

# A demonstration of the RaSEn package

Ye Tian and Yang Feng

2022-03-01

We provide a detailed demo of the usage for the **RaSEn** package. This package implements the random subspace ensemble classification (RaSE) method (Tian and Feng (2021b)), the variable screening approach via RaSE (Tian and Feng (2021a)) and the super RaSE method (Zhu and Feng (2021)).

- Random Subspace Ensemble Classification
  - Introduction
  - Installation
  - How to fit a RaSE classifier for prediction
  - How to use RaSE for feature ranking
  - Super RaSE
- Variable Screening via Random Subspace Ensembles
  - Introduction
  - Variable screening via RaSE

## Random Subspace Ensemble Classification

### Introduction

Suppose we have training data  $\{\mathbf{x}_i, y_i\}_{i=1}^n \in \{\mathbb{R}^p, \{0, 1\}\}$ , where each  $\mathbf{x}_i$  is a  $1 \times p$  vector.

Based on training data, RaSE algorithm aims to generate  $B_1$  weak learners  $\{C_n^{S_j}\}_{j=1}^{B_1}$ , each of which is constructed in a feature subspace  $S_j \subseteq \{1, \dots, p\}$  instead using all  $p$  features. To obtain each weak learner,  $B_2$  candidates  $\{C_n^{S_{jk}}\}_{k=1}^{B_2}$  are trained based in subspaces  $\{S_{jk}\}_{k=1}^{B_2}$ , respectively. To choose the optimal one among these  $B_2$  candidates, some criteria need to be applied, including minimizing ratio information criterion (RIC, Tian and Feng (2021b)), minimizing extended Bayes information criterion (eBIC, Chen and Chen (2008), Chen and Chen (2012)), minimizing the training error, minimizing the validation error (if validation data is available), minimizing the cross-validation error, minimizing leave-one-out error etc. And the type of weak learner can be quite flexible.

To better adapt RaSE into the sparse setting, we can update the distribution of random feature subspaces according to the selected percentage of features in  $B_1$  subspaces in each round. This can be seen as an adaptive strategy to increase the possibility to cover the signals that contribute to our model, which can improve the performance of RaSE classifiers in sparse settings.

The selected percentage of each of  $p$  features in  $B_1$  subspaces can be used for feature ranking as well. And we could plot the selected percentage to intuitively rank the importance of each feature in a RaSE model.

## Installation

RaSEn can be installed from CRAN.

```
install.packages("RaSEn", repos = "http://cran.us.r-project.org")
```

Then we can load the package:

```
library(RaSEn)
```

## How to Fit a RaSE Classifier for Prediction

We will show in this section how to fit RaSE classifiers based on different types of base classifiers. First we generate the data from a binary guassian mixture model (model 1 in Tian and Feng (2021b))

$$\mathbf{x} \sim (1 - y)N(\boldsymbol{\mu}^{(0)}, \Sigma) + yN(\boldsymbol{\mu}^{(1)}, \Sigma),$$

where  $\boldsymbol{\mu}^{(0)}, \boldsymbol{\mu}^{(1)}$  are both  $1 \times p$  vectors,  $\Sigma$  is a  $p \times p$  symmetric positive definite matrix. Here  $y$  follows a bernoulli distribution:

$$y \sim \text{Bernoulli}(\pi_1),$$

where  $\pi_1 \in (0, 1)$  and we denote  $\pi_0 = 1 - \pi_1$ .

