

# MixAll: Learning mixture models

Serge Iovleff  
University Lille 1

---

## Abstract

The MixALL package can also be used in order to learn mixture models when the labels class are known. This short vignette assume that you have already read the vignette "Clustering With MixAll" (Iovleff (2016)).

*Keywords:* R, C++, STK++, Learning, missing values.

---

## 1. Introduction

It is possible to perform supervised learning with MixAll when the labels of the individuals are known. Let us recall the notations defined in the Iovleff (2016) vignette.  $\mathcal{X}$  denote an arbitrary measurable space,  $\mathcal{Z} = \{1, \dots, K\}$  is the label set and  $(\mathbf{x}, \mathbf{z}) = \{(\mathbf{x}_1, \mathbf{z}_1), \dots, (\mathbf{x}_n, \mathbf{z}_n)\}$  represents  $n$  independent vectors in  $\mathcal{X} \times \mathcal{Z}$  such that each  $\Pr(\mathbf{z}_i = k) = p_k$  and such that conditionnaly to  $\mathbf{z}_i = k$ ,  $\mathbf{x}_i$  arises from a probability distribution with density

$$h(\mathbf{x}_i | \boldsymbol{\lambda}_k, \boldsymbol{\alpha}) \tag{1}$$

parameterized by  $\boldsymbol{\lambda}_k$  and  $\boldsymbol{\alpha}$ .

Given the matrix of observation  $\mathbf{x}$ , and the vector of labels  $\mathbf{z}$ , the learning methods will estimate the unknown parameters  $\boldsymbol{\lambda}_k$  and  $\boldsymbol{\alpha}$ .

## 2. Learning with MixAll

Learning analysis can be performed with the functions

1. `learnDiagGaussian` for diagonal Gaussian mixture models,
2. `learnCategorical` for Categorical mixture models,
3. `learnPoisson` for Poisson mixture models,
4. `learnGamma` for gamma mixture models,
5. `learnMixedData` for MixedData mixture models.

These functions have a common set of parameters with default values given in the table 1.

| Input Parameter        | Description                                                                                                                                                                                                                                                                      |
|------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <code>data</code>      | A matrix (or a list of matrix for mixed data) with the data to learn.                                                                                                                                                                                                            |
| <code>labels</code>    | A vector with the classes of each individuals. Values must be between 1 and $K$ .                                                                                                                                                                                                |
| <code>models</code>    | A vector with the models to adjust to each data set in case of mixed data, or a set of models to try to adjust. Default is <code>cluster*Names()</code> where '*' stands for <code>DiagGaussian</code> , <code>Poisson</code> , <code>Gamma</code> or <code>Categorical</code> . |
| <code>prop</code>      | A vector of size $K$ with the proportions of each class. If <code>prop</code> is <code>NULL</code> then the proportions are computed using the empirical distribution of the <code>labels</code> .                                                                               |
| <code>algo</code>      | A string defining the algorithm to use for the missing values. Possible values <code>"impute"</code> , <code>"simul"</code> .                                                                                                                                                    |
| <code>nbIter</code>    | maximal number of iteration to perform. Default value is 100. Note that if there is no missing values, it should be 1.                                                                                                                                                           |
| <code>epsilon</code>   | threshold to use in order to stop the iterations (not used by the <code>"simul"</code> algorithm). Default value <code>1e-08</code> .                                                                                                                                            |
| <code>criterion</code> | A string defining the model selection criterion to use. The best model is the one with the lowest criterion value. Possible values: <code>"AIC"</code> , <code>"BIC"</code> , <code>"ICL"</code> . Default is <code>"ICL"</code> .                                               |
| <code>nbCore</code>    | An integer defining the number of processor to use. Default is 1, 0 for all cores.                                                                                                                                                                                               |

Table 1: List of common parameters of the learning functions.

The `learnKernel` function has two more arguments described in table 2.

| Input Parameter               | Description                                                                                                                                                                                   |
|-------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <code>kernelName</code>       | A string defining the kernel to use. Use a <code>"gaussian"</code> kernel by default. Possible values are <code>"gaussian"</code> , <code>"polynomial"</code> or <code>"exponential"</code> . |
| <code>kernelParameters</code> | A vector with the kernel parameter value(s). Default value is 1.                                                                                                                              |

Table 2: List of all the specific parameters of the `learnKernel` function.

