

The MetabolicSurv R package

Predicting Risk Groups for Survival of Cancer Patients Using a Robust Metabolomic Signature

Olajumoke Evangelina Owokotomo¹, Ziv Shkedy¹, Adetayo Kasim²

¹ *Intersersity Institute for Biostatistics and Statistical Bioinformatics, CenStat, Universiteit Hasselt, Belgium*

² *Wolfson Research Institute, Durham University, University Boulevard, Thornaby, Stockton-on-Tees, United Kingdom*



Interuniversity Institute for Biostatistics
and statistical Bioinformatics

INTRODUCTION

- Recent developments in technology has lead to more and more increase in the availability of high throughput medical and biological datasets. In clinical practice, availability of this type of datasets, has lead to identifying signatures which can serves as biomarkers.
- In contrast to the ordinary survival prediction and classification of cancer patients based on prognostic factors measured at baseline alone, recent literature now focuses on integrating high dimensional data in addition to improve prediction and classification into high and low risk-group for survival.
- While methods for the development of high dimensional signatures for survival of cancer patients exist in various context in the literature software to conduct this analysis is not available. The **MetabolicSurv** aimed to close this gap and provides a user friendly data analysis tool, both for modelling and visualization for this type of applications using metabolomics data.

DATA STRUCTURE

The i^{th} observation is denoted by $(T_i, \delta_i, X_{i1}, \dots, X_{ip}, P_1, \dots, P_m)$

$$\begin{bmatrix} T_1 & \delta_1 \\ T_2 & \delta_2 \\ T_3 & \delta_3 \\ \vdots & \vdots \\ T_N & \delta_N \end{bmatrix} \quad \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ X_{31} & X_{32} & \cdots & X_{3p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{N1} & X_{N2} & \cdots & X_{Np} \end{bmatrix} \quad \begin{bmatrix} P_{11} & P_{12} & \cdots & P_{1m} \\ P_{21} & P_{22} & \cdots & P_{2m} \\ P_{31} & P_{32} & \cdots & P_{3m} \\ \vdots & \vdots & \ddots & \vdots \\ P_{N1} & P_{N2} & \cdots & P_{Nm} \end{bmatrix}$$

Survival

Metabolomics profile

Prognostic factors

MetabolicSurv PACKAGE

Category	Functions	Description
Basic	MSpecificCoxPh	Metabolite by metabolite Cox proportional hazard analysis.
	SurvPcaClass	Classifier based on first PCA.
	SurvPlsClass	Classifier based on first PLS.
	Majorityvotes	Classification for Majority Votes.
	Lasoelacox	Wrapper function for glmnet.
	MSData	Generate Artificial Metabolic Survival Data.
Advance	CVLasoelacox	Cross Validations for Lasso Elastic Net and Classification.
	CVSim	Cross-validation for Top K_1, \dots, K_n metabolites.
	CVPcaPls	Cross-validations for PCA and PLS based methods.
	CvMajorityvotes	Cross-validation for majority votes.
	MetFreq	Frequency of Selected Metabolites from MSpecificCoxPh.
	QuantileAnalysis	Sensitivity of the quantile used for classification.
	Icvlasoel	Inner and outer cross-validations for shrinkage methods.
	DistHR	Null distribution of the estimated HR.
	SIMet	Sequentially increase the number of top K metabolites.

METHODOLOGY

Cross-validation

- Assess how results of the statistical analysis can be generalize to an independent data set.
- A small proportion of the dataset is used to validate the result from the bigger proportion that was used to fit the model and label the patients.
- Distribution of the risk-group HR on the training set and testing set is recorded.