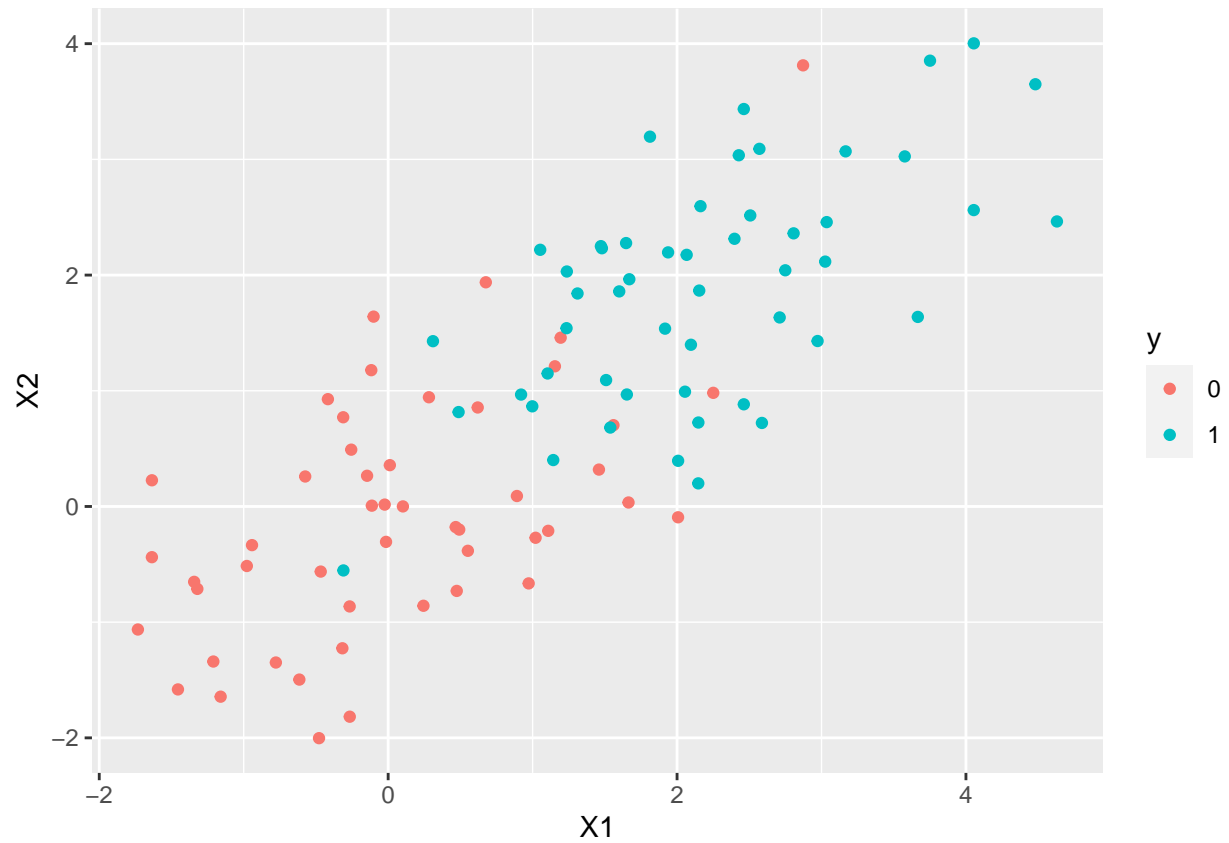
Here we follow from the setting of Mai, Zou, and Yuan (2012), letting  $\Sigma = (0.5^{|i-j|})_{p \times p}$ ,  $\boldsymbol{\mu}^{(0)} = \mathbf{0}_{p \times 1}$ ,  $\boldsymbol{\mu}^{(1)} = \Sigma^{-1} \times 0.556(3, 1.5, 0, 0, 2, \mathbf{0}_{1 \times (p-5)})^T$ . Let  $n = 100, p = 50$ . According to the definition of minimal discriminative set in Tian and Feng (2021b), here the minimal discriminative set  $S^* = \{1, 2, 5\}$ , which contribute to the classification.

Apply function `RaModel` to generate training data and test data of size 100 with dimension 50.

```
set.seed(0, kind = "L'Ecuyer-CMRG")
train.data <- RaModel("classification", 1, n = 100, p = 50)
test.data <- RaModel("classification", 1, n = 100, p = 50)
xtrain <- train.data$x
ytrain <- train.data$y
xtest <- test.data$x
ytest <- test.data$y
```