## 2.1. Learning Multivariate (diagonal) Gaussian Mixture Models

Multivariate Gaussian mixture models (without correlations) can be learned using the `learnDiagGaussian` function. We illustrate this function with the well known geysers data set (Azzalini and Bowman (1990), Härdle (1991)).

```
> data(iris);
> x <- as.matrix(iris[,1:4]); z <- as.vector(iris[,5]); n <- nrow(x); p <- ncol(x);
```

```
> indexes <- matrix(c(round(runif(5,1,n)), round(runif(5,1,p))), ncol=2);  
> cbind(indexes, x[indexes]) # display true values
```

```
      [,1] [,2] [,3]  
[1,]   29   4 0.2  
[2,]  106   1 7.6  
[3,]   86   4 1.6  
[4,]   26   2 3.0  
[5,]  142   3 5.1
```

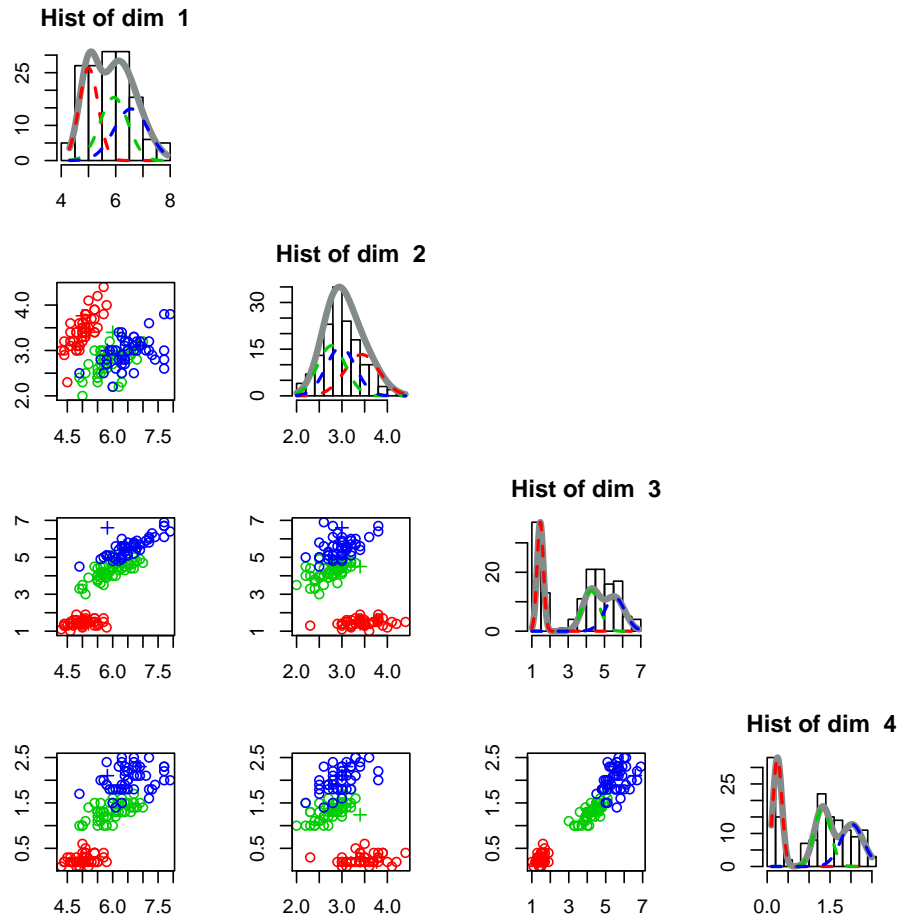
```
> x[indexes] <- NA;          # set them as missing  
> model <- learnDiagGaussian(data=x, labels = z, models = clusterDiagGaussianNames(prop =  
> summary(model)
```

```
*****  
* model name      = gaussian_p_sjk  
* nbSample       = 150  
* nbCluster      = 3  
* lnLikelihood   = -1023.296  
* nbFreeParameter= 70  
* criterion name = ICL  
* criterion value= 2397.336  
*****
```

```
> missingValues(model)
```

```
  row col  value  
1 106  1 5.823991  
2  26  2 3.764123  
3 142  3 5.544055  
4  29  4 0.267857  
5  86  4 1.237870
```

```
> plot(model)
```



## 2.2. Learning Multivariate categorical Mixture Models

Categorical (nominal) data can be learned using the `learnCategorical` function.