Figure 1: The ℓ -fold cross-validation

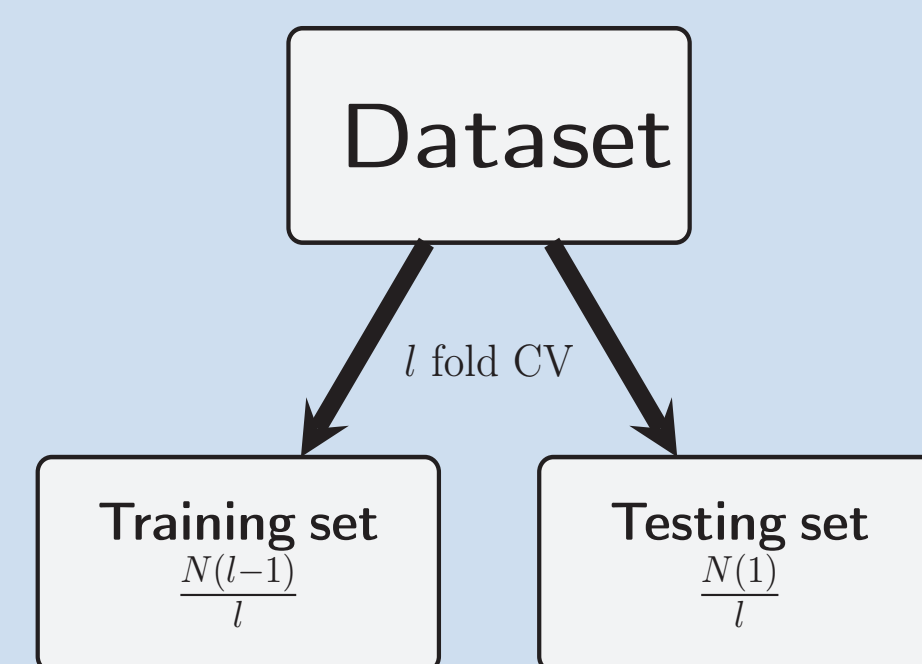
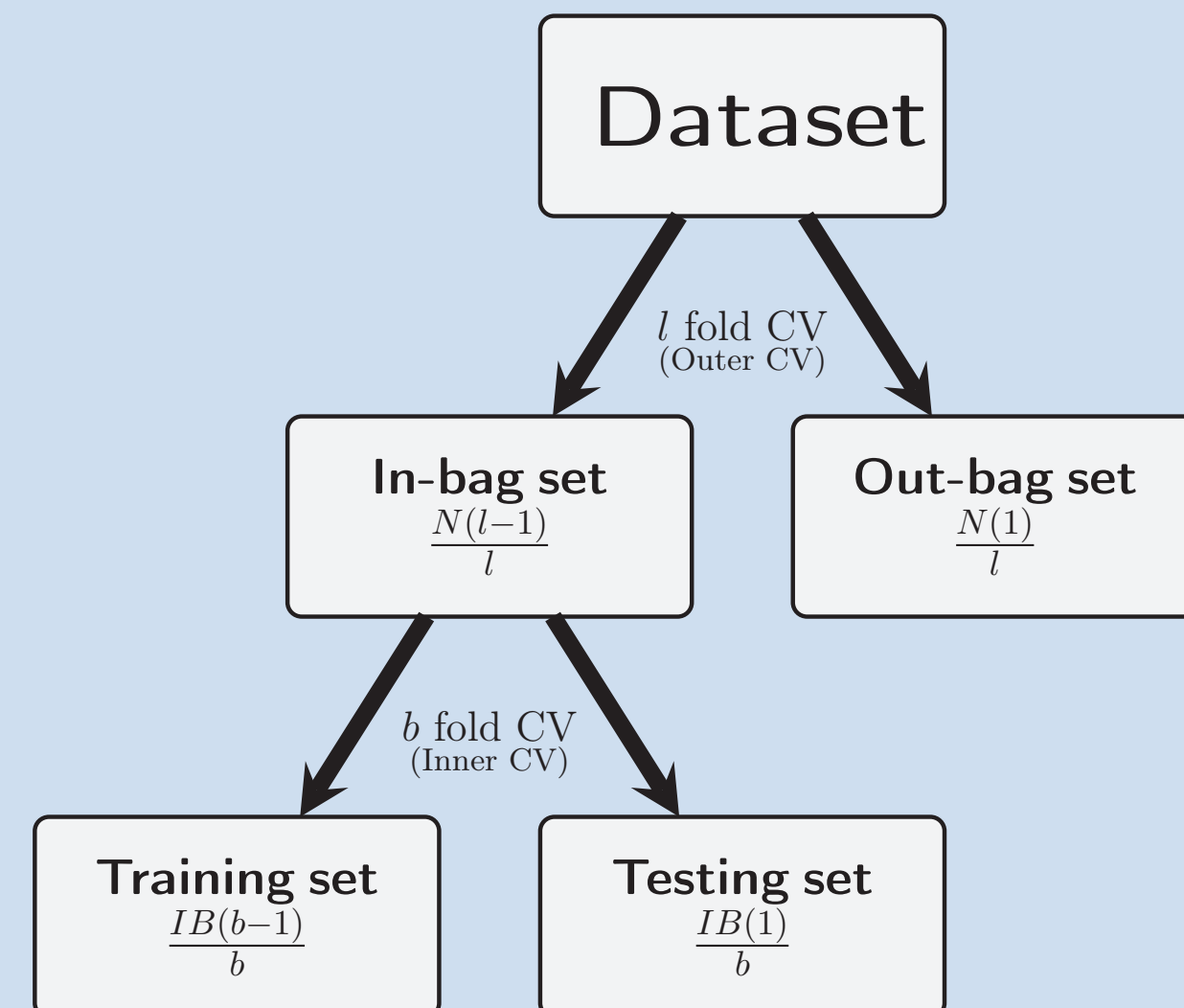