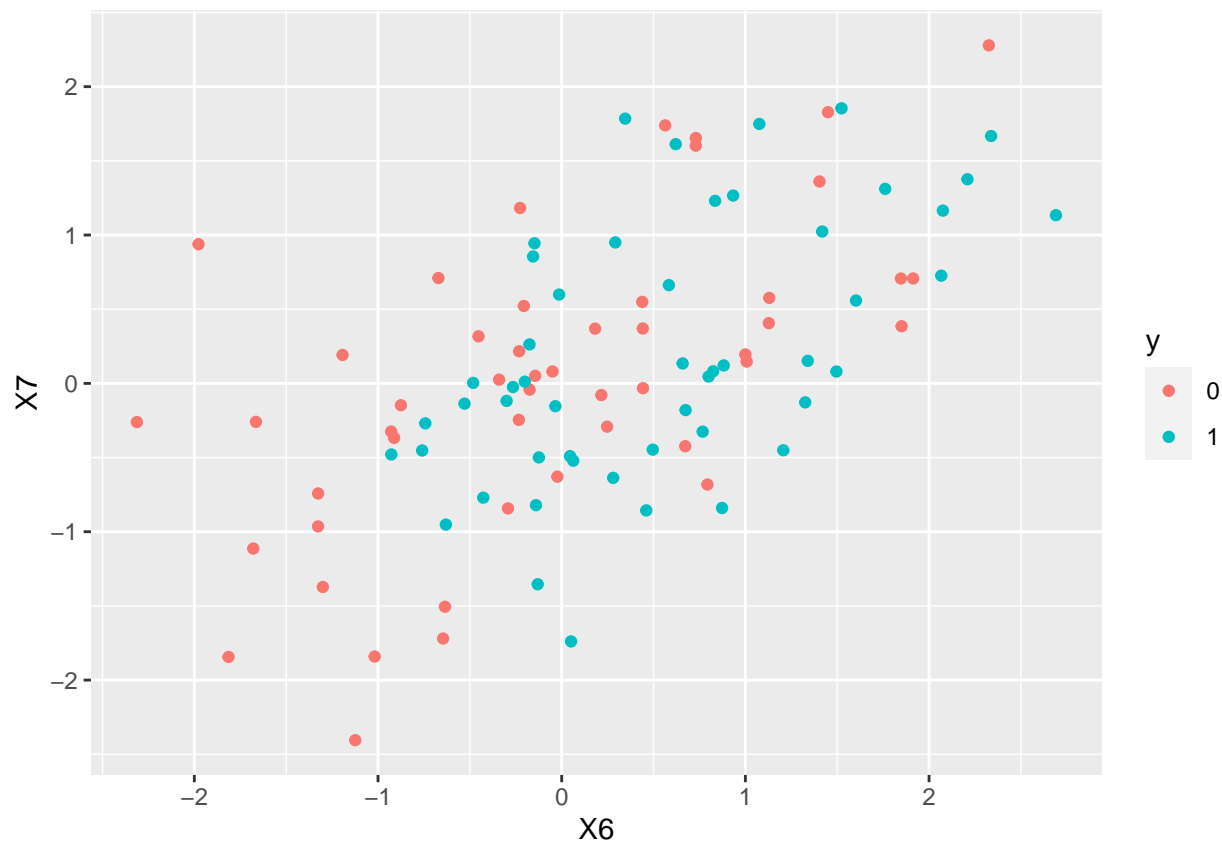
We can visualize the first two dimensions or feature 1 and 2 as belows:

```
library(ggplot2)
ggplot(data = data.frame(xtrain, y = factor(ytrain)), mapping = aes(x = X1,
  y = X2, color = y)) + geom_point()
```



Similarly, we can also visualize the feature 6 and 7:

```
ggplot(data = data.frame(xtrain, y = factor(ytrain)), mapping = aes(x = X6, y = X7,  
  color = y)) + geom_point()
```



