplot pro clr transformovaná data.

- 1 Vzdálenost mezi vrcholy šipek jsou úměrné směrodatné odchylce logaritmu podílu příslušných složek.
- 2 Šipky představují clr proměnné.
- 3 Délka šipky je úměrná směrodatné odchylce clr proměnné.
- 4 Kosinus úhlu mezi dvěma vrcholy šipek dává hodnotu korelačního koeficientu mezi odpovídajícími logaritmy podílu.



# Míry variability pro CoDa

Základní míra variability náhodné kompozice  $\mathbf{x} = (x_1, \dots, x_D)'$  je **variační matice** definovaná jako

$$\mathbf{T} = \left\{ \text{var} \left( \ln \frac{x_i}{x_j} \right) \right\}_{i,j=1}^D.$$

- $\mathbf{T}$  je symetrická matice, která má na hlavní diagonále 0.
- Prvky matice  $\mathbf{T}$  popisují variabilitu logaritmu podílu složek  $x_i$  a  $x_j$ .

(Normovaný) součet prvků variační matice  $\mathbf{T}$  se nazývá **celková variance**,

$$\text{totvar}(\mathbf{x}) = \frac{1}{2D} \sum_{i=1}^D \sum_{j=1}^D \text{var} \left( \ln \frac{x_i}{x_j} \right),$$

vyjadřuje celkovou variabilitu v kompozičním datovém souboru.

- Platí, že  $\text{totvar}(\mathbf{x}) = \sum_{i=1}^D \text{var}(y_i)$ , kde

$$\text{var}(y_i) = \frac{D-1}{D^2} \sum_{j=1}^D \text{var} \left( \ln \frac{x_j}{x_i} \right) - \frac{1}{2D^2} \sum_{\substack{j=1 \\ j \neq i}}^D \sum_{\substack{l=1 \\ l \neq i}}^D \text{var} \left( \ln \frac{x_j}{x_l} \right),$$

⇒ Existuje celkem silný vztah mezi  $\text{var}(y_i)$  a součtem prvků v  $i$ -tém řádku (sloupce) příslušné variační matici  $\mathbf{T}$ .

## Věta

Nechť jsou dány clr proměnné  $y_i$  a  $y_j$ ,  $i \neq j$ ,  $i, j = 1, \dots, D$ . Potom,  $\text{var}(y_i) \geq \text{var}(y_j)$ , právě tehdy, když

$$\sum_{p=1}^D \text{var} \left( \ln \frac{x_i}{x_p} \right) \geq \sum_{p=1}^D \text{var} \left( \ln \frac{x_j}{x_p} \right).$$

# Navržená kroková procedura

Uvažujme kompozici  $\mathbf{x} = (x_1, \dots, x_D)'$ , takovou, že

$$\text{var}(y_1) \geq \dots \geq \text{var}(y_D) \quad (2)$$

$\Leftrightarrow$

$$\sum_{p=1}^D \text{var} \left( \ln \frac{x_1}{x_p} \right) \geq \sum_{p=1}^D \text{var} \left( \ln \frac{x_2}{x_p} \right) \geq \dots \geq \sum_{p=1}^D \text{var} \left( \ln \frac{x_D}{x_p} \right). \quad (3)$$

## Algoritmus:

- 1 Vynecháme složku  $x_D$ , jejíž rozptyl odpovídající clr proměnné je nejmenší.  
Uvažujme tedy podkompozici  $\mathbf{x}_1 = (x_1, \dots, x_{D-1})'$  a provedeme clr transformaci na  $\mathbf{x}_1$ .  
Vypočítáme rozptyly clr transformovaných proměnných z  $\mathbf{x}_1$ .
- 2 Opakujeme krok 1.
- 3 STOP nejpozději po  $D - 2$  krocích.

# Navržená kroková procedura

- Bude dodrženo pořadí variancí clr proměnných při přechodu z  $D$ -složkové k  $D - 1$ - složkové podkompozici?

⇒ Ano, ale jen za splnění předpokladu

$$\text{var} \left( \ln \frac{x_1}{x_D} \right) \geq \text{var} \left( \ln \frac{x_2}{x_D} \right) \geq \dots \geq \text{var} \left( \ln \frac{x_{D-1}}{x_D} \right).$$

- Kdy máme skončit s výběrem kompozičních složek?
- ⇒ Poté co použijeme stop kritérium, které porovnává celkovou varianci podkompozici  $\mathbf{x}_i$ , získanou v  $i$ -tém kroku algoritmu  $i = 1, \dots, D - 2$ , s celkovou variancí podkompozici  $\mathbf{x}_{i-1}$  z předchozího kroku.