Figure 2: Inner & Outer cross-validation



METHODOLOGY

Estimation of the metabolic risk score in the **MetabolicSurv** package: Data reduction methods, Univariate analysis and Majority voting.

$$\begin{aligned} h(t) &= h_o(t) \exp(\beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \gamma_1 P_1 + \dots + \gamma_m P_m), \\ &= h_o(t) \exp(\mu.), \\ &= \beta \mathbf{U}(\mathbf{X}) + \sum_{m=1}^m \gamma_m P_m \end{aligned}$$

In the **MetabolicSurv** package, $\mathbf{U}(X)$ is a linear combination of the metabolic profile obtained by;

- Scores of the first component of a supervised PCA.
- Scores of the first component of a supervised PLS analysis.
- Linear combination of metabolites with coefficients using LASSO or Elastic Net
- A single metabolite profile.

Data reduction/Variable selection

$$\mathbf{R}(x) = \begin{cases} U(X_r)\beta_{PCA1}, & \text{first PCA is used,} \\ U(X_r)\beta_{PLS1}, & \text{first PLS is used,} \\ U(X_k)\beta_j, & \text{metabolite j is used,} \\ \sum_{j=1}^J X_j\beta_j^s, & \text{shrinkage methods are used and } J \ll p. \end{cases}$$

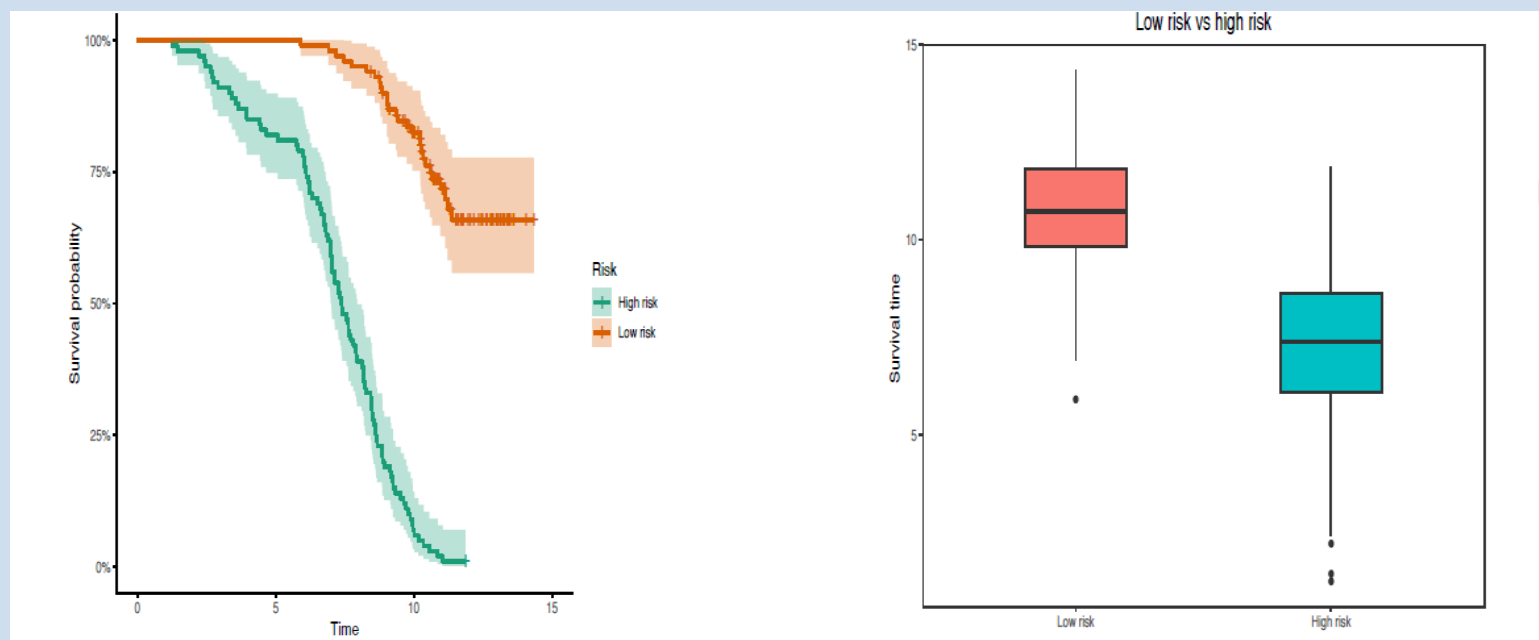
Let τ be a pre defined threshold, for example the median of $R_i(x)$, then $\mathbf{R}_i(x) > \tau$ patient i is labelled as having high risk for mortality, otherwise having low risk for mortality. Let \mathbf{Z}_i be the indicator variable which is given by;