It's obvious to see that in dimension 1 and 2 the data from two classes are more linearly separate than in dimension 6 and 7. Then we call **RaSE** function to fit the RaSE classifier with LDA, QDA and logistic regression base classifiers with criterion of minimizing RIC and RaSE classifier with knn base classifier with criterion of minimizing leave-one-out error. To use different types of base classifier, we set **base** as "lda," "qda," "knn" and "logistic," respectively. **B1** is set to be 100 to generate 100 weak learners and **B2** is set to be 100 as well to generate 100 subspace candidates for each weak learner. Without using iterations, we set **iteration** as 0. **criterion** is set to be "ric" for RaSE classifier with LDA, QDA and logistic regression while it is "loo" for RaSE classifier with knn base classifier. To speed up the computation, we apply parallel computing with 2 cores by setting **cores** = 2.

```
fit.lda <- RaSE(xtrain, ytrain, B1 = 100, B2 = 50, iteration = 0, base = "lda",
  cores = 2, criterion = "ric")
fit.qda <- RaSE(xtrain, ytrain, B1 = 100, B2 = 50, iteration = 0, base = "qda",
  cores = 2, criterion = "ric")
fit.knn <- RaSE(xtrain, ytrain, B1 = 100, B2 = 50, iteration = 0, base = "knn",
  cores = 2, criterion = "loo")
fit.logistic <- RaSE(xtrain, ytrain, B1 = 100, B2 = 50, iteration = 0,
  base = "logistic", cores = 2, criterion = "ric")
```

We can print the summarized results of RaSE model by calling **print** function. For instance, we print the RaSE model with LDA base classifier:

```
print(fit.lda)
```

```
## Marginal probabilities:
## class 0 class 1
```

```
##      0.49      0.51
## Type of base classifiers: lda
## Criterion: ric
## B1: 100
## B2: 50
## D: 10
## Cutoff: 0.5321693
## Selected percentage of each feature appearing in B1 subspaces:
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26
## 95 20 10 15 53 30 11  6  9 12 21  7 18 14  7 13 12 13  9 16 21  6  9 10  8  9
## 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50
## 11  8  6 13 19 16 12 13  9 13 11 13 14 11 15 13 11  8  7 10 13 15 12 15
```

To evaluate the performance of four different models, we calculate the test error on test data:

```
er.lda <- mean(predict(fit.lda, xtest) != ytest)
er.qda <- mean(predict(fit.qda, xtest) != ytest)
er.knn <- mean(predict(fit.knn, xtest) != ytest)
er.logistic <- mean(predict(fit.logistic, xtest) != ytest)
cat("LDA:", er.lda, "QDA:", er.qda, "knn:", er.knn, "logistic:", er.logistic)
```

```
## LDA: 0.1 QDA: 0.12 knn: 0.14 logistic: 0.11
```

And the output of `Rase` function is an object belonging to S3 class “RaSE.” It contains:

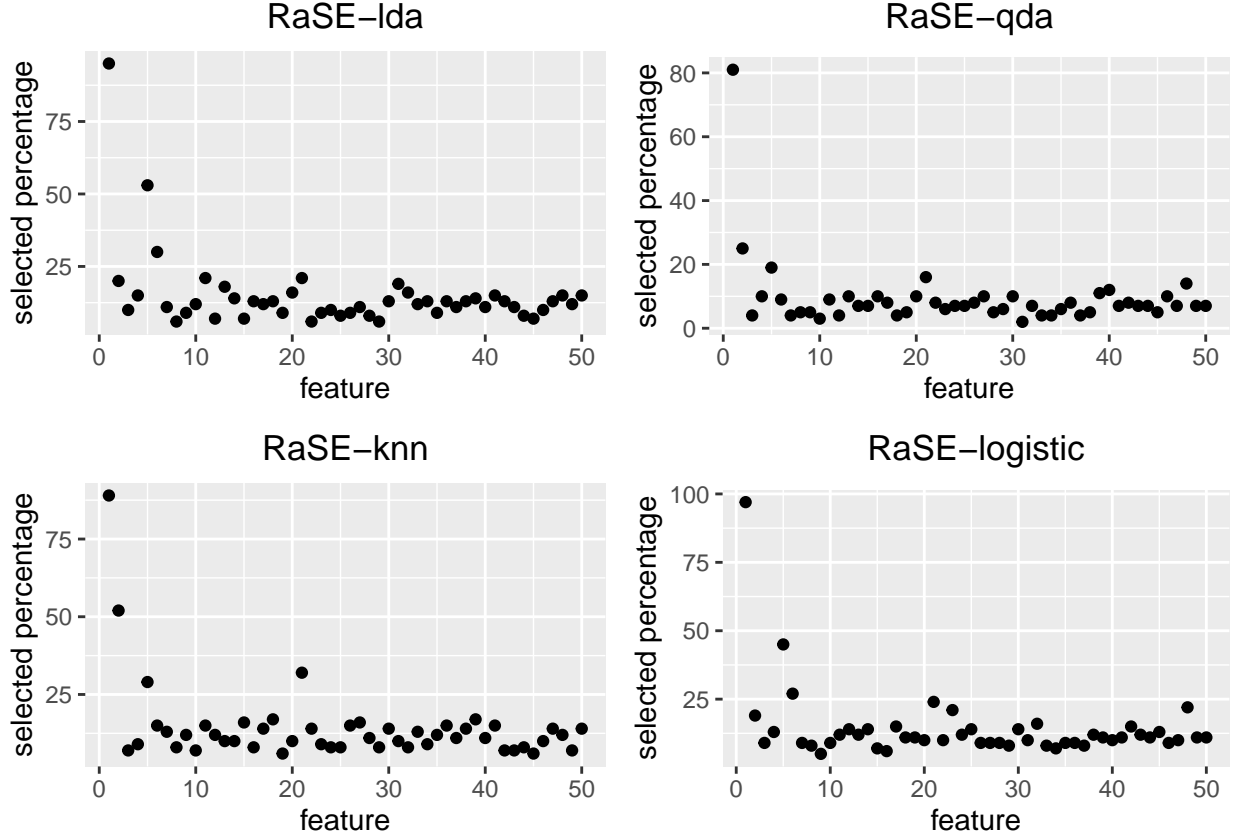
- marginal: the marginal probability for each class.
- base: the type of base classifier.
- criterion: the criterion to choose the best subspace for each weak learner.
- B1: the number of weak learners.
- B2: the number of subspace candidates generated for each weak learner.
- D: the maximal subspace size when generating random subspaces.
- iteration: the number of iterations.
- fit.list: a list of B1 fitted base classifiers.
- cutoff: the empirically optimal threshold.
- subspace: a list of subspaces corresponding to B1 weak learners.
- ranking: the selected percentage of each feature in B1 subspaces.
- scale: a list of scaling parameters, including the scaling center and the scale parameter for each feature.

