

# Subselect 0.9-99: Selecting variable subsets in multivariate linear models

A. PEDRO DUARTE SILVA<sup>(1)(\*)</sup>

JORGE CADIMA<sup>(2)</sup>

MANUEL MINHOTO<sup>(3)</sup>

JORGE ORESTES CERDEIRA<sup>(2)</sup>

**THE PROBLEM:** Finding a k-variable subset that is a good surrogate for a full p-variable data set

**CONTEXT:**

Exploratory data analysis – Subselect 0.1-- 0.9

Multivariate Linear Models – Subselect 0.9-99

(Cadima, Cerdeira, Duarte Silva and Minhoto -- useR! 2004)

## A LINEAR HYPOTHESIS FRAMEWORK

$$X = A \Psi + U$$

$$H_0: C \Psi = 0$$

- SELECT COLUMNS OF X IN ORDER TO EXPLAIN H1

## PARTICULAR CASES:

CANONICAL CORRELATION ANALYSIS /

LINEAR REGRESSION ANALYSIS

$$A = [1 \mid Y] \quad C = [0 \mid I]$$

LINEAR DISCRIMINANT ANALYSIS

$$A = [1_g] \quad \Psi = [\mu_g] \quad C = \begin{bmatrix} 1 & -1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 1 & 0 & \dots & -1 \end{bmatrix}$$

MULTI-WAY MANOVA/MANCOVA EFFECTS

## COMPARISON CRITERIA:

MULTIVARIATE INDICES

$$ccr_1^2$$

$$\zeta^2 = 1 - \frac{r}{\sum_{i=1}^r (1 - ccr_i^2)^{-1}}$$

$$\tau^2 = 1 - \left( \prod_{i=1}^r (1 - ccr_i^2) \right)^{1/r}$$

$$\xi^2 = \frac{\sum_{i=1}^r ccr_i^2}{r}$$

$$\max ccr_1^2 \Leftrightarrow$$

$$\max \zeta^2 \Leftrightarrow$$

$$\max \tau^2 \Leftrightarrow$$

$$\max \xi^2 \Leftrightarrow$$

max Roy first root

max Lawley-Hotelling trace,

min Wilks  $\Lambda$

max Bartlett-Pillai trace,

$$ccr_i^2 = \text{Eigval}_i(T^{-1}H)$$

$$T = X' (I - P_\omega) X$$

$$\Omega = \mathcal{R}(A)$$

$$\omega = \mathcal{R}(A) \cap \mathcal{N}(C)$$

$$r = \dim(\Omega) - \dim(\omega)$$

$$H = X' (P_\Omega - P_\omega) X$$

## THE SUBSELECT PACKAGE

Search routines for (combinatorial) criteria optimization

Exact Algorithm:

leaps - based on Furnival and Wilson's leaps and bounds algorithm

- viable with up to 30 - 35 original variables

Heuristics:

anneal - simulated annealing

genetic - genetic algorithm

improve - restricted local improvement