$$\mathbf{Z}_i = \begin{cases} 1, & \text{low risk,} \\ 0, & \text{high risk.} \end{cases} \implies \implies \implies \implies h(t) = \sum_{m=1}^m \gamma_m P_m + \lambda Z_i.$$

SELECTED OUTPUT OF THE PACKAGE

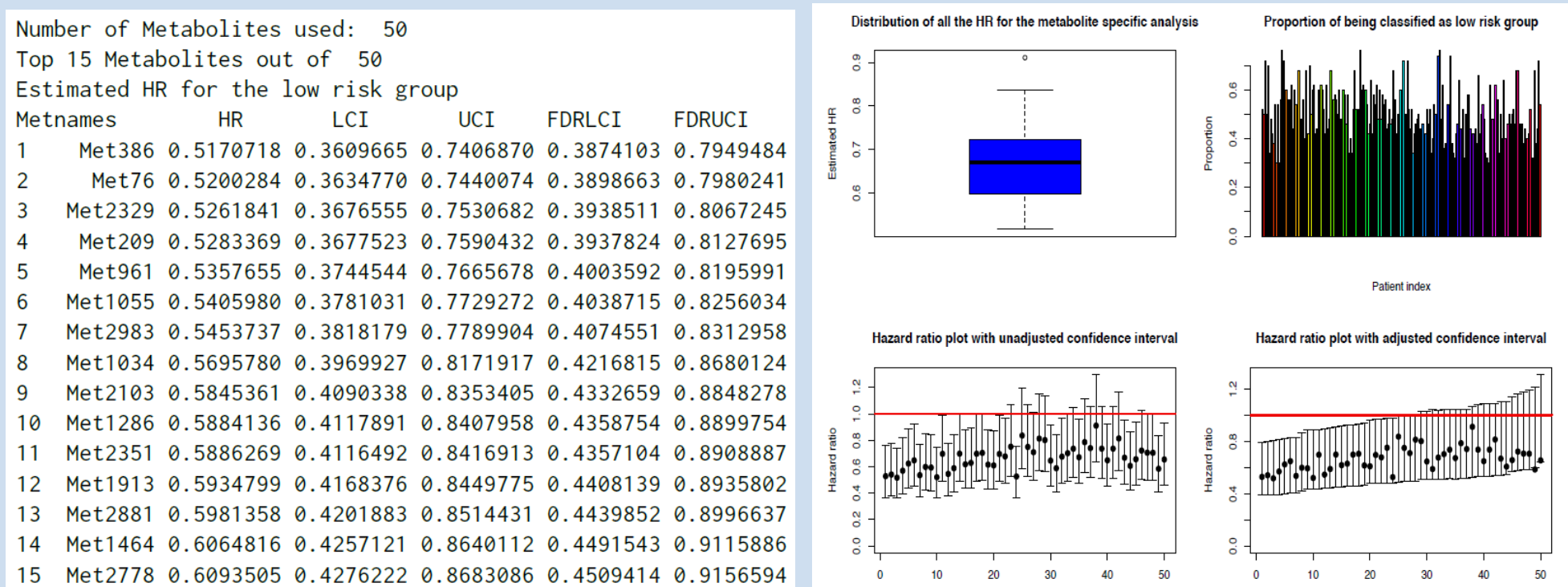
PCA and PLS based method

```
> Result <- SurvPcaClass(Survival = Survdata, Mdata = Metdata, Censor =  
  Censordata, Reduce = TRUE, Select = 150, Prognostic = Progdata, Plots =  
  TRUE, Quantile = 0.5)  
> Result$KMplot; Result$SurvBPlot
```