## How to Use RaSE for Feature Ranking

The selected percentage of features in  $B_1$  subspaces for four RaSE classifiers are contained in the output, which can be used for feature ranking. We can plot them by using `RaPlot` function:

```
library(gridExtra)
plot_lda <- RaPlot(fit.lda)
plot_qda <- RaPlot(fit.qda)
plot_knn <- RaPlot(fit.knn)
plot_logistic <- RaPlot(fit.logistic)

grid.arrange(plot_lda, plot_qda, plot_knn, plot_logistic, ncol = 2)
```



From four figures, it can be noticed that feature 1, 2 and 5 obtain high selected percentage among all  $p = 50$  features under LDA, QDA and  $k$ NN models, implying their importance in classification model. We can set a positive iteration number to increase the selected percentage of three signals among  $B_1$  subspaces, which may improve the performance.

## Super RaSE

In RaSE, we consider only a single type of base classifiers (e.g. LDA, QDA,  $k$ NN, etc.). Zhu and Feng (2021) extends the idea of RaSE by combining classifiers of different types. For each of the  $B_1 B_2$  weak learners, the base classifier type is sampled randomly from some given types with corresponding probabilities. For iterative super RaSE, not only the feature sampling probability will be updated based on the feature selected frequencies, but the classifier type sampling probability will be updated according to the type selected frequencies as well. Note that here the feature sampling probability is updated on the basis of feature frequencies in the last iteration for corresponding base classifier type. The user can also fix the classifier type sampling probability. It can be controlled by component **base.update** in parameter **super**. The super RaSE will be fitted when the parameter **base** is a string vector of base classifier types or a numeric probability vector with classifier type names. In the first case, the base classifier type will be sampled uniformly, while in the second case, it will be sampled according to the provided probability.

