

# pathClass: SVM-based classification with prior knowledge on feature connectivity

(Version 0.8.0)

## User's Guide

Marc Johannes  
German Cancer Research Center  
Heidelberg, Germany

February 8, 2011

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>What data do we need</b>	<b>2</b>
2.1	The class labels . . . . .	2
2.2	The data matrix . . . . .	2
2.3	The graph . . . . .	2
2.4	The mapping . . . . .	3
<b>3</b>	<b>Which classification methods are available</b>	<b>4</b>
3.1	Rewighted Recursive Feature Elimination . . . . .	4
3.2	network-based SVM . . . . .	5
3.3	graph SVM . . . . .	6
<b>4</b>	<b>Showing the results</b>	<b>6</b>

## 1 Introduction

The package `pathClass` was developed for classification tasks with the usage of *prior* knowledge about the feature connectivity. At the German Cancer Research Center we are dealing mostly with biological data. Thus, in this vignette we demonstrate the usage of the package and its functions using biologically data.

The package can be loaded by typing:

```
> library(pathClass)
```

## 2 What data do we need

For a *standard* classification task one needs a data matrix to train on as well as class labels which tell the algorithm to what class a sample belongs to. However, we now have an additional source of knowledge, i.e. a graph structure. For the algorithm to know which feature in the data matrix corresponds to which node in the graph we need a mapping as well. In the following sections we will describe the structure of these data objects and give examples how to create and use them.

### 2.1 The class labels

In this vignette we are going to use the combined test and training sets from the Golub paper which is part of the R package `golubEsets` available on Bioconductor:

```
> library(golubEsets)
> data(Golub_Merge)
```

We are going to predict whether the patient had AML or ALL. Hence, we create the class labels as:

```
> y <- pData(Golub_Merge)$ALL
```

From this output we can see that `y` contains 72 entries. That is, 72 patients are used for the analysis.

### 2.2 The data matrix

Next, we need the corresponding expression data as data matrix  $\mathbf{D}^{n \times p}$  with  $n$  samples of  $p$  measurements.

```
> x <- exprs(Golub_Merge)
```

This data set contains 7129 features measured in 72 samples. However, we need the transposed version of it:

```
> x <- t(x)
> dim(x)
```

```
[1] 72 7129
```

### 2.3 The graph

As a next step we have to create a adjacency matrix that represents the connectivity of the features in `x`. Therefore, we download from <http://www.hprd.org/download> the file of binary protein-protein interactions in tab delimited format. After extracting the archive we use `pathClass` to read the tab-delimited file:

```
> hprd <- read.hprd('BINARY_PROTEIN_PROTEIN_INTERACTIONS.txt')
```

Since most classification algorithm can “only” use those features that are present in both, the data matrix `x` and the `hprd` adjacency matrix, we have to match both objects to each other. Therefore, we need a mapping containing the information which protein of `hprd` matches to which probe set in `x`.

As an alternative, the user can load a small fake network using the command:

```
> data(adjacency.matrix)
> hprd <- adjacency.matrix
```

## 2.4 The mapping

For most microarrays there is a annotation package available. Since we are dealing with expression data from chip hu6800 we load the corresponding annotation package and create a mapping from probe set ID to protein ID:

```
> ann <- annotation(Golub_Merge)
> library(paste(ann,'db',sep='.'), character.only=TRUE)
> graphIDs <- "REFSEQ"
> rs <- get(paste(ann, graphIDs, sep=''))
> refseq <- mget(featureNames(Golub_Merge), rs)
> times <- sapply(refseq, length)
> mapping <- data.frame(probesetID=rep(names(refseq), times=times),
+                       graphID=unlist(refseq),
+                       row.names=NULL,
+                       stringsAsFactors=FALSE)
> nas <- which(is.na(mapping[, 'graphID']))
> mapping <- mapping[-nas,]
> mapping <- unique(mapping)
> head(mapping)
```

	probesetID	graphID
35	AFFX-HUMISGF3A/M97935_5_at	NM_007315
36	AFFX-HUMISGF3A/M97935_5_at	NM_139266
37	AFFX-HUMISGF3A/M97935_5_at	NP_009330
38	AFFX-HUMISGF3A/M97935_5_at	NP_644671
39	AFFX-HUMISGF3A/M97935_MA_at	NM_007315
40	AFFX-HUMISGF3A/M97935_MA_at	NM_139266