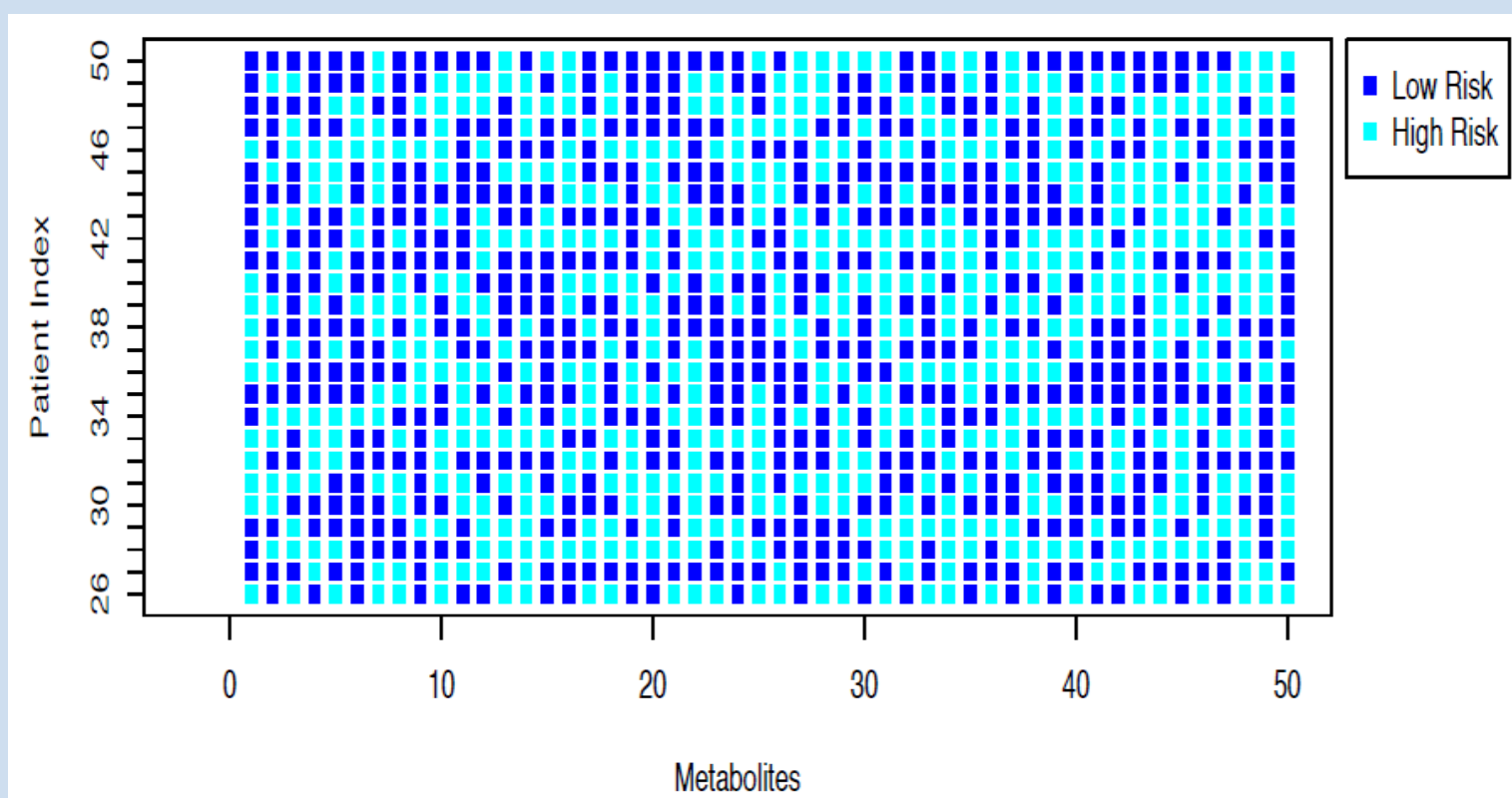
Feature Specific CoxPh Modelling

```
> Result1 <- MSpecificCoxPh(Survival = Survdata, Mdata = Metdata,  
  Censor = Censordata, Reduce = TRUE, Select = 50, Prognostic = Progdata,  
  Quantile = 0.5)  
> summary(Result1)
```