We illustrate this function with the birds data set.

```
> data(birds)
> ## add 10 missing values
> x <- as.matrix(birds[,2:5]); z <- as.vector(birds[,1]); n <- nrow(x); p <- ncol(x);
> indexes <- matrix(c(round(runif(5,1,n)), round(runif(5,1,p))), ncol=2);
> cbind(indexes, x[indexes]) # display true values
```

```
      [,1] [,2] [,3]
[1,] "14" "2"  "dotted"
[2,] "14" "3"  "black & white"
[3,] "58" "3"  "white"
[4,] "13" "3"  "black"
[5,] "42" "1"  "pronounced"
```

```
> x[indexes] <- NA;          # set them as missing
> model <- learnCategorical( data=x, labels=z
```

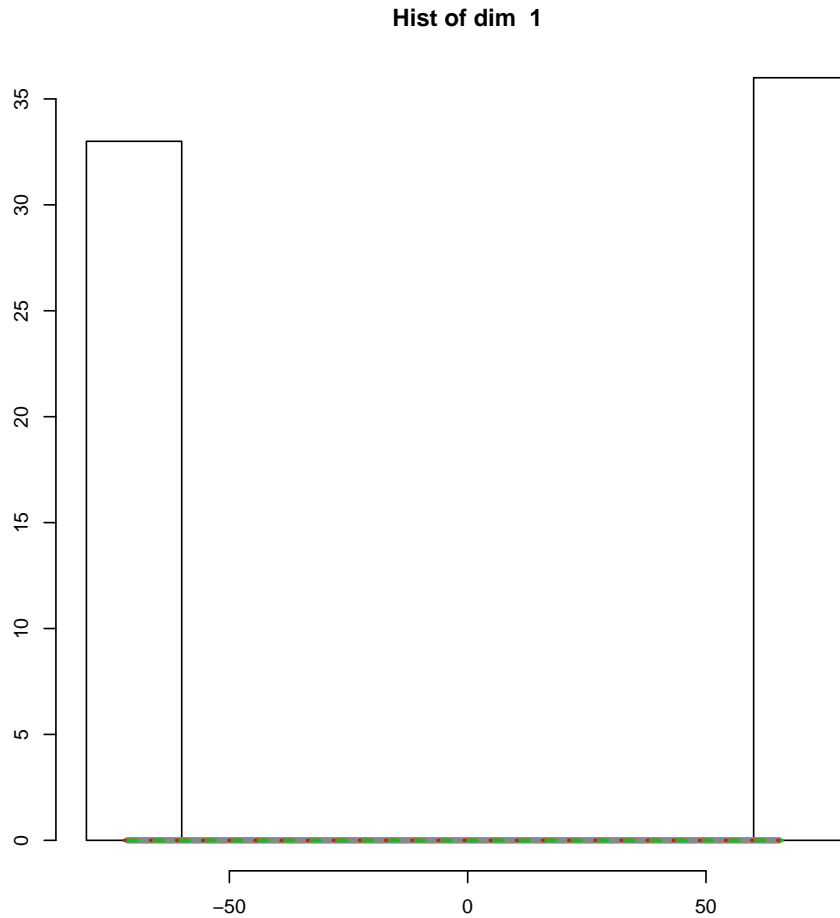
```
+           , models = clusterCategoricalNames(prop = "equal")
+           , algo="simul", nbIter = 2)
> summary(model)
```

```
*****
* model name      = categorical_p_pjk
*****
* nbModalities   = 4
*****
* levels =
[1] "none           , poor pronounced, pronounced       , very pronounced"
[2] "dotted, none   "
[3] "black          , black & WHITE, black & white, white       "
[4] "few , many, none"
*****
* nbSample       = 69
* nbCluster      = 2
* lnLikelihood   = -553.9514
* nbFreeParameter= 25
* criterion name = ICL
* criterion value= 1213.756
*****
```

```
> missingValues(model)
```

```
   row col value
1  42  1    3
2  14  2    1
3  13  3    4
4  14  3    4
5  58  3    4
```

```
> plot(model)
```



### 2.3. Learning Multivariate Gamma Mixture Models

Gamma data can be learned using the `learnGamma` function.