The first line in the above code-chunk identifies the annotation of the expression set. The second line load the corresponding annotation package. Line three defines the kind of IDs that are present in the graph structure, must be one of `ls(paste('package:', ann, '.db', sep=■))`. Lines four and five extract the graph IDs, in this case REFSEQ, that match to the `featureNames` of our expression set. The remaing code puts everything into the order needed by `pathClass`. Now we have a mapping with 23686 rows. It is important that this mapping has at least two columns named `graphID` and `probesetID` since those names are needed internally when `pathClass` makes use of the mapping.

In a next step we can make use of the function `matchMatrices()` to match the data matrix `x` to the `hprd`:

```
> matched <- matchMatrices(x=x, adjacency=hprd, mapping=mapping)
```

The list `matched` contains copies of `x`, `hprd` and `mapping` however with matching dimensions. Thus, these objects can now be used for classification.

### 3 Which classification methods are available

That far, all classification algorithms we implemented are based on the support vector machine (SVM, [Vapnik and Cortes 1995](#)). As a standard tool we provide the recursive feature elimination (SVM-RFE, [Guyon et al. 2002](#)) algorithm for the SVM. This algorithm performs a feature selection, however it makes no use of *prior* knowledge. In addition to SVM-RFE we implemented three other SVM-based algorithm that use *prior* knowledge:

1. Reweighted Recursive Feature Elimination (RRFE, [Johannes et al. 2010](#))
2. Network-based SVM ([Zhu et al., 2009](#))
3. Graph SVM ([Rapaport et al., 2007](#))

The functions to train these methods are called: `fit.rfe`, `fit.rrfe`, `fit.networkBasedSVM` and `fit.graph.svm`, respectively. The user can use these functions directly to obtain a fit object of the corresponding algorithm or use the wrapper-function `crossval()` to perform a  $x$  times repeated  $y$ -fold cross-validation. Additionally the `crossval()` function is able to make use of the multicore architecture of modern PCs or a computing cluster. To use the parallel version of the method the user has to load the library `multicore` prior to calling `crossval()` and to set the parameter `parallel` to `TRUE`. It is, however, worth mentioning that for parallel use all data objects are copied for each of the CPU processes. Therefore, one has to ensure that the object fits into the memory of the server.

For the purpose of reproducibility we initialize the random number generator prior to calling the individual algorithms. This also ensures that each algorithm uses the same splits within the cross-validation.

#### 3.1 Reweighted Recursive Feature Elimination

The RRFE method can be run without using the mapping created above. The reason for this is, that the method can use all features if the user sets the parameter `useAllFeatures` to `TRUE`. Therefore, this method has its own, internal mapping routine. RRFE has a tuning parameter  $d \in (0, 1)$  that controls the influence of the graph structure on the ranking of the genes. A value of  $d \rightarrow 1$  puts more weight on the connectivity information whereas  $d \rightarrow 0$  relies more on the expression data. To use the RRFE method one can use:

```

> set.seed(12345)
> res.rrfe <- crossval(x,
+                      y,
+                      DEBUG=TRUE,
+                      theta.fit=fit.rrfe,
+                      folds=10,
+                      repeats=5,
+                      parallel=TRUE,
+                      Cs=10^(-3:3),
+                      mapping=mapping,
+                      Gsub=hprd,
+                      d=1/2)

```

or, to use all features:

```

> res.rrfe <- crossval(x,
+                      y,
+                      DEBUG=TRUE,
+                      theta.fit=fit.rrfe,
+                      folds=10,
+                      repeats=5,
+                      parallel=TRUE,
+                      Cs=10^(-3:3),
+                      useAllFeatures=TRUE,
+                      mapping=mapping,
+                      Gsub=hprd,
+                      d=1/2)

```

Please, have a look into the help files or the paper ([Johannes et al., 2010](#)) for more information on the `useAllFeatures` option.