The following example shows how to fit a super RaSE classifier which mixes  $k$ NN, LDA and logistic regression.

```
fit.super <- RaSE(xtrain = xtrain, ytrain = ytrain, B1 = 100, B2 = 50,
  base = c("knn", "lda", "logistic"), super = list(type = "separate",
    base.update = T), criterion = "cv", cv = 5, iteration = 0, cores = 2)
ypred <- predict(fit.super, xtest)
mean(ypred != ytest)
```

```
## [1] 0.09
```

We can look at the sampling percentage of each feature from each classifier type and the base classifier type selected percentages.

```
fit.super$ranking.feature
```

```
##           1           2           3           4           5           6           7
## knn      77.14286 51.42857 0.00000 11.42857 34.28571 8.571429 17.142857
## lda      88.88889 38.88889 11.11111 16.66667 22.22222 16.66667 5.555556
## logistic 82.75862 34.48276 10.34483 17.24138 44.82759 20.689655 6.896552
##           8           9          10          11          12          13          14
## knn      8.571429 5.714286 2.857143 11.42857 5.714286 11.428571 14.285714
## lda      8.333333 13.88889 8.333333 25.00000 5.555556 13.88889 8.333333
## logistic 13.793103 17.241379 10.344828 10.34483 13.793103 3.448276 10.344828
##          15          16          17          18          19          20          21
## knn      8.571429 5.714286 11.428571 11.428571 5.714286 11.428571 20.00000
## lda      8.333333 8.333333 8.333333 16.66667 11.11111 5.555556 19.44444
## logistic 10.344828 3.448276 3.448276 6.896552 13.793103 6.896552 20.68966
##          22          23          24          25          26          27          28
## knn      17.142857 8.571429 11.428571 8.571429 14.285714 17.142857 5.714286
## lda      8.333333 27.777778 11.111111 13.88889 2.777778 8.333333 8.333333
## logistic 3.448276 31.034483 3.448276 13.793103 3.448276 10.344828 10.344828
##          29          30          31          32          33          34          35
## knn      8.571429 2.857143 17.142857 20.000000 17.14286 2.857143 11.42857
## lda      13.88889 8.333333 11.11111 2.777778 11.11111 13.88889 0.00000
## logistic 13.793103 3.448276 6.896552 6.896552 0.00000 17.241379 13.79310
##          36          37          38          39          40          41          42
## knn      8.571429 8.571429 8.571429 20.00000 11.42857 20.00000 17.142857
## lda      8.333333 13.88889 13.88889 22.22222 16.66667 11.11111 16.66667
## logistic 6.896552 6.896552 13.793103 17.24138 10.34483 10.34483 6.896552
##          43          44          45          46          47          48          49
## knn      11.42857 8.571429 17.142857 11.428571 14.285714 8.571429 2.857143
## lda      11.11111 8.333333 8.333333 2.777778 2.777778 19.444444 8.333333
## logistic 24.13793 6.896552 10.344828 27.586207 10.344828 6.896552 6.896552
##          50
## knn      14.285714
## lda      13.88889
## logistic 3.448276
```

```
fit.super$ranking.base
```

```
##      knn      lda logistic
##      35      36      29
```

# Variable Screening via Random Subspace Ensemble

In this section, we describe how to apply RaSE for variable screening.

## Introduction

We follow the aforementioned notations. Although Tian and Feng (2021b) only discusses the classification problem, RaSE framework can be immediately applied for continuous response with no extra effort. As before, we would like to select  $B_1$  subspaces  $\{S_{j*}\}_{j=1}^{B_1}$ , for each of which  $B_2$  candidate subspaces  $\{S_{jk}\}_{k=1}^{B_2}$  are generated. Some specific criterion is required for subspace selection. The selected percentage of features in  $\{S_{j*}\}_{j=1}^{B_1}$  can be used for variable screening (Tian and Feng (2021a)).

## Variable screening via RaSE

We will present how to do variable screening through RaSE framework in this section. First we generate the data from the following model (model 1 in Tian and Feng (2021b), model II in Fan and Lv (2008)).

$$y = 5x_1 + 5x_5 + 5x_3 - \frac{15}{\sqrt{2}}x_4 + \epsilon,$$

where  $\mathbf{x} = (x_1, \dots, x_p)^T \sim N(\mathbf{0}, \Sigma)$ ,  $\Sigma = (\sigma_{ij})_{p \times p}$ ,  $\sigma_{ij} = 0.5\mathbb{1}(i \neq j)$ ,  $\epsilon \sim N(0, 1)$ , and  $\epsilon \perp \mathbf{x}$ . The signal set  $S^* = \{1, 2, 3, 4\}$ .