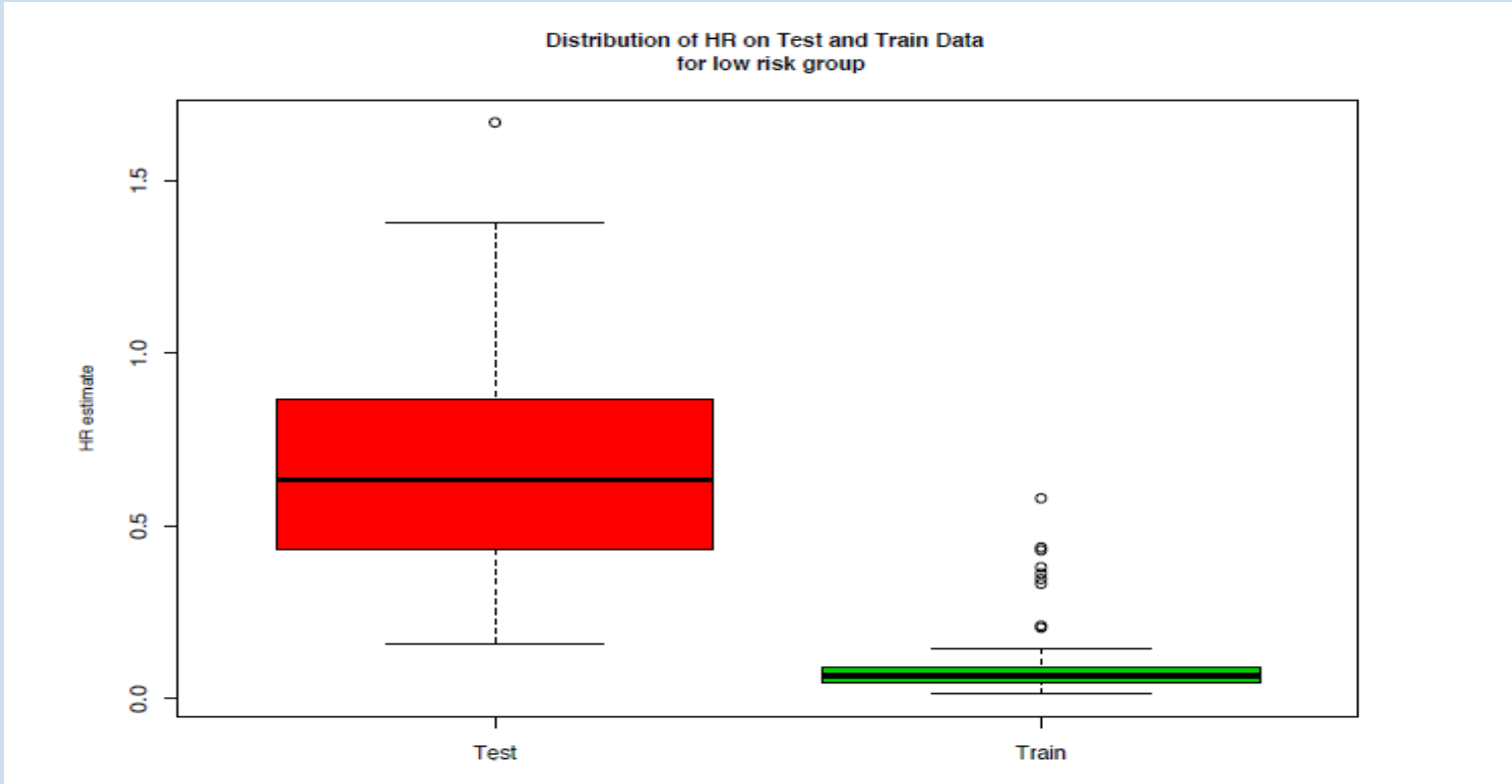
Majority Votes

```
> Result2=Majorityvotes(Result1,Progdata,Survdata, Censordata,J=2)
```



