

**DISCARDING VARIABLES IN PRINCIPAL COMPONENT ANALYSIS:
ALGORITHMS FOR ALL-SUBSETS COMPARISONS
BASED ON THE RV COEFFICIENT**

ANTÓNIO PEDRO DUARTE SILVA



**FACULDADE DE CIÊNCIAS ECONÓMICAS E EMPRESARIAIS
UNIVERSIDADE CATÓLICA PORTUGUESA
CENTRO REGIONAL DO PORTO**

Discarding and/or Selecting Variables in PCA

WHY DISCARD VARIABLES ?

- **PARSIMONITY**
- **REDUCING COSTS OF DATA COLLECTION**
- **HELPING THE INTERPRETATION OF A CLASSICAL PCA**
 - **THE TRADITIONAL APPROACH OF IGNORING VARIABLES WITH SMALL LOADINGS, IMPLICITLY DISCARDS THEM**
 - **ASSESSING VARIABLE IMPORTANCE ONLY BASED ON LOADINGS CAN BE MISLEADING (CADIMA AND JOLLIFFE 1995)**
- **AS A DIRECT APPROACH TO THE PROBLEM OF DIMENSIONALITY REDUCTION (McCabe 1984)**

CRITERIA FOR COMPARING VARIABLE SUBSETS

McCABE APPROACH (PRINCIPAL VARIABLES 1984):

$$\begin{array}{ccccccc} \mathbf{Y} = \mathbf{A} \mathbf{X} ; & \mathbf{Z} = \mathbf{B} \mathbf{Y} & & (\mathbf{A} = [\mathbf{I}_1 \mid \mathbf{0}] ; \mathbf{X} = [\mathbf{X}_1 \mid \mathbf{X}_2] \\ & & & \Rightarrow \mathbf{Y} = \mathbf{X}_1) \\ \mathbf{(n \times k)} & \mathbf{(n \times p)} & \mathbf{(n \times p)} & \mathbf{(n \times k)} \end{array}$$

$$(1) \text{ MIN } |\mathbf{Z} - \mathbf{X}| \Leftrightarrow \text{ MIN } |\mathbf{S}_{22|1}| = \prod \theta_i$$

$$(2) \text{ MIN } \text{Tr} (\mathbf{Z} - \mathbf{X}) \Leftrightarrow \text{ MIN } \text{Tr} \mathbf{S}_{22|1} = \sum \theta_i$$

$$(3) \text{ MIN } \|\mathbf{Z} - \mathbf{X}\|^2 \Leftrightarrow \text{ MIN } \|\mathbf{S}_{22|1}\|^2 = \sum \theta_i^2$$

$$(4) \text{ MIN } (\mathbf{Z} - \mathbf{X})^T \mathbf{S}_X^{-1} (\mathbf{Z} - \mathbf{X}) \Leftrightarrow \text{ MIN } \text{Tr} (\mathbf{S}_{22}^{-1} \mathbf{S}_{22|1}) = (p-k) - \sum \rho_i^2$$

$$(\mathbf{S}_{22|1} = \mathbf{S}_{22} - \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12})$$

CRITERIA BASED ON CORRELATIONAL MEASURES :

$$r_1(A,B) = \text{Tr}(A^T B) / [\text{Tr}(A^T A) \text{Tr}(B^T B)]^{1/2}$$

$$A = P_A \Lambda_A Q_A^T \quad B = P_B \Lambda_B Q_B^T$$

YANAI' S GENERALIZED CORRELATION COEFFICIENT

$$\text{GCD}(A,B) = r_1(P_A P_A^T, P_B P_B^T)$$

$$\text{MAX GCD}(X_1, \text{PC}_q(X)) \quad (\text{Cadima and Jolliffe 1996})$$

ESCOUFIER'S RV COEFFICIENT

$$\text{RV}(A,B) = r_1(P_A \Lambda_A^2 P_A^T, P_B \Lambda_B^2 P_B^T)$$

$$\text{MAX RV}(X, M X_1) \quad (\text{Robert and Escoufier 1976})$$

ALL-SUBSETS COMPARISONS IN LINEAR REGRESSION

$$y = X' \hat{\beta} + e$$

FURNIVAL'S ALGORITHM (1971):

$$\begin{bmatrix} X'X & X'y \\ y'X & y'y \end{bmatrix} \quad \longrightarrow \quad \begin{bmatrix} -(X'X)^{-1} & \hat{\beta} \\ \hat{\beta}' & e'e \end{bmatrix}$$

