

Metoda dílčích nejmenších čtverců pro kompoziční data s aplikací v metabolomice

Alžběta Kalivodová^{a,b}

Karel Hron^{a,b}, Peter Filzmoser^c, Lukáš Najdekr^d

^a Katedra matematické analýzy a aplikací matematiky

^b Katedra geoinformatiky

PŘF UPOL

^c Department of Statistics and Probability Theory

Vienna University of Technology

^d Laboratoř metabolomiky

Ústav molekulární a translační medicíny, UPOL a FNOL

Robust 2014, 19. - 24.1. 2014

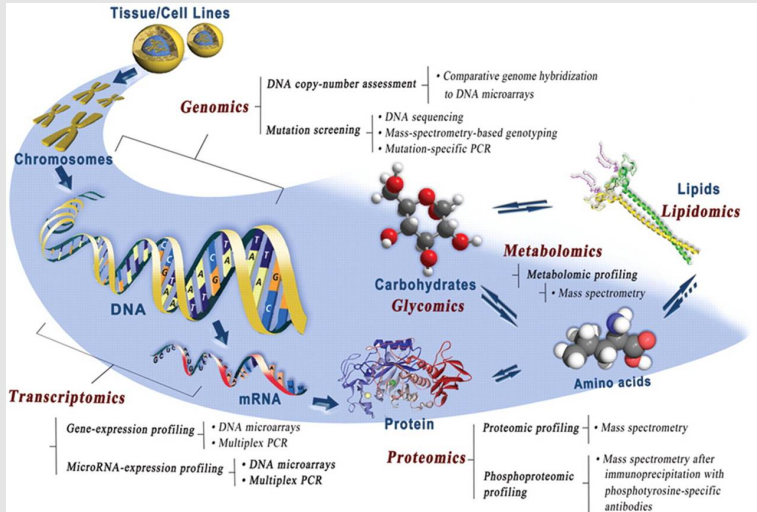
Obsah

- 1 Úvod
- 2 Metabolomika
- 3 Kompoziční data
- 4 PLS-DA
- 5 Praktický příklad
- 6 Závěr

Úvod

- Metabolomická data mají speciální strukturu.
- Na metabolomická data aplikujeme přístup tzv. kompozičních dat.
- Speciální struktura dat → užití metody dílčích nejmenších čtverců - diskriminační analýzy (PLS-DA).
- Teoretické aspekty jsou dále ukázány na konkrétních datových souborech.
- Porovnání standardního a kompozičního přístupu.

Metabolomika



(Wu et al., 2011)

Metabolomika

- **Metabolom** - soubor organických sloučenin, které jsou obsaženy v daném biologickém materiálu. Jejich velikost je na úrovni molekul.
- **Metabolity** - molekuly tvořící metabolom.
- **Metabolomika** - analýza metabolomu za daných podmínek.
- Metabolity jsou měřeny na biologických materiálech (buňky, krev, moč, ...) a nesou pouze relativní informaci.

Metabolomický datový soubor

- Pozorování - **pacienti vs. kontroly** . Proměnné - jednotlivé **metabolity**.
- Velký datový soubor.
- Více proměnných než pozorování.
- Malý počet pacientů (nejvíce desítky, spíše méně). × Velký počet metabolitů (stovky) → užití vhodných metod.

Kompoziční data

- Vícerozměrný statistický soubor, jehož prvky jsou kvantitativně vyjádřené části celku.

Definice

Sloupcový vektor $\mathbf{x} = (x_1, \dots, x_D)^T$ se nazývá D -složková kompozice, jestliže jsou všechny jeho prvky kladná reálná čísla nesoucí pouze relativní informaci.

- Možnost reprezentovat kompozice jako data s konstantním součtem (k).
- Např. procenta ($k = 100$) nebo podíly na celku velikosti k (nejčastěji $k = 1$).
- **Příklad:** koncentrace chemických sloučenin v horninách.

