

Selecting variables in multivariate linear models with Subselect

A. PEDRO DUARTE SILVA^{(1)(*)}

JORGE CADIMA⁽²⁾

MANUEL MINHOTO⁽³⁾

JORGE ORESTES CERDEIRA⁽²⁾

(1) FEG/CEGE – UNIV. CATÓLICA PORTUGUESA – C.R. PORTO

(2) I.S. AGRONOMIA – UNIV. TÉCNICA DE LISBOA

(3) DEP. MATEMÁTICA – UNIVERSIDADE DE ÉVORA

(*) Supported by: FEDER / POCI 2010




**Ciência.Inovação
2010**

Programa Operacional Ciência e Inovação 2010
MINISTÉRIO DA CIÊNCIA, TECNOLOGIA E ENSINO SUPERIOR

FCT Fundação para a Ciência e a Tecnologia
MINISTÉRIO DA CIÊNCIA E DA TECNOLOGIA

The *Subselect* Package

OVERVIEW

1. GOALS AND ISSUES OF VARIABLE SELECTION
2. VARIABLE SELECTION IN LINEAR REGRESSION
3. VARIABLE SELECTION IN LINEAR MODELS WITH MULTIPLE RESPONSES
4. VARIABLE SELECTION IN GENERALIZED LINEAR MODELS
5. The  *Subselect* PACKAGE
6. FINAL REMARKS

The *Subselect* Package

GOALS OF VARIABLE SELECTION

- PARCIMONY / INTERPETABILITY
- BETTER PREDICTIONS
- EFFICIENT ESTIMATION / INFERENCE

ISSUES IN VARIABLE SELECTION

- CHOICE OF SELECTION CRITERIA
- EFFICIENCY / EFFECTIVENESS OF SELECTION ALGORITHMS
- EFFECTS OF SELECTION BIAS

The *Subselect* Package

VARIABLE SELECTION IN LINEAR REGRESSION

$$y = Xb + e \quad \text{RSS} = e^T e \quad R^2 = 1 - \frac{\text{RSS}}{(n-1)s_y^2}$$

• RSS (OR R^2) - BASED SELECTION CRITERIA :

ADJUSTED $R^2 \rightarrow R_a^2 = 1 - \frac{\text{RSS}}{(n-k-1)s_y^2}$ (DESCRIP.INDEX)

MALLOWS $C_p \rightarrow C_k = \frac{\text{RSS}}{\hat{\sigma}^2} - (n - 2(k+1))$ (PRED.CRITERION)

AKAIKE AND
SCHWARZ \rightarrow

$$\text{AIC} = n \ln \left(\frac{2\pi \text{RSS}}{n-k-1} \right) + n + 2k$$
$$\text{BIC} = -\frac{1}{2} \left(n \ln \left(\frac{2\pi \text{RSS}}{n-k-1} \right) + n + k \ln(n) \right)$$

(INFORMATION CRITERIA)

The *Subselect* Package

VARIABLE SELECTION IN LINEAR REGRESSION

- PROBLEMS OF SELECTION BIAS:

1. PARAMETER ESTIMATES MAY BE OVERBLOWN
2. CLASSICAL INFERENCE METHODS ARE NOT VALID
3. GOODNESS OF FIT MEASURES ARE TOO OPTIMISTIC

- REMEDIES

1. SHRINKAGE METHODS
2. MULTIPLE-TESTS INFERENCE
3. RESAMPLING AND/OR CROSS-VALIDATION

The *Subselect* Package

VARIABLE SELECTION IN LINEAR REGRESSION

OTHER APPROACHES :

- BAYES FACTORS AND BAYES AVERAGING
- GARROTE, LASSO, ELASTIC NET AND RELATIVES
- GENERAL TO SPECIFIC MODELLING

The *Subselect* Package

LINEAR MODELS WITH MULTIPLE RESPONSES

$$\mathbf{Y} = \mathbf{X} \mathbf{B} + \mathbf{U} \quad r = \min(\dim(\mathbf{X}), \dim(\mathbf{Y}))$$

$$\mathbf{T} = \mathbf{S}_{\mathbf{X}\mathbf{X}} \quad \mathbf{H} = \mathbf{S}_{\mathbf{X}\mathbf{X}} - \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{-1} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \quad \text{ccr}_i^2 = \text{Eigval}_i(\mathbf{T}^{-1}\mathbf{H})$$

Comparison Criteria:

Multivariate Indices

$$\text{ccr}_1^2$$

($\max \text{ccr}_1^2 \Leftrightarrow \max \text{Roy } \lambda_1$)

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

($\max \tau^2 \Leftrightarrow \min \text{Wilks } \Lambda$)

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

($\max \zeta^2 \Leftrightarrow \max \text{Lawley-Hotelling trace}$)

$$\xi^2 = \frac{\sum_{i=1}^r \text{ccr}_i^2}{r}$$

($\max \xi^2 \Leftrightarrow \max \text{Bartlett-Pillai trace}$)

The *Subselect* Package

A LINEAR HYPOTHESIS FRAMEWORK:

$$X = A \Psi + U \quad H_0: C \Psi = 0$$

→ SELECT COLUMNS OF X IN ORDER TO EXPLAIN H1

PARTICULAR CASES:

- 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

$$\Omega = \mathcal{R}(A) \quad \omega = \mathcal{R}(A) \cap \mathcal{N}(C (A^T A)^- A^T) \quad r = \dim(\Omega) - \dim(\omega)$$

$$ccr_i^2 = \text{Eigval}_i(T^{-1}H) \quad T = X' (I - P_\omega) X \quad H = X' (P_\Omega - P_\omega) X$$

The *Subselect* Package

GENERALIZED LINEAR MODELS

$$g(\mathbf{y}) = \mathbf{X} \mathbf{B} + \mathbf{u} = [\mathbf{X}_1 \mathbf{X}_2] [\mathbf{B}_1^T \mathbf{B}_2^T]^T + \mathbf{u}$$

Comparison Criteria:

$$\lambda = 2 \left[\ln L(\hat{\mathbf{B}}) - \ln L(\hat{\mathbf{B}}_{1R}, 0) \right] \quad \text{(Likelihood ratio statistic)}$$

$$\lambda^* = \hat{\mathbf{B}}_2^T (\hat{\mathbf{I}}_{22} - \hat{\mathbf{I}}_{21} \hat{\mathbf{I}}_{11}^{-1} \hat{\mathbf{I}}_{12}) \hat{\mathbf{B}}_2 \quad \hat{\mathbf{I}} = - \frac{\partial^2 \ln L(\mathbf{B})}{\partial(\mathbf{B}) \partial(\mathbf{B}^T)} \Big|_{\mathbf{B}=\hat{\mathbf{B}}} = \begin{bmatrix} \hat{\mathbf{I}}_{11} & \hat{\mathbf{I}}_{12} \\ \hat{\mathbf{I}}_{21} & \hat{\mathbf{I}}_{22} \end{bmatrix} \quad \text{(Wald statistic)}$$

$$\lambda^{**} = 2 \left[\ln L(\hat{\mathbf{B}}) - \ln L(\tilde{\mathbf{B}}_1, 0) \right] \quad \tilde{\mathbf{B}}_1 = \hat{\mathbf{B}}_1 + \hat{\mathbf{I}}_{11}^{-1} \hat{\mathbf{I}}_{12} \hat{\mathbf{B}}_2$$

$$\text{AIC} = \lambda^\# + 2(k + m(k)) \quad \lambda^\# = \lambda, \lambda^* \text{ OR } \lambda^{**}$$

The *Subselect* Package

The Subselect Package

Search routines for (combinatorial) criteria optimization

Exact Algorithm:

- leaps - based on Furnival and Wilson's leaps and bounds algorithm for linear regression
 - viable with up to 30 - 35 original variables

Heuristics:

- anneal - simulated annealing
- genetic - genetic algorithm
- improve - restricted local improvement

The *Subselect* Package

Comparison criteria → functions of $T^{-1} H$ eigenvalues
 T, H symmetric $\text{rank}(H) = r$

Linear Regression ($r=1$) → $\max R^2$

Linear Models with Multiple Responses → \max
“*ccr12*”, “*tau2*”, “*xi2*” or “*zeta2*”

Generalized Linear Models ($r=1$) → \min
Wald statistic (λ^*)

The *Subselect* Package

Search functions arguments :

- Matrices T and H
- Comparison criterion (models with multiple responses)
- Tuning parameters for heuristics
- Maximum time allowed for exact search
- Variables forcibly included or excluded in the selected subsets
- Number of solutions by subset dimensionality
- Numerical tolerance for detecting singular or non-symmetrical matrices

The *Subselect* Package

Auxiliary functions:

- lmHmat** - creates H and T matrices for linear regression/canonical correlation analysis
- ldaHmat** - creates H and T matrices for linear discriminant analysis
- glhHmat** - creates H and T matrices for an analysis based on a linear hypothesis specified by the user

The *Subselect* Package


Expanding a Subselect Analysis

Other issues in variable selection can be tackled by combining subselect with the wider  system:

- Searches across different dimensionalities can be handled by post-processing the results of subselect
- Selection biases can be estimated by resampling or cross-validation procedures that include the basic searches within larger cycles
- Relative importance of individual variables can be accessed by collecting statistics on groups of “good” subsets of different dimensionalities

The *Subselect* Package

Final Remarks

- **Subselect implements effective searches for variable subsets in a wide range of models**
- **Exact searches are viable if the number of original variables is not much larger than 30**
- **Random searches are usually good and much better than greedy (“stepwise”) algorithms**
- **Many subsets lead to similar criterion values. Finding the exact optimum may not be so important**
- **Subselect is most useful when used in combination with other utilities of the  system**

References

Cadima J, Cerdeira JO and Minhoto M (2004). Computational Aspects of Algorithms for Variable Selection in the Context of Principal Components. *Computational Statistics and Data Analysis* **47**: 225-236.

Duarte Silva, A.P. (2001). Efficient Variable Screening for Multivariate Analysis. *Journal of Multivariate Analysis* **76**, 35-62.

Furnival, G.M. & Wilson, R.W. (1974). Regressions by Leaps and Bounds. *Technometrics* **16**: 499-511.

Lawless, J.F and Singhal, K. (1978). Efficient Screening of NonNormal Regression Models. *Biometrics*. **34**: 318-327.

Lawless, J.F and Singhal, K. (1987). ISMOD: An All-Subsets Regression Program for Generalized Linear Models. I. Statistical and Computational Background. *Computing Methods and Programs in Biomedecine* **24**: 125-134