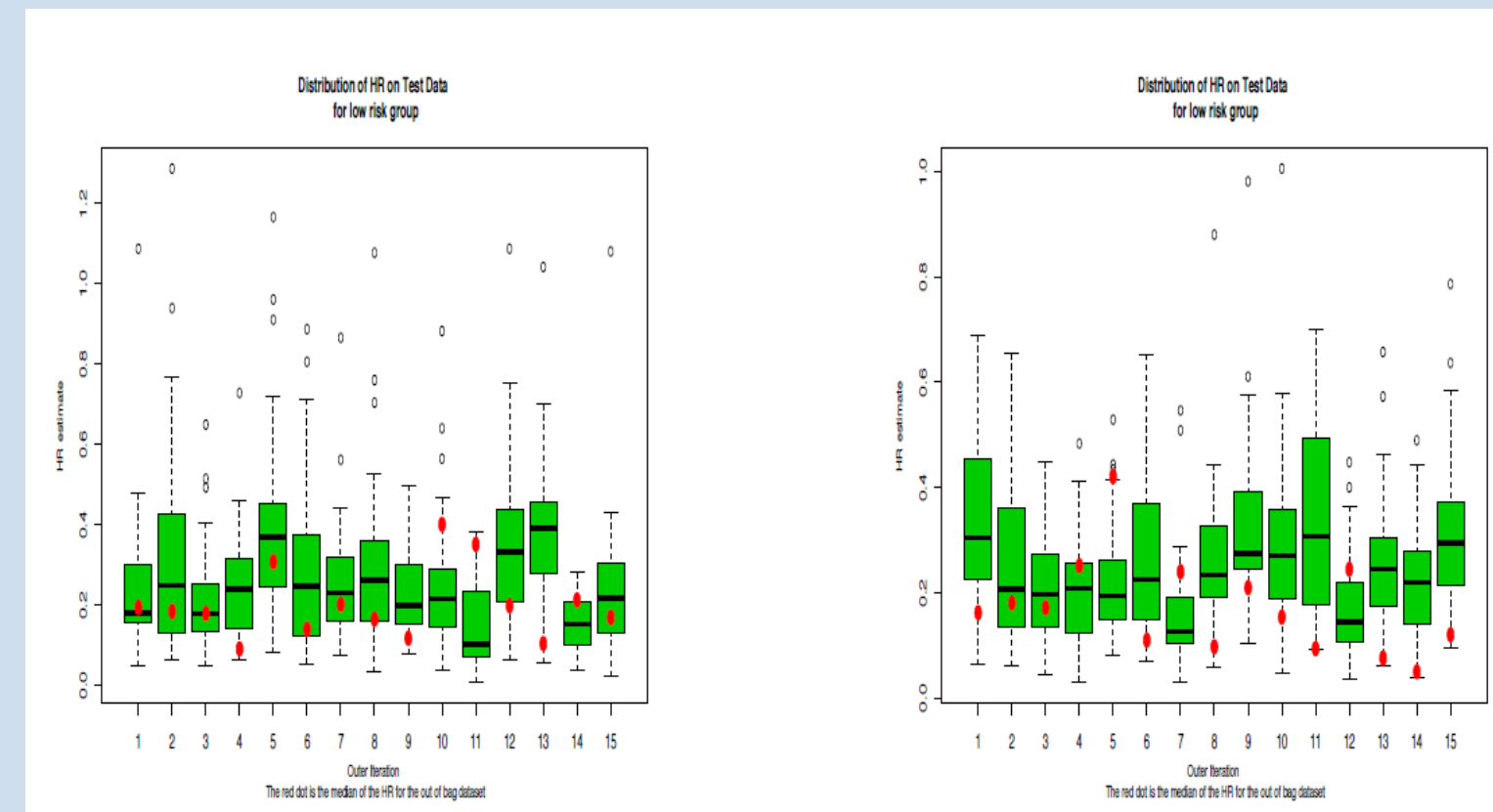
Cross validations for Lasso and ElasticNet

```
> Result3 = CVLasoelacox(Survdata,Censordata, Mdata = Metdata,Progdata,  
  Quantile = 0.5, Metlist = NULL,Standardize = TRUE, Reduce=TRUE,  
  Select=150, Alpha = 1,Fold = 4,Ncv = 100,nlambd = 100)  
> summary(Result1)
```