Centrované logratio (clr) souřadnice

- Geometrie pro práci s kompozičními daty - tzv. Aitchisonova geometrie.
- Převod na známou euklidovskou geometrii přes tzv. logratio souřadnice.
- Clr souřadnice (centred logratio coordinations) - zachovává vzdálenosti mezi daty, vede k singulární varianční matici.
- Clr souřadnice vyjádřené po složkách:

$$\text{clr}(\mathbf{x}) = \mathbf{y} = \left(\ln \frac{x_1}{g(\mathbf{x})}, \ln \frac{x_2}{g(\mathbf{x})}, \dots, \ln \frac{x_D}{g(\mathbf{x})} \right)^T,$$

kde $g(\mathbf{x}) = \sqrt[D]{\prod_{i=1}^D x_i}$ je geometrický průměr složek kompozice \mathbf{x} .

Isometrické logratio (ilr) souřadnice

- Principem je vytvořit ortonormální bázi pro kompozici $\mathbf{x} = (x_1, \dots, x_D)'$ na simplexu (výběrovém prostoru).

Ilr souřadnice

$(D - 1)$ -dimenzionální reálný vektor $ilr(\mathbf{x}) = \mathbf{z} = (z_1, \dots, z_{D-1})'$, jehož složky jsou definovány jako

$$z_i = \sqrt{\frac{D-i}{D-i+1}} \ln \frac{x_i}{\sqrt[D-i]{\prod_{j=i+1}^D x_j}}, \quad i = 1, \dots, D-1. \quad (1)$$

- Proměnná z_1 nese všechnu relevantní informaci o části kompozice x_1 , vysvětluje všechny podíly mezi x_1 a ostatními částmi \mathbf{x} .

Isometrická logratio (ilr) souřadnice

- Chceme vytvořit ortonormální bázi, ve které první ilr souřadnice vysvětluje veškerou důležitou informaci o zvolené složce.
- Konstruujeme D různých ilr souřadnic, kde je D -tice (x_1, \dots, x_D) v (1) nahrazena pro $l = 1, \dots, D$ pomocí $(x_l, x_1, \dots, x_{l-1}, x_{l+1}, \dots, x_D) =:$
 $(x_1^{(l)}, x_2^{(l)}, \dots, x_l^{(l)}, x_{l+1}^{(l)}, \dots, x_D^{(l)})$ [3].

$$z_i^{(l)} = \sqrt{\frac{D-i}{D-i+1}} \ln \frac{x_i^{(l)}}{\sqrt[D-i]{\prod_{j=i+1}^D x_j^{(l)}}}, \quad i = 1, \dots, D-1. \quad (2)$$

Vztah mezi clr a ilr souřadnicemi

$$\mathbf{y} = \mathbf{V}\mathbf{z}. \quad (3)$$

Matice $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_{D-1})$ má dimenzi $D \times (D - 1)$ a její sloupce jsou tvořeny clr souřadnicovými vektory, které tvoří ortonormální bázi (vzhledem k Aitchisonově geometrii). Pro $i = 1, \dots, D - 1$

$$\mathbf{v}_{D-i} = \sqrt{\frac{D-i}{D-i+1}} \left(0, \dots, 0, 1, -\frac{1}{D-i}, \dots, -\frac{1}{D-i} \right)'. \quad (4)$$

PLS-DA

- Metoda dílčích nejmenších čtverců - lineární diskriminace (partial least squares regression - discriminant analysis).
- Třída metod, která slouží k modelování vztahů mezi pozorováními (vysvětlujícími proměnnými) a skupinou tzv. latentních proměnných.
- Kombinace metody hlavních komponent a mnohonásobné regrese.
- **Cíl** - predikovat množinu závislých proměnných na základě informace ze skupiny latentních proměnných (prediktorů). Predikce je prováděna extrakcí ortogonálních faktorů z prediktorů.
- **Cíl** - maximalizovat kovarianci mezi skupinami závislých a latentních proměnných.

PLS-DA

- Závislé a latentní proměnné - blokové matice, centrovány vzhledem k Aitchisonově geometrii.
- Regrese je prováděna na latentních proměnných.
- PLS1 - pouze jedna závisle proměnná, PLS2 - více závislých proměnných (užívá se častěji).