We illustrate this function with the iris data set.

```
> data(iris)
> x <- as.matrix(iris[,1:4]); z <- as.vector(iris[,5]); n <- nrow(x); p <- ncol(x);
> indexes <- matrix(c(round(runif(5,1,n)), round(runif(5,1,p))), ncol=2);
> cbind(indexes, x[indexes]) # display true values
```

```
      [,1] [,2] [,3]
[1,]   75   2  2.9
[2,]  144   3  5.9
[3,]  100   4  1.3
[4,]  120   1  6.0
[5,]  116   3  5.3
```

```
> x[indexes] <- NA;          # set them as missing
> model <- learnGamma( data=x, labels= z,
+                      , models = clusterGammaNames(prop = "equal")
```

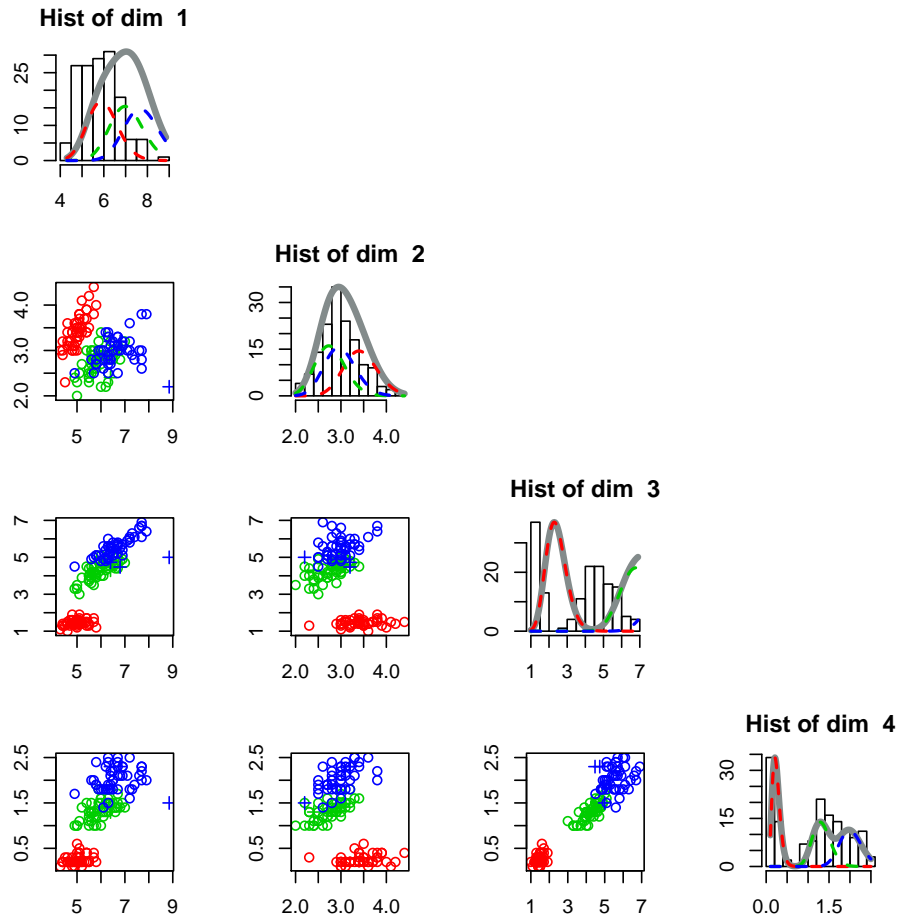
```
+           , algo = "simul", nbIter = 2, epsilon = 1e-08
+         )
> summary(model)
```

```
*****
* model name      = gamma_p_ajk_bj
* nbSample       = 150
* nbCluster      = 3
* lnLikelihood   = -15417.01
* nbFreeParameter= 142
* criterion name = ICL
* criterion value= 31545.53
*****
```

```
> missingValues(model)
```

```
  row col  value
1  120   1 8.850716
2   75   2 2.401010
3  116   3 4.730563
4  144   3 4.474997
5  100   4 1.294791
```

```
> plot(model)
```



## 2.4. Learning Multivariate Poisson Models

Poisson data (count data) can be learned using the `learnPoisson` function.