Let  $n = 100$  and  $p = 100$ . Call function `RaModel` to generate the data.

```
train.data <- RaModel("screening", 1, n = 100, p = 100)
xtrain <- train.data$x
ytrain <- train.data$y
```

As Tian and Feng (2021b) describes, here  $x_4$  is marginally independent of  $y$ . We first apply RaSE equipped with linear regression model and BIC by calling function `RaScreen`. Set  $B_1 = B_2 = 100$ , `model = lm` and `criterion = bic`. To demonstrate the power of iterations, we set `iteration = 1`.

```
fit.bic <- RaScreen(xtrain, ytrain, B1 = 100, B2 = 100, model = "lm", criterion = "bic",
  cores = 2, iteration = 1)
```

The output of `RaScreen` is a list including the selected percentage of variables, the model we use and other information. Note that the selected percentage of variables are stored in a list (when `iteration = 0`, it is a vector) called “selected.perc.” All results from different iteration rounds are available. Function `RaRank` provides a convenient and automatic way to select variables from the output of `RaScreen`. We set `selected.num = n/logn` to select  $\lceil n/\log n \rceil = 21$  variables. Let’s compare the results from vanilla RaSE (no iteration) and RaSE with 1 iteration round as follows.

```
RaRank(fit.bic, selected.num = "n/logn", iteration = 0)
```

```
## [1] 2 3 1 4 62 46 51 22 58 80 20 54 57 63 83 98 99 7 12 14 43
```

```
RaRank(fit.bic, selected.num = "n/logn", iteration = 1)
```

```
## [1] 2 3 4 1 58 62 46 13 51 45 55 80 98 10 20 43 63 54 57 74 83
```



Observe that RaSE with linear regression model and BIC captures the signals very well. Iteration indeed improves vanilla RaSE by ranking four signals on the top.

Note that there could be some variables with zero selected percentage, especially in iterative RaSE. In this case, such variables are indistinguishable. When the user requests more variables beyond the number of variables with positive selected percentages, **RaRank** function will randomly sample from variables with zero selected percentage and pop up a warning message as below.

```
RaRank(fit.bic, selected.num = "n-1", iteration = 1)
```

```
## Warning in RaRank(fit.bic, selected.num = "n-1", iteration = 1): Only 85
## variables have positive selected percentage but request 99 ones. The last 14
## variables are randomly sampled!
```

```
## [1] 2 3 4 1 58 62 46 13 51 45 55 80 98 10 20 43 63 54 57
## [20] 74 83 89 11 19 22 33 48 49 78 12 36 53 42 47 86 88 93 6
## [39] 7 14 16 21 26 27 29 31 40 41 56 59 68 79 81 91 97 100 15
## [58] 24 30 32 61 70 73 76 92 95 17 28 39 67 72 85 94 9 18 25
## [77] 35 37 65 84 44 52 64 71 99 5 82 77 34 96 66 23 90 60 75
## [96] 8 87 69 50
```

## Reference

- Chen, Jiahua, and Zehua Chen. 2008. "Extended Bayesian Information Criteria for Model Selection with Large Model Spaces." *Biometrika* 95 (3): 759–71.
- . 2012. "Extended BIC for Small-n-Large-p Sparse GLM." *Statistica Sinica*, 555–74.
- Fan, Jianqing, and Jinchi Lv. 2008. "Sure Independence Screening for Ultrahigh Dimensional Feature Space." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 70 (5): 849–911.
- Mai, Qing, Hui Zou, and Ming Yuan. 2012. "A Direct Approach to Sparse Discriminant Analysis in Ultra-High Dimensions." *Biometrika* 99 (1): 29–42.
- Tian, Ye, and Yang Feng. 2021a. "RaSE: A Variable Screening Framework via Random Subspace Ensembles." *Journal of the American Statistical Association*, no. just-accepted: 1–30.
- . 2021b. "RaSE: Random Subspace Ensemble Classification." *Journal of Machine Learning Research* 22 (45): 1–93.
- Zhu, Jianan, and Yang Feng. 2021. "Super RaSE: Super Random Subspace Ensemble Classification." <https://www.preprints.org/manuscript/202110.0042>.