Datové matice $\mathbf{X}_{n \times D}$ a $\mathbf{Y}_{n \times q}$. Matice \mathbf{X} je transformována na \mathbf{Z} pomocí ilr souřadnic (2). Výsledkem je lineární vztah

$$\mathbf{Y} = \mathbf{Z}\mathbf{\Gamma} + \mathbf{E}, \quad (5)$$

kde $\mathbf{\Gamma}$ je $(D - 1) \times q$ dimenzionální matice regresních koeficientů a \mathbf{E} je matice chyb. \mathbf{Y} je tvořena binárními proměnnými reprezentujícími příslušnost ke skupinám.

PLS-DA

Místo řešení rovnice (5) jsou matice \mathbf{Z} a \mathbf{Y} vyjádřeny pomocí latentních proměnných:

$$\begin{aligned}\mathbf{Z} &= \mathbf{TP}^T + \mathbf{E}_X, \\ \mathbf{Y} &= \mathbf{UQ}^T + \mathbf{E}_Y.\end{aligned}$$

- \mathbf{E}_X , \mathbf{E}_Y - matice reziduí
- \mathbf{T} , \mathbf{U} - matice skóru
- \mathbf{P} , \mathbf{Q} - matice zátěží
- Všechny matice mají a sloupců, $a \leq \min(D, q, n)$ je počet PLS komponent.

PLS-DA

- Kovariance mezi x a y je maximalizována.
- Kovariance mezi \mathbf{t} a \mathbf{u} : $\text{cov}(\mathbf{t}, \mathbf{u}) = \mathbf{t}^T \mathbf{u} / (n - 1)$, za podmínky $\|\mathbf{t}\| = \|\mathbf{u}\| = 1$.
- Vektory vah \mathbf{w} a \mathbf{c} : $\mathbf{t} = \mathbf{Z}\mathbf{w}$, $\mathbf{u} = \mathbf{Y}\mathbf{c}$.

Maximalizace:

$$\text{cov}(\mathbf{t}, \mathbf{u}) = \text{cov}(\mathbf{Z}\mathbf{w}, \mathbf{Y}\mathbf{c}) \rightarrow \max_{\|\mathbf{t}\|=\|\mathbf{u}\|=1} .$$

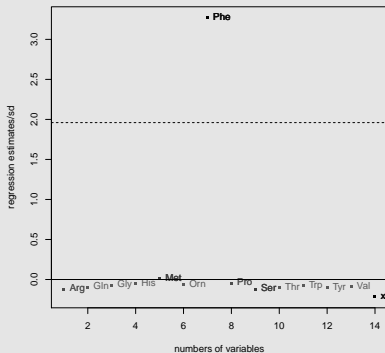
Řešení: první dva vektory skóru \mathbf{t}_1 and \mathbf{u}_1 .

PLS-DA - další postup

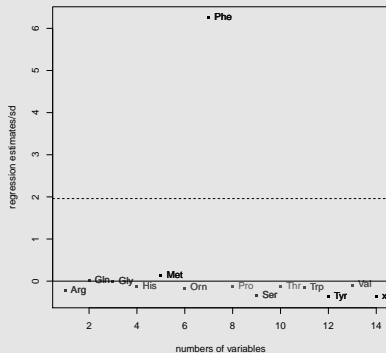
- Vyšetřování statistické významnosti jednotlivých regresních parametrů pomocí metody **bootstrap**.
- V každém kroku je proveden náhodný výběr s opakováním prvků z jednotlivých skupin.
- Velikost výběru = počet prvků ve skupinách.
- Na tomto výběru je následně proveden odhad regresních parametrů.
- Postup je opakován, následně jsou zjištěny směrodatné odchylky výsledných odhadů parametrů.

Praktický příklad I.

Standardized regression estimates – standard approach

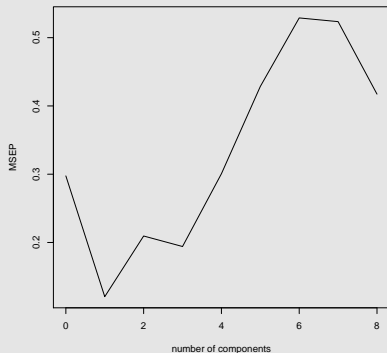


Standardized regression estimates – logratio approach

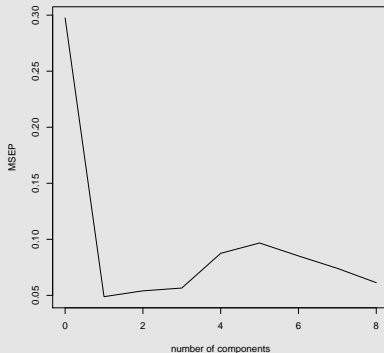


Praktický příklad I.