ELEMENTARY MATRIX OPERATIONS (SYMMETRIC SWEEPS)

COMMENTS:

- MATRIX SIMMETRY IS ALWAYS PRESERVED
- ONLY THE RIGHT-LOWER CORNER NEEDS TO BE UPDATED AT EACH STEP
- BY SEQUECING THE SUBSETS EVALUATIONS PROPERLY:
 - 1/2 OF THE SWEEPS UPDATE (1*1) SUBMATRICES
 - 1/4 OF THE SWEEPS UPDATE (2*2) SUBMATRICES
 - ...
 - ONLY ONE SWEEP UPDATES THE FULL (P*P) MATRIX

**TOTAL NUMBER OF FLOATING
POINT OPERATIONS:**

$$\approx 6 * 2^p$$

$$(6 (2^p) - p(p+7)/2 - 6)$$

FURNIVAL AND WILSON'S "LEAPS AND BOUNDS"
(1974):

$$S_A \subset S_B \Rightarrow (e^T e)_A \geq (e^T e)_B$$

DOUBLE TREE SEARCH

ϕ

X_1, X_2, \dots, X_p

↓ ↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓ ↓

ADDING VARIABLES

REMOVING VARIABLES

CREATE BOUNDS

$e^T e > \text{BOUND}$

\Rightarrow **PRUNE BRANCH**

ALL-SUBSETS COMPARISONS BASED ON THE RV COEFFICIENT

$$\begin{aligned} \text{MAX RV}(X, M X_1) &\Leftrightarrow \text{MAX Tr} \{ [(S_{11})^{-1} (S^2)_{11}]^2 \} = \\ &= \sum_{i,j} \{ [(S_{11})^{-1} (S^2)_{11}](i,j) * [(S_{11})^{-1} (S^2)_{11}](j,i) \} \end{aligned}$$

SYMMETRIC SWEEPING OF COVARIANCE MATRICES:

$$A = \begin{bmatrix} -(S_{11})^{-1} & -(S_{11})^{-1} S_{12} \\ -S_{21} (S_{11})^{-1} & S_{22|1} \end{bmatrix} \quad B = \begin{bmatrix} -(S_{1'1'})^{-1} & -(S_{1'1'})^{-1} S_{1'2'} \\ -S_{2'1'} (S_{1'1'})^{-1} & S_{2'2'|1'} \end{bmatrix}$$

$$X_{1'} = X_1 \cup \{X_u\} ; \quad X_{2'} = X_2 \setminus \{X_u\}$$

$$B(u,u) = -1 / A(u,u)$$

$$B(i,u) = A(i,u) * B(u,u) \quad (i \neq u)$$

$$B(i,j) = A(i,,j) + A(i,u) * B(j,u) \quad (i \neq u, i \neq u)$$

UPDATING THE ELEMENTS OF $[(S_{11})^{-1} (S^2)_{11}]$

$$\left((S_{11'})^{-1} (S^2)_{11'} \right) (\mathbf{u}, \mathbf{j}) = \frac{\left(S^2 \right) (\mathbf{u}, \mathbf{j}) - \sum_{\mathbf{a}: X_{\mathbf{a}} \in X_1} \left[\left((S_{11})^{-1} S_{12} \right) (\mathbf{a}, \mathbf{u}) * \left(S^2 \right) (\mathbf{a}, \mathbf{j}) \right]}{S_{22|1} (\mathbf{u}, \mathbf{u})}$$