## 3.2 network-based SVM

The network-based support vector machine ([Zhu et al., 2009](#)) needs the mapping from above, since the dimensions of the data objects have to match exactly. However, instead of an adjacency matrix it needs an adjacency list which we have to create before:

```

> ad.list <- as.adjacencyList(matched$adjacency)
> set.seed(12345)
> res.nBSVM <- crossval(matched$x,
+                      y,
+                      theta.fit=fit.networkBasedSVM,
+                      folds=10,
+                      repeats=5,
+                      DEBUG=TRUE,
+                      parallel=FALSE,

```

```

+             adjacencyList=ad.list,
+             lambdas=10^(-1:2),
+             sd.cutoff=150)

```

Since, the algorithm internally uses `lpSolve`, it has to calculate a constraints-matrix. Thus, when having lots of features this matrix can become very big. Therefore, we added the parameter `sd.cutoff` which only keeps genes with standard deviation  $\geq$  `sd.cutoff`. Further, we recommend not to run this algorithm in parallel, since the constraints matrix is created by each process, which might result in memory overflow.

### 3.3 graph SVM

Rapaport et al. (2007) developed a supervised classification framework which we refer to as “graph SVM”. This method makes use of a so-called diffusion kernel, which has to be calculated before using this method:

```

> dk <- calc.diffusionKernel(L=matched$adjacency,
+                           is.adjacency=TRUE,
+                           beta=0)
> set.seed(12345)
> res.gSVM <- crossval(matched$x,
+                      y,
+                      theta.fit=fit.graph.svm,
+                      folds=10,
+                      repeats=5,
+                      DEBUG=TRUE,
+                      parallel=FALSE,
+                      Cs=10^(-3:3),
+                      mapping=matched$mapping,
+                      diffusionKernel=dk)

```

Where `beta` is a tuning parameter that controls the extent of diffusion. This parameter should be optimized.

## 4 Showing the results

We can have a look on the individual results by typing:

```

> plot(res.rrfe, toFile=F)

```

We get a boxplot for each repeat of the cross-validation showing the distribution of AUC's obtained by the classifiers trained in the repeat as well as a receiver operator characteristic (ROC) curve showing the overall performance.

We can, however, also combine all results into one ROC curve by using the `ROCR` package:

```

> cv.labels <- matrix(rep(y,5), ncol=5)
> pred.rrfe <- prediction(res.rrfe$cv, labels=cv.labels)
> auc.rrfe <- round(mean(unlist(performance(pred.rrfe, 'auc')@y.values)),3)
> plot(performance(pred.rrfe, measure = "tpr", x.measure = "fpr"),
+       col='red',
+       main='Benchmark of the algorithms',
+       avg = "threshold")
> pred.nBSVM <- prediction(res.nBSVM$cv, labels=cv.labels)
> auc.nBSVM <- round(mean(unlist(performance(pred.nBSVM, 'auc')@y.values)),3)
> plot(performance(pred.nBSVM, measure = "tpr", x.measure = "fpr"),
+       add=TRUE,
+       col='blue',
+       avg = "threshold")
> pred.gSVM <- prediction(res.gSVM$cv, labels=cv.labels)
> auc.gSVM <- round(mean(unlist(performance(pred.gSVM, 'auc')@y.values)),3)
> plot(performance(pred.gSVM, measure = "tpr", x.measure = "fpr"),
+       add=TRUE,
+       col='green',
+       avg = "threshold")
> legend('bottomright',
+       c(paste('RRFE (AUC=',auc.rrfe,')',sep=''),
+         paste('network based SVM (AUC=',auc.nBSVM,')',sep=''),
+         paste('graph SVM (AUC=',auc.gSVM,')',sep='')),
+       text.col=c('red','blue','green'),
+       col=c('red','blue','green'),
+       lty=1,
+       bty='n',
+       cex=1.3)
> abline(b=1,a=0,col='gray')

```

It is important to note that these results can not be generalized to be true for all data sets. We provide this package, that the user can easily evaluate all algorithms and based on these result choose the best one.