MSEP – standard approach

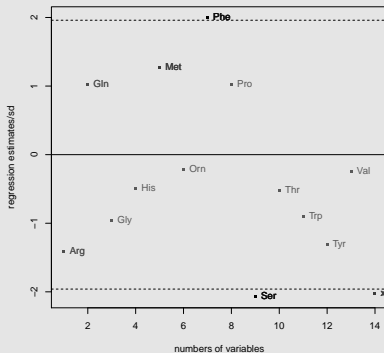


MSEP – logratio approach

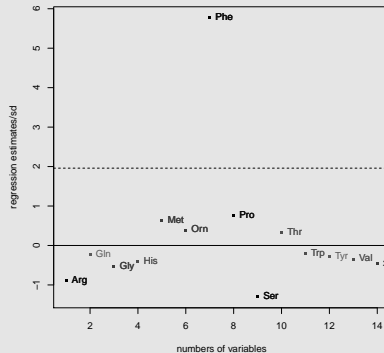


Praktický příklad I.

Standardized regression estimates – standard approach

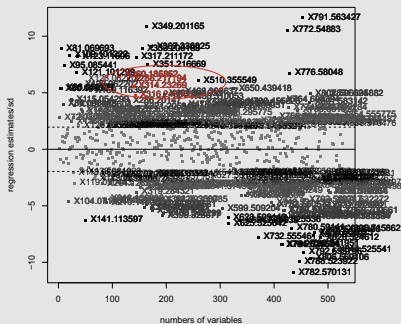


Standardized regression estimates – logratio approach

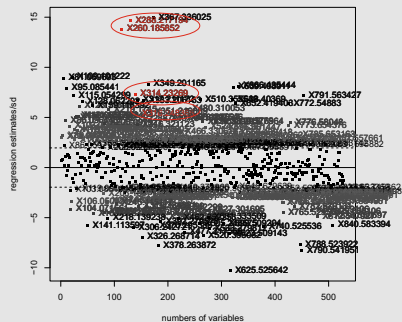


Praktický příklad II.

Standardized regression estimates – standard approach







Standardized regression estimates – logratio approach





Závěr

- Při zpracování metabolomických dat je potřeba vzít v potaz jejich specifické vlastnosti.
- Pro praktické účely je třeba zajistit interpretovatelné ortonormální souřadnice.
- Úprava PLS-DA pro kompoziční data probíhá v logratio souřadnicích.
- Kompoziční přístup k PLS-DA dává velmi dobré výsledky - ukazuje signifikantní metabolity lépe než standardní metodika.

Literatura

-  Roux, A., Lison, D., Junot, Ch. and Heilier, J. (2011) *Applications of liquid chromatography coupled to mass spectrometry-based metabolomics in clinical chemistry and toxicology: A review*. *Clinical Biochemistry* 44 , 119–135.
-  Aitchison, J. (1986) *The Statistical Analysis of Compositional Data*. Chapman & Hall, London.
-  Hron, K., Filzmoser, P., Thompson, K. (2012) *Linear regression with compositional explanatory variables*. *Journal of Applied Statistics* 39 (5), 1115–1128.
-  Egozcue, J. J., Pawlowsky-Glahn, V. (2005) *Groups of Parts and Their Balances in Compositional Data Analysis*. *Mathematical Geology*, 37 (7), 795–828.

Literatura

-  Filzmoser, P., Hron, K., Reimann, C. (2012) *Interpretation of multivariate outliers for compositional data*. *Computers & Geosciences* 39, 77–85.
-  Roux, A., Lison, D., Junot, Ch., Heilier, J. (2011) *Applications of liquid chromatography coupled to mass spectrometry-based metabolomics in clinical chemistry and toxicology: A review*. *Clinical Biochemistry* 44 , 119–135.