We illustrate this function with the `debTrivedi` data set.

```
> data(DebTrivedi)
> x <- DebTrivedi[, c(1, 6, 8, 15)]; z <- DebTrivedi$medicaid; n <- nrow(x); p <- ncol(x);
> indexes <- matrix(c(round(runif(5,1,n)), round(runif(5,1,p))), ncol=2);
> cbind(indexes, x[indexes]) # display true values
```

```
      [,1] [,2] [,3]
[1,] 2550   3   5
[2,] 2974   4   8
[3,] 1744   3   1
[4,]   88   2   1
[5,] 1990   3   2
```

```
> x[indexes] <- NA; # set them as missing
> model <- learnPoisson( data=x, labels=z
+ , models = clusterPoissonNames(prop = "equal")
```



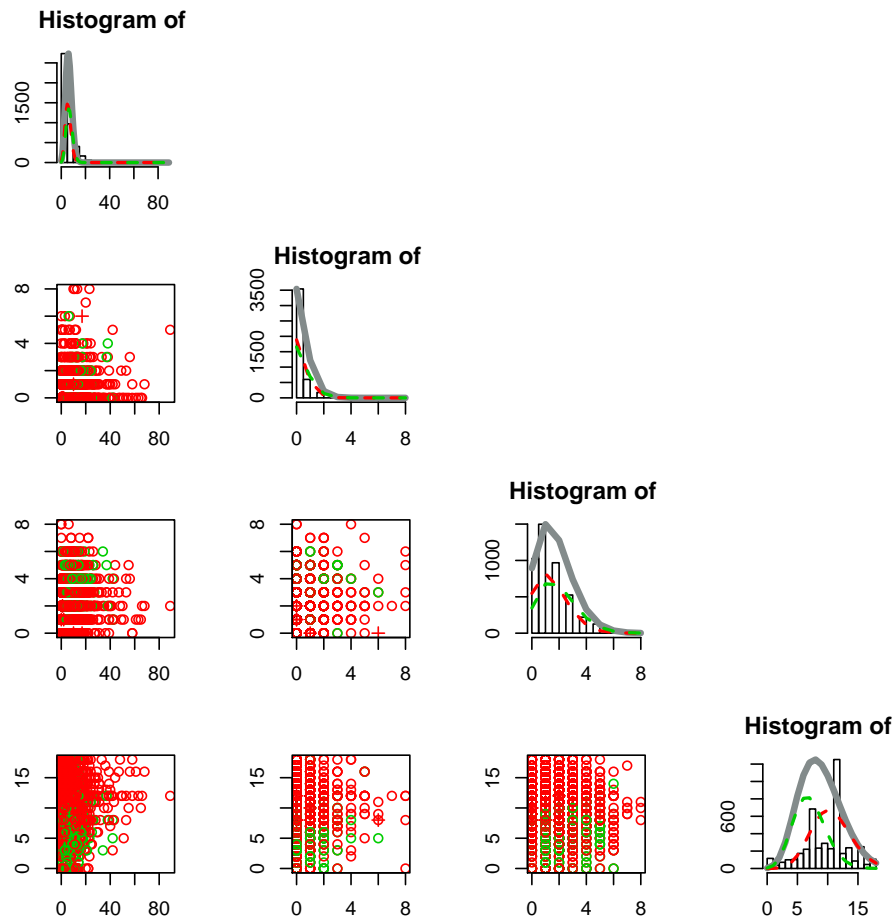
```
+           , algo="simul", nbIter = 2, epsilon = 1e-08
+ )
> summary(model)
```

```
*****
* model name      = poisson_p_ljk
* nbSample       = 4406
* nbCluster      = 2
* lnLikelihood   = -161275.9
* nbFreeParameter= 17
* criterion name = ICL
* criterion value= 322694.5
*****
```

```
> missingValues(model)
```

```
   row col value
1   88  2     1
2 1744  3     1
3 1990  3     2
4 2550  3     0
5 2974  4     8
```

```
> plot(model)
```



## 2.5. Learning Mixed data sets

Mixed data sets can be learned using the `learnMixedData` function. The original mixed data set has to be splitted in multiple homogeneous data sets and each one associated to a mixture model name.