These commands produce figure 1. Additionally we can extract the features which have been chosen by the classifier by using the following function:

```

> extractFeatures(res.rrfe, toFile=T, fName='OurFeatures.csv')

```

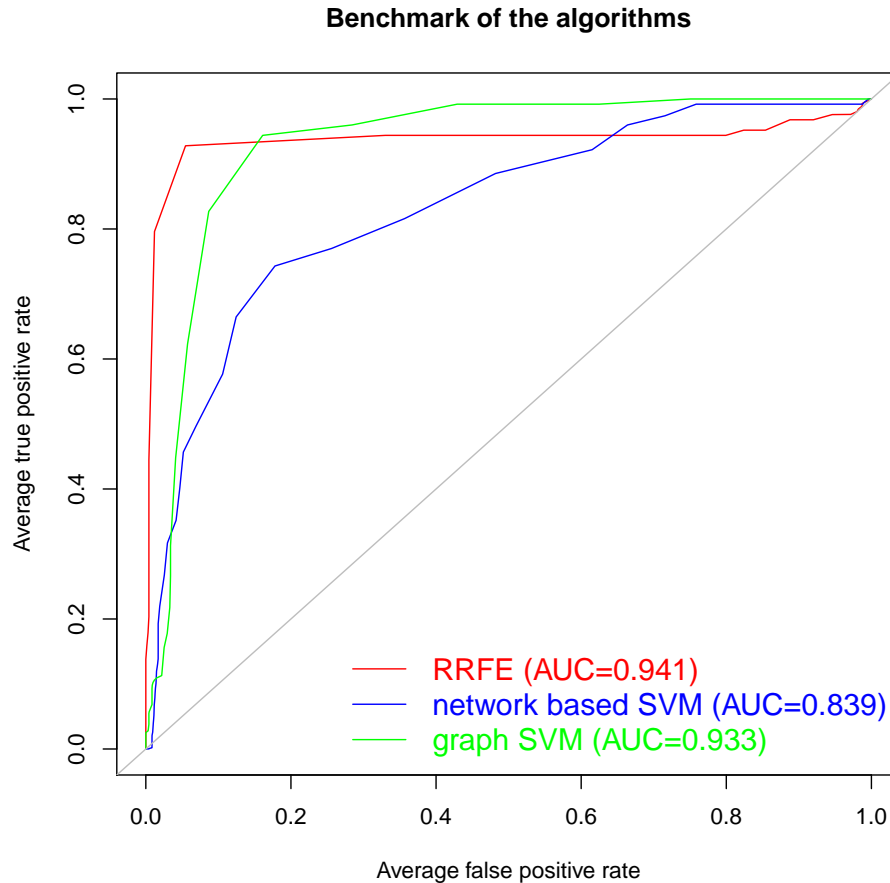


Figure 1: ROC curves for all three algorithms

## References

- I. Guyon, J. Weston, S. Barnhill, and V. Vapnik. Gene selection for cancer classification using support vector machines. *Machine Learning*, 46(1-3):389–422, 2002. URL <http://www.springerlink.com/index/W68424066825VR3L.pdf>.
- M. Johannes, J. C. Brase, H. Fröhlich, S. Gade, M. Gehrmann, M. Fälth, H. Sültmann, and T. Beißbarth. Integration of pathway knowledge into a reweighted recursive feature elimination approach for risk stratification of cancer patients. *Bioinformatics*, 26(17):2136–2144, Jun 2010. doi: 10.1093/bioinformatics/btq345. URL <http://dx.doi.org/10.1093/bioinformatics/btq345>.
- F. Rapaport, A. Zinovyev, M. Dutreix, E. Barillot, and J.-P. Vert. Classification of microarray data using gene networks. *BMC Bioinformatics*, 8:35, 2007. doi: 10.1186/1471-2105-8-35. URL <http://dx.doi.org/10.1186/1471-2105-8-35>.
- V. Vapnik and C. Cortes. Support-vector networks. *Machine Learning*, Jan 1995. URL <http://www.springerlink.com/index/K238JX04HM87J80G.pdf>.
- Y. Zhu, X. Shen, and W. Pan. Network-based support vector machine for classification of microarray



samples. *BMC Bioinformatics*, 10 Suppl 1:S21, 2009. doi: 10.1186/1471-2105-10-S1-S21. URL <http://dx.doi.org/10.1186/1471-2105-10-S1-S21>.