# Navržená kroková procedura – STOP kritérium

$H_0 : \text{totvar}(\mathbf{x}_i) = \text{totvar}(\mathbf{x}_{i-1})$  v.s.  $H_A : \text{totvar}(\mathbf{x}_i) < \text{totvar}(\mathbf{x}_{i-1})$

- Za předpokladu platnosti nulové hypotézy

$$U_i^+ = \frac{\widehat{\text{totvar}}(\mathbf{x}_i) - \text{totvar}(\mathbf{x}_{i-1})}{\sqrt{\frac{2}{n-1} \text{tr}(\widehat{\Sigma}_i^2)}} \sim N(0, 1), \text{ pro } i \rightarrow \infty,$$

kde  $\widehat{\Sigma}_i$  je výběrová varianční matice kompozice  $\mathbf{x}_i$  v (libovolných) ilr souřadnicích.

- $H_0$  se zamítá na hladině významnosti  $\alpha$  jestliže

$$u_i^+ \in W = (-\infty, u_\alpha),$$

kde  $u_i^+$  je realizace  $U_i^+$  a  $u_\alpha$  je  $\alpha$ - kvantil normovaného normálního rozdělení.

- V každém kroku spočítáme  $U_i^+$  proceduru ukončíme když  $u_i^+ \in W$ .

Datový soubor Kola je výsledkem velkého geochemického projektu, který byl prováděn od 1992 do 1998 Geologickým průzkumným ústavem Finska a Norska, Centrální Kola expedicí, Rusko.

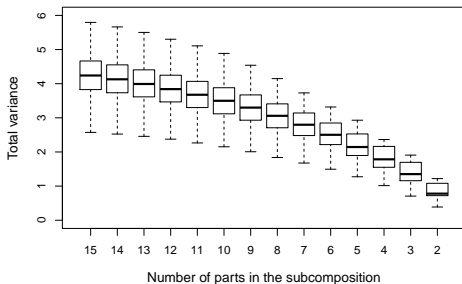
- Zkoumaná oblast: 188 000  $km^2$  na poloostrově Kola v Severní Evropě.
- Celkem je odebráno cca 600 vzorků ze 4 různých půdních horizontů (horizont O, horizont A, horizont B, horizont C).
- Analyzovaný jsou více než 50 chemických prvků v půdních horizontech.
- Data jsou dostupná v balíčku StatDA softwaru R (R Development Core Team, 2012).

# Příklad – Data Kola – I. Experiment

- Náhodně vybereme 15 proměnných z 30 prvků z horizontu O.
- Algoritmus je aplikovaný dokud nedocílíme 2-složkovou podkompozici.
- V každém kroku je vypočítaná celková variance.
- Celá procedura je opakovaná 1000 krát.

# Příklad – Data Kola – I. Experiment

- Náhodně vybereme 15 proměnných z 30 prvků z horizontu O.
- Algoritmus je aplikovaný dokud nedocílíme 2-složkovou podkompozici.
- V každém kroku je vypočítaná celková variance.
- Celá procedura je opakovaná 1000 krát.



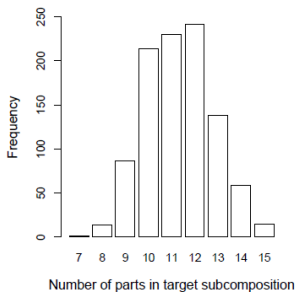
Obr: Celková variance podkompozic obdržena z postupné krokové procedury.

## Příklad – Data Kola – II. experiment

- Znovu vybíráme náhodně 15 proměnných z 30 prvků z horizontu  $O$ .
- Postupná procedura je aplikovaná dokud testová statistika nenavrhne zastavení procesu.
- Celá procedura je opakovaná 1000 krát.

# Příklad – Data Kola – II. experiment

- Znovu vybíráme náhodně 15 proměnných z 30 prvků z horizontu O.
- Postupná procedura je aplikovaná dokud testová statistika nenavrhne zastavení procesu.
- Celá procedura je opakovaná 1000 krát.



Obr: Sloupcový graf počtů složek výsledných podkompozic z postupné krokové procedury užitím STOP kritéria.

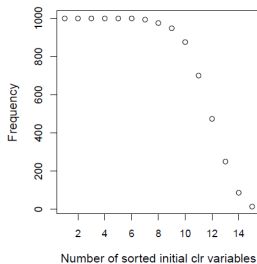
Obsahuje výsledná kompozice složky, které měly vysokou clr varianci v původní D-složkové kompozici?

- Složky ze všech 1000 počátečních podkompozic jsou vybrané dle jejich klesajících hodnot clr variací.
- Počítáme, jak často jsou top  $k$  clr proměnné zahrnuty do výsledné kompozice, kde  $k = 1, \dots, 15$ .

# Příklad – Data Kola – II. experiment

Obsahuje výsledná kompozice složky, které měly vysokou clr varianci v původní D-složkové kompozici?