We illustrate this function with the HeartDisease data set.

```
> data(HeartDisease.cat)
> data(HeartDisease.cont)
> data(HeartDisease.target)
> ldata = list(HeartDisease.cat, HeartDisease.cont);
> models = c("categorical_pk_pjk", "gaussian_pk_sjk")
> z<-HeartDisease.target[[1]];
> model <- learnMixedData(ldata, models, z, algo="simul", nbIter=2)
> summary(model)
```

```
*****
* model name = categorical_pk_pjk
* model name = gaussian_pk_sjk
* nbSample   = 303
```

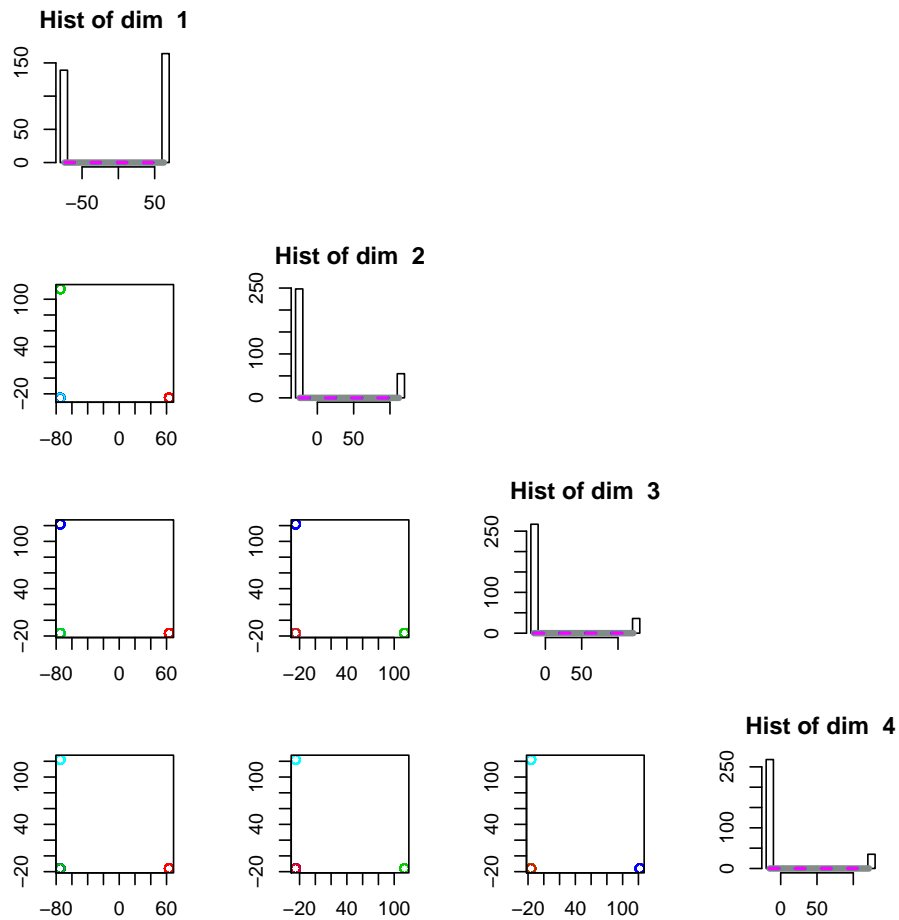
```
* nbCluster      = 5
* lnLikelihood   = -7531.688
* nbFreeParameter= 129
* criterion name = ICL
* criterion value= 15800.45
*****
```

```
> missingValues(model)
```

```
[[1]]
      row col value
[1,] 167  7     1
[2,] 193  7     2
[3,] 288  7     2
[4,] 303  7     1
[5,]  88  8     3
[6,] 267  8     3
```

```
[[2]]
      row col value
```

```
> plot(model)
```



## References

- Azzalini A, Bowman AW (1990). “A look at some data on the Old Faithful geyser.” *Applied Statistics*, pp. 357–365.
- Härdle W (1991). *Smoothing techniques: with implementation in S*. Springer Science & Business Media.
- Iovleff S (2016). *Clustering With MixAll*. R package version 1.1.1, URL <http://CRAN.R-Project.org/package=MixAll>.

### Affiliation:

Serge Iovleff  
 Univ. Lille 1, CNRS U.M.R. 8524, Inria Lille Nord Europe  
 59655 Villeneuve d’Ascq Cedex, France  
 E-mail: [Serge.Iovleff@stkpp.org](mailto:Serge.Iovleff@stkpp.org)

URL: <http://www.stkpp.org>