

VarSelLCM

Variable Selection for Model-Based Clustering of Mixed-Type Data Set with Missing Values.

References:

- Marbac, M. and Sedki, M. (2017), Variable selection for model-based clustering using the integrated complete-data likelihood, *Statistics and Computing*, Volume 27, Issue 4, pp 1049–1063.
- Marbac, M., Patin, E. and Sedki, M. (2018), Variable selection for mixed data clustering: Application in human population genomics, *Journal of Classification*, forthcoming.
- Marbac, M. and Sedki, M. (2018), VarSelLCM: an R/C++ package for variable selection in model-based clustering of mixed-data with missing values., *Bioinformatics*, forthcoming.

Introduction

VarSelLCM permits a full model selection (detection of the relevant features for clustering and selection of the number of clusters) in model-based clustering, according to classical information criteria (BIC, MICL or AIC).

Data to analyzed can be composed of continuous, integer and/or categorical features. Moreover, missing values are managed, without any pre-processing, by the model used to cluster with the assumption that values are missing completely at random. Thus, VarSelLCM can also be used for data imputation via mixture models.

An R-Shiny application is implemented to easily interpret the clustering results

Clustering for cytological diagnosis

This section performs the whole analysis of the Golub data set. The first column indicates the type of leukemia. Clustering is made on the other (continuous) features.

```
library(VarSelLCM)
```

Attaching package: 'VarSelLCM'

The following object is masked from 'package:stats':

```
predict
# please install the package multtest to get the data
# source("https://bioconductor.org/biocLite.R")
# biocLite("multtest")
data(golub, package = "multtest")
# one row = one individual
golub <- t(golub)
```

Clustering is done with and without variable selection. Here ICL and MICL criteria are used because the number of observations is less than the number of features (thus, BIC is not relevant).

```
# Please indicate the number of cores you wan to use for parallelization
nb.CPU <- 4
# clustering without variable selection (about than 10/20 sec on 4 CPU)
```

```
res.noselec <- VarSelCluster(golub,
                             gvals = 2,
                             crit.varsel = "ICL",
                             vbleSelec = FALSE,
                             nbcores = nb.CPU)
# clustering with variable selection
res.selec <- VarSelCluster(golub,
                             gvals = 2,
                             vbleSelec= TRUE,
                             crit.varsel = "MICL",
                             nbcores = nb.CPU)
```

Summary of the results: variable selection increases the values of information criteria

```
summary(res.noselec)
```

Model:

Number of components: 2

Model selection has been performed according to the ICL criterion

```
summary(res.selec)
```

Model:

Number of components: 2

Model selection has been performed according to the MICL criterion

Variable selection has been performed, 553 (18.13 %) of the variables are relevant for clustering

Evaluation of the partition accuracy: Adjusted Rand Index (ARI) is computed between the true partition (ztrue) and its estimators. The expectation of ARI is zero if the two partitions are independent. The ARI is equal to one if the partitions are equals. Variable selection permits to improve the ARI. Note that ARI cannot be used for model selection in clustering, because there is no true partition.

```
ARI(golub.cl, fitted(res.noselec))
```

```
[1] 0.4550916
```

```
ARI(golub.cl, fitted(res.selec))
```

```
[1] 0.7927409
```

Mixed-type data analysis

Clustering

This section performs the whole analysis of the *Heart* data set. *Warning the univariate margin distribution are defined by class of the features: numeric columns imply Gaussian distributions, integer columns imply Poisson distribution while factor (or ordered) columns imply multinomial distribution*

```
library(VarSelLCM)
# Data loading:
# x contains the observed variables
# z the known status (i.e. 1: absence and 2: presence of heart disease)
data(heart)
ztrue <- heart[, "Class"]
x <- heart[, -13]
```

```
# Add a missing value artificially (just to show that it works!)
x[1,1] <- NA
```

Clustering is performed with variable selection. Model selection is done with BIC because the number of observations is large (compared to the number of features). The number of components is between 1 and 3. Do not hesitate to use parallelization (here only four cores are used).

```
# Cluster analysis without variable selection
res_without <- VarSelCluster(x, gvals = 1:3, vbleSelec = FALSE, crit.varsel = "BIC")

# Cluster analysis with variable selection (with parallelisation)
res_with <- VarSelCluster(x, gvals = 1:3, nbcores = 4, crit.varsel = "BIC")
```