$$\left((S_{11'})^{-1} (S^2)_{11'} \right) (\mathbf{i}, \mathbf{j}) = \left((S_{11})^{-1} (S^2)_{11} \right) (\mathbf{i}, \mathbf{j}) - \frac{\left((S_{11})^{-1} S_{12} \right) (\mathbf{i}, \mathbf{u})}{S_{22|1} (\mathbf{u}, \mathbf{u})} *$$

$$* \left\{ \left(S^2 \right) (\mathbf{u}, \mathbf{j}) - \sum_{\mathbf{a}: X_{\mathbf{a}} \in X_1} \left[\left((S_{11})^{-1} S_{12} \right) (\mathbf{a}, \mathbf{u}) * \left(S^2 \right) (\mathbf{a}, \mathbf{j}) \right] \right\} \quad (\mathbf{i} \neq \mathbf{u})$$

NUMBER OF FLOATING POINT OPERATIONS

PER SWEEP:

$$(5/2) K_1^2 + (11/2) K_1 + 2$$

TOTAL:

$$(5/2) K_1^2 + (11/2) K_1 + 2$$

OPERATORS FOR REMOVING VARIABLES

$$\mathbf{A} = \begin{bmatrix} (\mathbf{S}_{1'1'})^{-1} & (\mathbf{S}_{1'1'})^{-1} \mathbf{S}_{1'2'} \\ \mathbf{S}_{2'1'} (\mathbf{S}_{1'1'})^{-1} & -\mathbf{S}_{2'2'|1'} \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} (\mathbf{S}_{11})^{-1} & (\mathbf{S}_{11})^{-1} \mathbf{S}_{12} \\ \mathbf{S}_{21} (\mathbf{S}_{11})^{-1} & -\mathbf{S}_{22|1} \end{bmatrix}$$

$$\mathbf{X}_{1'} = \mathbf{X}_1 \cup \{\mathbf{X}_u\} ; \quad \mathbf{X}_{2'} = \mathbf{X}_2 \setminus \{\mathbf{X}_u\}$$

$$\mathbf{B}(u,u) = -1 / \mathbf{A}(u,u)$$

$$\mathbf{B}(i,u) = \mathbf{A}(i,u) * \mathbf{B}(u,u) \quad (i \neq u)$$

$$\mathbf{B}(i,j) = \mathbf{A}(i,,j) + \mathbf{A}(i,u) * \mathbf{B}(j,u) \quad (i \neq u, i \neq u)$$

$$\begin{aligned} \left((\mathbf{S}_{11})^{-1} (\mathbf{S}^2)_{1\bullet} \right) (i,j) &= \left((\mathbf{S}_{1'1'})^{-1} (\mathbf{S}^2)_{1'\bullet} \right) (i,j) - (\mathbf{S}_{1'1'})^{-1} (i,u) * (\mathbf{S}^2)(j,u) - \\ &- \frac{(\mathbf{S}_{1'1'})^{-1} (i,u)}{(\mathbf{S}_{1'1'})^{-1} (u,u)} * \sum_{\mathbf{a}: \mathbf{X}_a \in \mathbf{X}_1} \left[(\mathbf{S}_{1'1'})^{-1} (\mathbf{a},u) * (\mathbf{S}^2)(\mathbf{a},j) \right] \end{aligned}$$

CONCLUSIONS

ALL-SUBSETS COMPARISONS ARE POSSIBLE:

(a) IF THE NUMBER OF ORIGINAL VARIABLES IS
NOT TOO LARGE (SAY $P \leq 25$)

OR

(b) IF A "FEW" SUBSETS ARE CLEARLY
PREFERABLE TO THE OTHERS

FURTHER RESEARCH

**WHEN ALL-SUBSETS COMPARISONS ARE NOT POSSIBLE
CAN HEURISTIC METHODS BE ALMOST OPTIMAL ???**

SOFTWARE IMPLEMENTATION

LOOK AT:

<http://www.porto.ucp.pt/psilva>