- Složky ze všech 1000 počátečních podkompozic jsou vybrány dle jejich klesajících hodnot clr variací.
- Počítáme, jak často jsou top  $k$  clr proměnné zahrnuty do výsledné kompozice, kde  $k = 1, \dots, 15$ .



Obr: Clr proměnné z původní kompozice uspořádané dle klesajících variací, vs. počet, kolikrát příslušné kompoziční složky byly zahrnuty ve výsledné podkompozici.

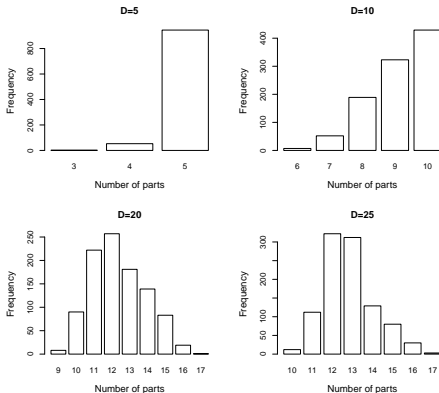


## Příklad – Data Kola – III. experiment

- Použijeme stejné simulační nastavení jako před tím, ale vybereme 5, 10, 20 a 25 - složkové počáteční kompozice z datového souboru Kola.
- Cela procedura je opakovaná 1000 krát.

# Příklad – Data Kola – III. experiment

- Použijeme stejné simulační nastavení jako před tím, ale vybereme 5, 10, 20 a 25 - složkové počáteční kompozice z datového souboru Kola.
- Cela procedura je opakovaná 1000 krát.



Obr: Sloupcové grafy počtů složek výsledné podkompozice z krokové procedury užitím STOP – kritéria.

# Příklad – Baltický průzkum půdy

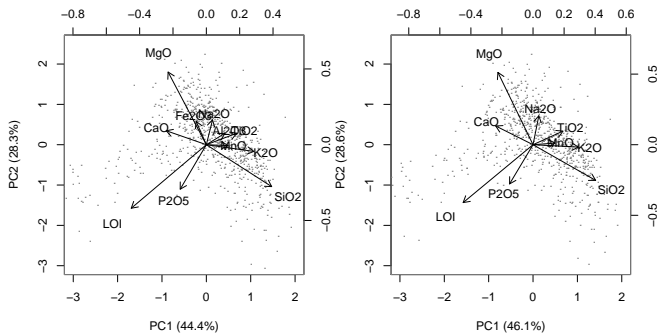
BSS datový soubor (Reimann et al., 2003), je získány z rozsáhlého geochemického projektu který, je realizovaný v Severní Evropě v oblasti rozsahu cca 1 800 000 km<sup>2</sup>.

- 769 vzorků zemědělské půdy.
- Vzorky pocházejí ze dvou vrstev: horní vrstva (0-25cm) a spodní vrstva (50-75 cm).
- U všech vzorků je analyzovaná koncentrace z více než 40 chemických sloučenin.
- Používáme významné sloučeniny ( $\text{Al}_2\text{O}_3$ ,  $\text{Fe}_2\text{O}_3$ ,  $\text{K}_2\text{O}$ ,  $\text{MgO}$ ,  $\text{MnO}$ ,  $\text{CaO}$ ,  $\text{TiO}_2$ ,  $\text{Na}_2\text{O}$ ,  $\text{P}_2\text{O}_5$  a  $\text{SiO}_2$ ), plus LOI (Loss on ignition) z horní vrstvy tj., 11-ti složkovou CoDa.
- Datový soubor z horní a dolní vrstvy půdy je přístupný v balíčku `mvoutlier` v R.






# Příklad – Baltický průzkum půdy

Vliv postupné procedury je zobrazen pomocí CoDa biplotu na BSS datech.

- Při použití testového kritéria pro  $\alpha = 0.05$  jsme získali 9-složkovou podkompozici ( $\text{Al}_2\text{O}_3$  a  $\text{Fe}_2\text{O}_3$  jsou vyloučeny).



Obr: Biplot pro BSS data pro všechny významné sloučeniny (vlevo) a po odstranění  $\text{Al}_2\text{O}_3$  a  $\text{Fe}_2\text{O}_3$  (vpravo).

-  Aitchison J (1986) The statistical analysis of compositional data. Chapman and Hall, London.
-  Egozcue JJ (2009) Reply to "On the Harker Variation Diagrams; ..." by J. A. Cortés. Math Geosci 41:829–834.
-  Filzmoser P, Hron K, Reimann C (2012) Interpretation of multivariate outliers for compositional data. Computers & Geosciences 39: 77–85.
-  Hron K, Filzmoser P, Donevska S, Fišerová E (2013) Covariance based variable selection for compositional data. Mathematical geosciences 45:487–498.
-  Hron K, Kubáček L (2011) Statistical properties of the total variation estimator for compositional data. Metrika 74: 221–230.