Comparison of the BIC for both models: variable selection permits to improve the BIC

```
BIC(res_without)
```

```
[1] -6516.216
```

```
BIC(res_with)
```

```
[1] -6509.506
```

Evaluation of the partition accuracy: Adjusted Rand Index (ARI) is computed between the true partition (ztrue) and its estimators. The expectation of ARI is zero if the two partitions are independent. The ARI is equal to one if the partitions are equals. Variable selection permits to improve the ARI. Note that ARI cannot be used for model selection in clustering, because there is no true partition.

```
ARI(ztrue, fitted(res_without))
```

```
[1] 0.2218655
```

```
ARI(ztrue, fitted(res_with))
```

```
[1] 0.2661321
```

To obtain the partition and the probabilities of classification

```
# Estimated partition
fitted(res_with)
```

```
[1] 2 2 1 2 2 2 2 2 2 2 2 2 1 2 1 2 2 2 2 1 2 1 1 1 1 2 1 2 2 2 2 1 2 2
[36] 2 2 2 1 1 1 1 1 1 2 1 2 1 2 2 2 1 1 1 1 2 2 2 2 1 2 1 1 2 2 1 1 1 1
[71] 2 1 1 2 2 2 2 1 1 1 2 2 2 1 2 1 1 2 1 2 1 1 2 2 1 2 2 2 2 1 1 2 1 2 2
[106] 2 2 2 2 1 2 1 1 1 1 1 1 2 2 2 2 2 2 1 1 1 2 1 1 2 2 2 1 2 1 1 1 2 1 2
[141] 2 1 2 2 1 2 1 2 1 1 1 1 1 2 1 1 2 1 2 2 2 2 1 2 1 2 1 1 2 2 2 2 2 1
[176] 2 2 1 2 1 1 2 1 2 1 1 2 2 1 2 1 2 1 1 1 2 1 2 2 2 2 2 2 2 1 1 2 2 1
[211] 2 1 1 2 1 1 1 2 2 1 2 2 1 2 1 2 2 2 1 2 2 2 1 2 2 2 1 2 1 1 2 1 2 2 1
[246] 2 2 1 1 2 2 1 1 2 1 1 2 2 1 1 1 2 1 1 1 1 1 1 2 2
```

```
# Estimated probabilities of classification
head(fitted(res_with, type="probability"))
```

```
      class-1  class-2
[1,] 8.261580e-06 0.9999917
[2,] 3.665324e-01 0.6334676
[3,] 8.244673e-01 0.1755327
[4,] 4.443181e-08 1.0000000
[5,] 3.884759e-03 0.9961152
[6,] 4.521690e-02 0.9547831
```

To get a summary of the selected model.

```
# Summary of the best model
summary(res_with)
```

Model:

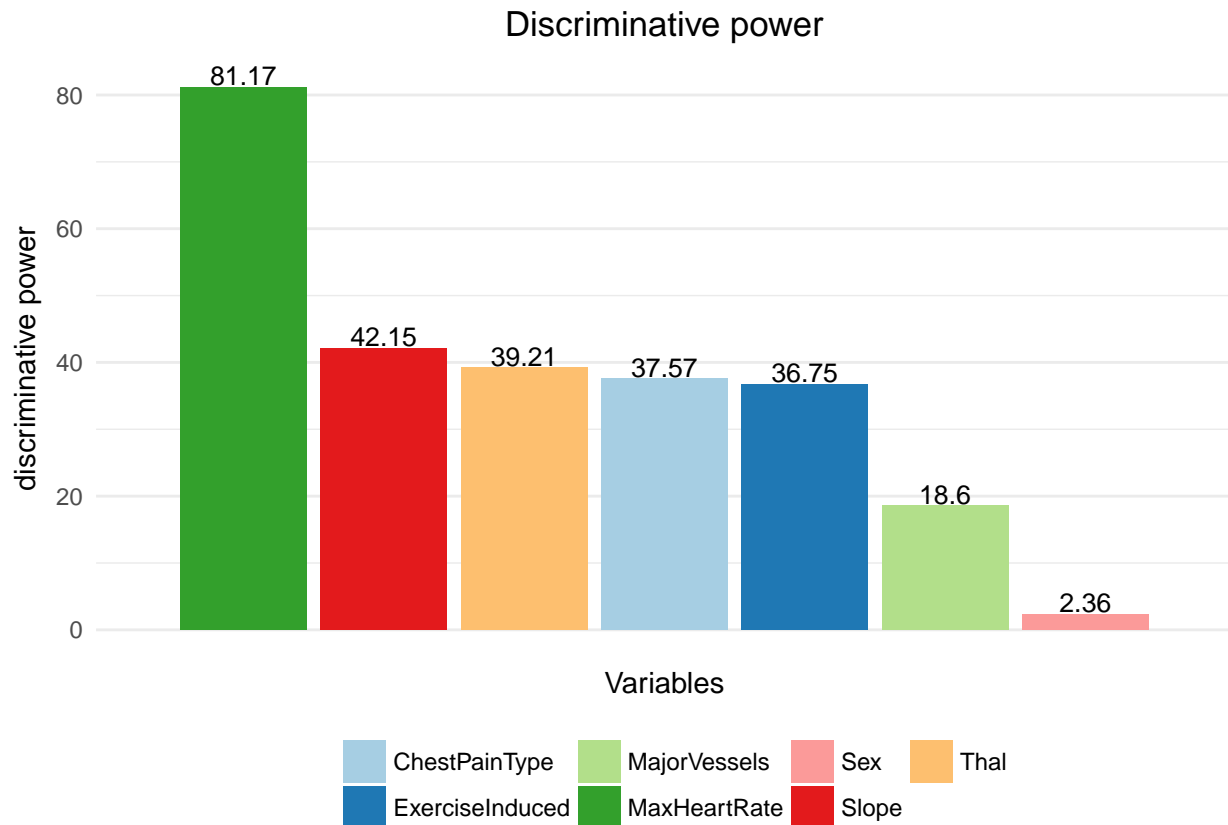
Number of components: 2

Model selection has been performed according to the BIC criterion

Variable selection has been performed, 8 (66.67 %) of the variables are relevant for clustering

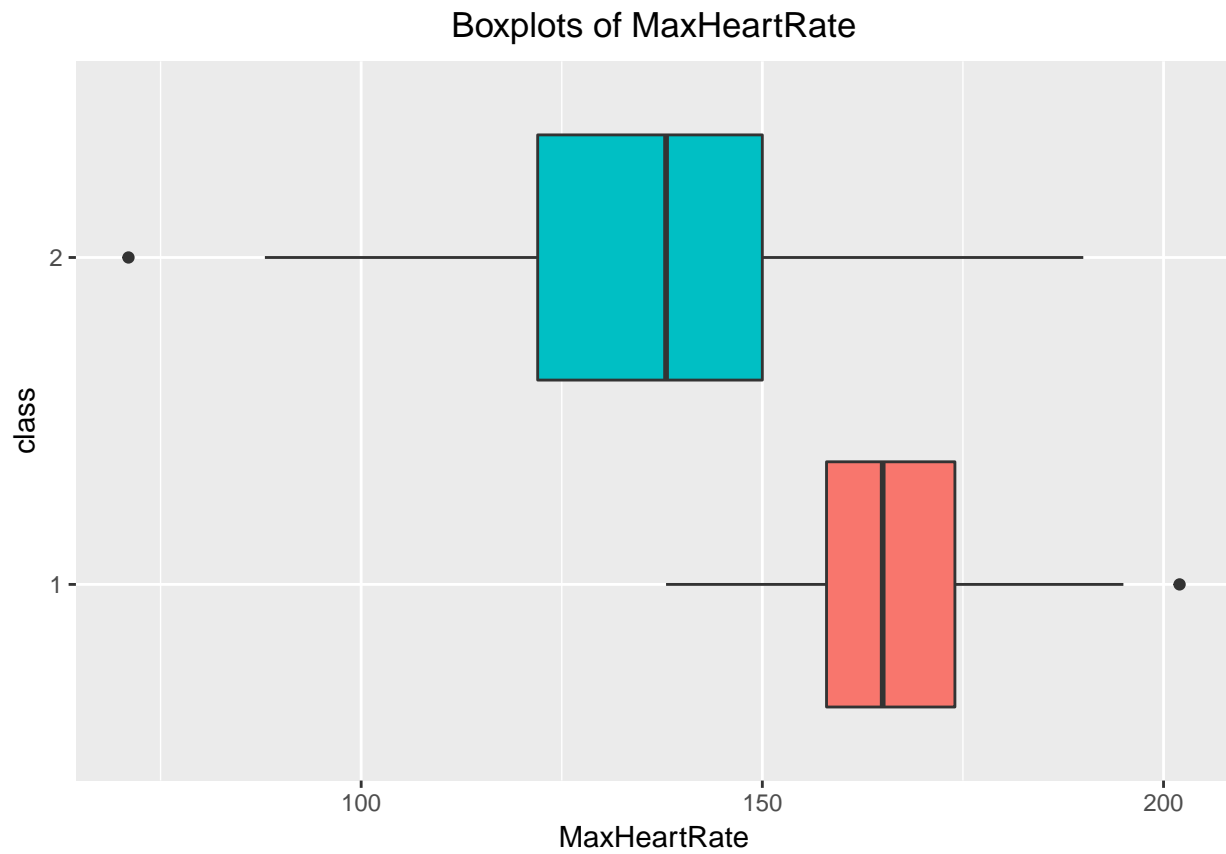
Discriminative power of the variables (here, the most discriminative variable is MaxHeartRate). The greater this index, the more the variable distinguishes the clusters.

```
plot(res_with)
```