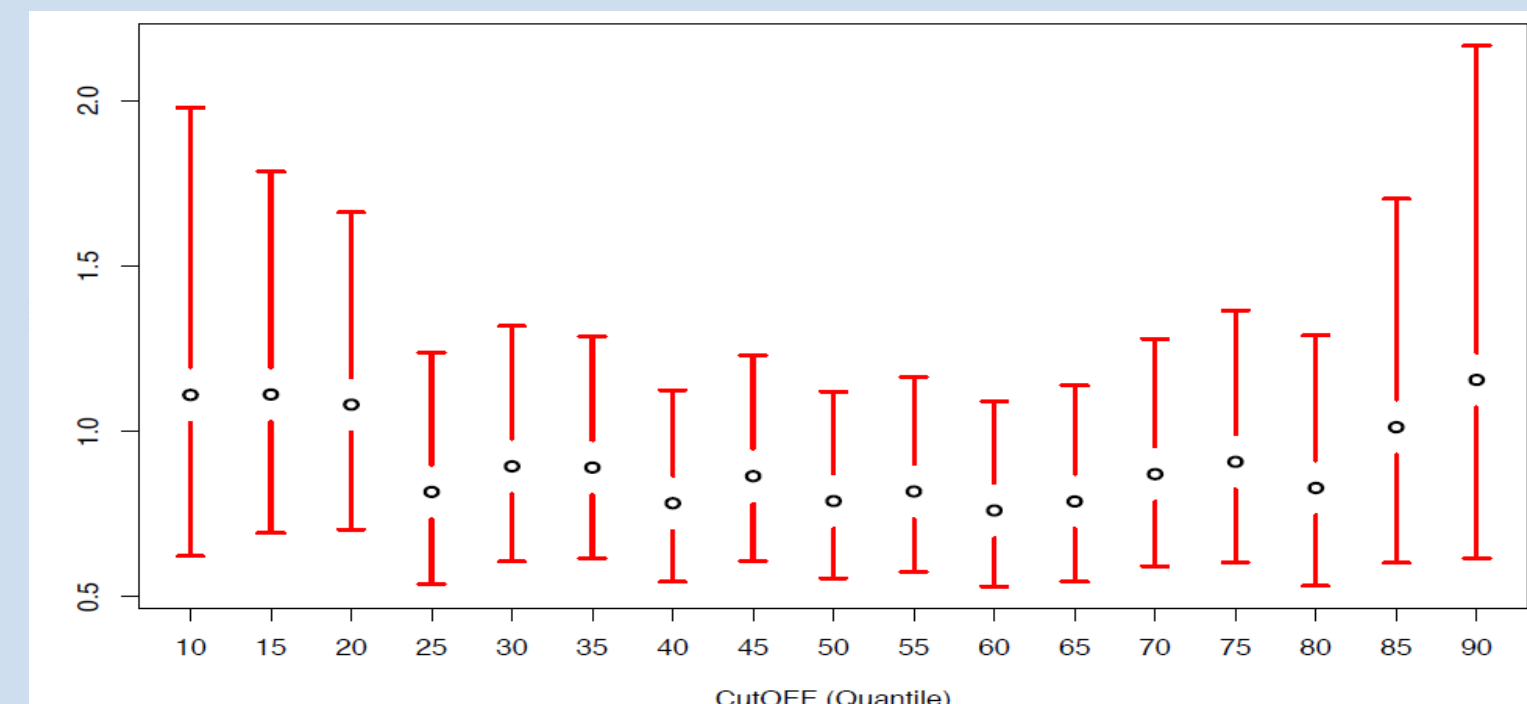
Inner & Outer Cross validations for Lasso

```
> Result4 = Icvlasoel(Survdata, Censordata, Progdata,Metdata, Fold =  
  3,Ncv = 15, Nicv = 30, Alpha = 1, TopK = Mets, Weights = FALSE)  
> Result5 = Icvlasoel(Survdata, Censordata, Progdata, Metdata, Fold =  
  3,Ncv = 15, Nicv = 30, Alpha = 1, TopK = Mets, Weights = TRUE)  
> plot(Result4, type = 1); plot(Result5, type = 1)
```



Sensitivity Analysis of the Cutoff value

```
> Result6 <- QuantileAnalysis(Survdata,Metdata, Censordata,Reduce=TRUE,  
  Select=150, Prognostic=Progdata,Plots = TRUE,DimMethod="PCA",Alpha=1)
```



DOWNLOAD

- MetabolicSurv: <https://cran.r-project.org/web/packages/MetabolicSurv/index.html>.
- More details available in <https://github.com/OlajumokeEvangelina/MetabolicSurv>.

REFERENCES

- [1] O. E. Owokotomo and Z. Shkedy, *MetabolicSurv: A Biomarker Validation Approach for Classification and Predicting Survival Using Metabolomics Signature*, 2019. R package version 1.0.0.