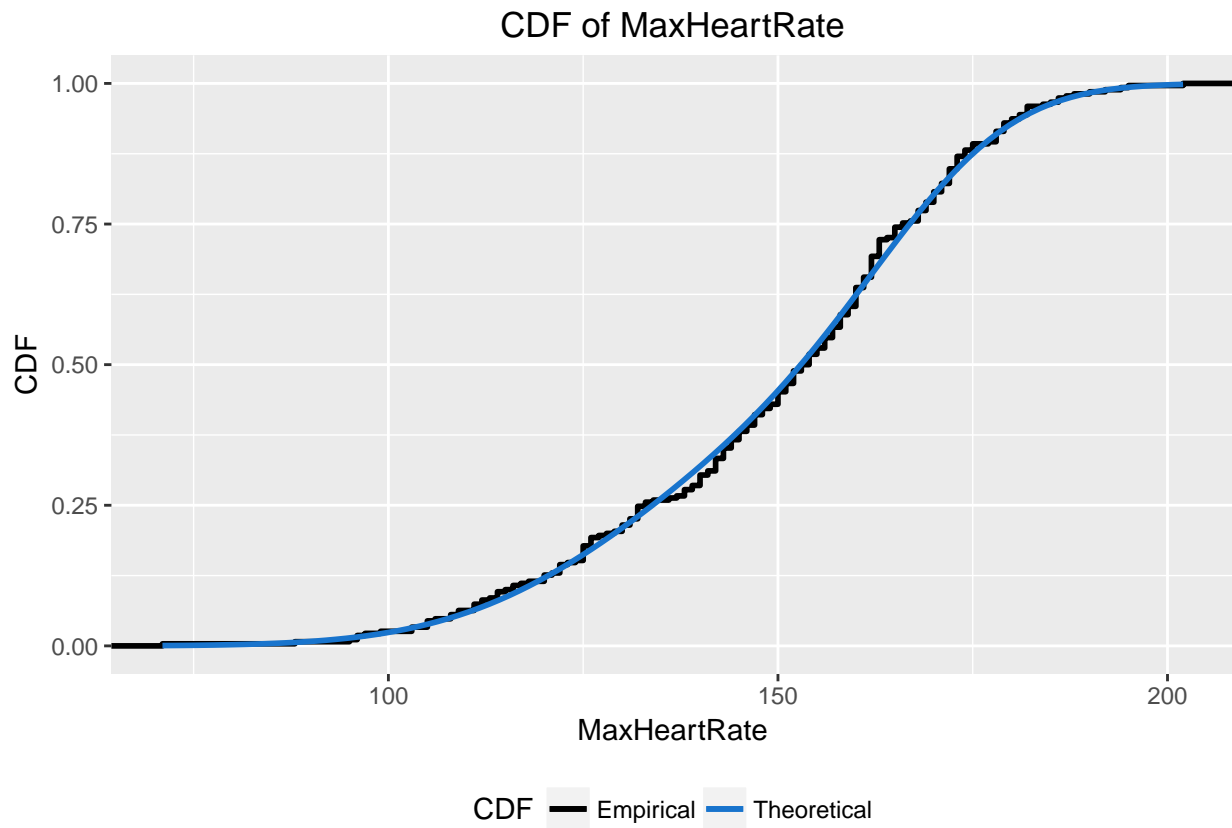
Distribution of the most discriminative variable per clusters

```
# Boxplot for the continuous variable MaxHeartRate
plot(x=res_with, y="MaxHeartRate")
```



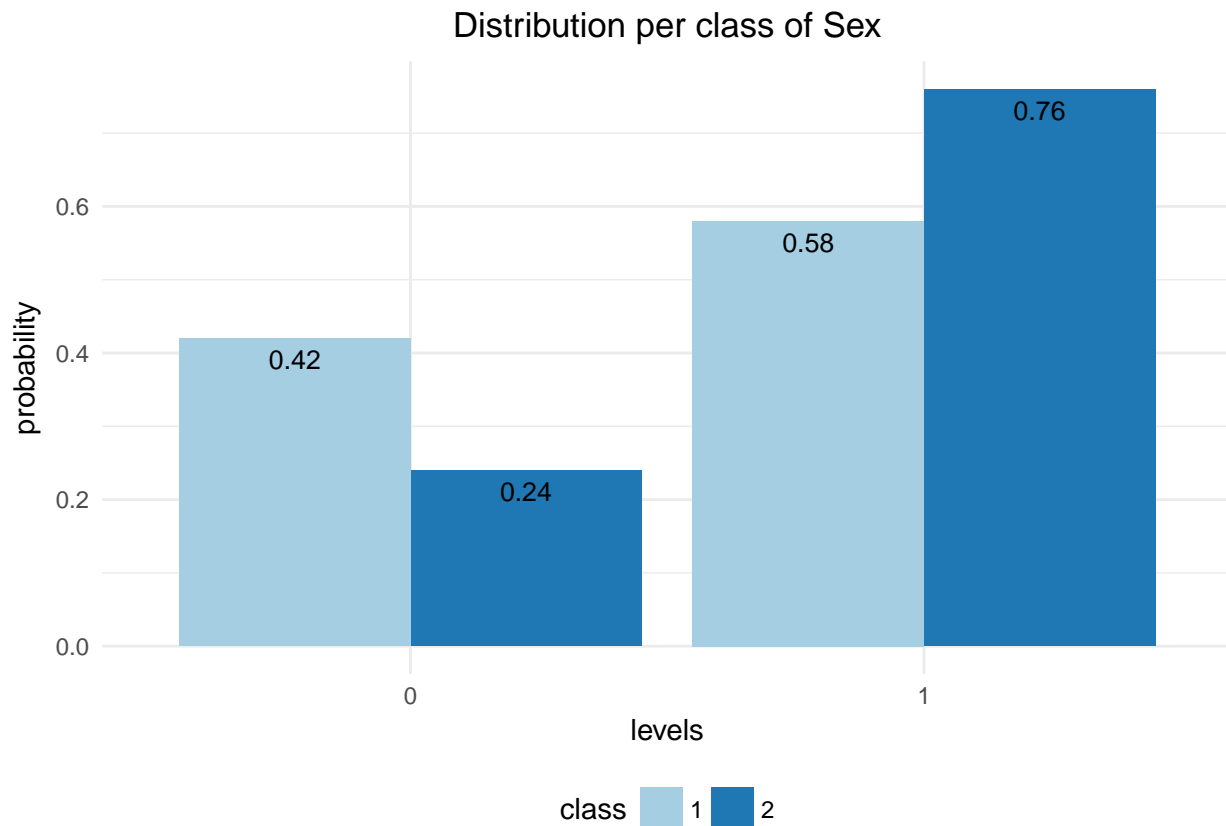
Empirical and theoretical distributions of the most discriminative variable (to check that the distribution is well-fitted)

```
# Empirical and theoretical distributions (to check that the distribution is well-fitted)  
plot(res_with, y="MaxHeartRate", type="cdf")
```



Distribution of a categorical variable per clusters

```
# Summary of categorical variable  
plot(res_with, y="Sex")
```



To have details about the selected model

```
# More detailed output
print(res_with)
```

Data set:

```
Number of individuals: 270
Number of continuous variables: 3
Number of count variables: 1
Percentile of missing values for the integer variables: 0.37
Number of categorical variables: 8
```

Model:

```
Number of components: 2
Model selection has been performed according to the BIC criterion
Variable selection has been performed, 8 ( 66.67 % ) of the variables are relevant for clustering
```

Information Criteria:

```
loglike: -6403.136
AIC: -6441.136
BIC: -6509.506
ICL: -6638.116
```

To print the parameters

```
# Print model parameter
coef(res_with)
```

An object of class "VSLCMparam"
Slot "pi":

```
class-1 class-2
0.4778869 0.5221131
```

```
Slot "paramContinuous":
An object of class "VSLCMparamContinuous"
Slot "pi":
numeric(0)
```

```
Slot "mu":
class-1 class-2
RestBloodPressure 131.3444 131.3444
SerumCholestoral 249.6593 249.6593
MaxHeartRate 165.2587 135.4167
```

```
Slot "sd":
class-1 class-2
RestBloodPressure 17.82850 17.82850
SerumCholestoral 51.59043 51.59043
MaxHeartRate 13.14846 20.98138
```

```
Slot "paramInteger":
An object of class "VSLCMparamInteger"
Slot "pi":
numeric(0)
```

```
Slot "lambda":
class-1 class-2
Age 50.32061 58.11337
```

```
Slot "paramCategorical":
An object of class "VSLCMparamCategorical"
Slot "pi":
numeric(0)
```

```
Slot "alpha":
$Sex
0 1
class-1 0.4166338 0.5833662
class-2 0.2358079 0.7641921
```

```
$ChestPainType
1 2 3 4
class-1 0.05752257 0.28954437 0.4223084 0.2306246
class-2 0.08922357 0.03291641 0.1738645 0.7039955
```

```
$FastingBloodSugar
0 1
class-1 0.8518519 0.1481481
class-2 0.8518519 0.1481481
```

```
$ResElectrocardiographic
0 1 2
```



```
class-1 0.4851852 0.007407407 0.5074074
class-2 0.4851852 0.007407407 0.5074074
```

\$ExerciseInduced

```
      0      1
class-1 0.9128098 0.0871902
class-2 0.4484670 0.5515330
```

\$Slope

```
      1      2      3
class-1 0.7599029 0.1933831 0.04671400
class-2 0.2266441 0.6884267 0.08492922
```

\$MajorVessels

```
      0      1      2      3
class-1 0.7915993 0.1402684 0.05987815 0.008254236
class-2 0.4104430 0.2830467 0.17928537 0.127224890
```

\$Thal

```
      3      6      7
class-1 0.8302562 2.020048e-13 0.1697438
class-2 0.3183111 9.931154e-02 0.5823774
```

Probabilities of classification for new observations

```
# Probabilities of classification for new observations
predict(res_with, newdata = x[1:3,])
```

```
      class-1  class-2
[1,] 8.635654e-06 0.9999914
[2,] 3.768739e-01 0.6231261
[3,] 8.307843e-01 0.1692157
```

The model can be used for imputation (of the clustered data or of a new observation)

```
# Imputation by posterior mean for the first observation
not.imputed <- x[1,]
imputed <- VarSelImputation(res_with, x[1,], method = "sampling")
rbind(not.imputed, imputed)
```

```
Age Sex ChestPainType RestBloodPressure SerumCholestoral
1 NA 1 4 130 322
2 55 1 4 130 322
FastingBloodSugar ResElectrocardiographic MaxHeartRate ExerciseInduced
1 0 2 109 0
2 0 2 109 0
Slope MajorVessels Thal
1 2 3 3
2 2 3 3
```

Shiny application

All the results can be analyzed by the Shiny application...

```
# Start the shiny application
VarSelShiny(res_with)
```

... but this analysis can also be done on \mathbb{R} .