

Using the BioSeqClass Package

Hong Li^{‡*}

May 3, 2016

[‡]Key Lab of Systems Biology
Shanghai Institutes for Biological Sciences
Chinese Academy of Sciences, P. R. China

Contents

1 Overview

There are now 863 completely sequenced genomes of cellular organisms in NCBI genome database. Nevertheless, functional annotation drops far behind sequencing because functional validation experiments are time-consuming and costly. Taken model organism *Homo sapiens*, *Mus musculus* and *Saccharomyces cerevisiae* as examples, only 16 and 18 annotations in Gene Ontology), respectively. Thus computational methods for predicting function is still a fundamental complement. The most common computation approach is biological sequence based classification, since sequence information is still the most abundant and reliable. Sequence based classification has been used in: discovering new microRNA candidates, predicting transcription factor binding sites, detecting protein post-translational modification sites, and so on.

Features and models are two basic factors for classification. Features generally are numerical values that can be used to distinguish different classes. Therefore it is preferable to select features that can achieve better and faster classification. Classification models are built from features by various algorithms, and it is necessary to evaluate its prediction ability by cross validation or jackknife test. For biological sequences, there are additional steps: one is to reduce homolog sequences which might result in overestimation of prediction accuracy, and then another most important step is to convert sequences into numerical features. Thus, the general workflow for sequence-based classifications includes (Figure ??): reduce homolog sequences; extract features from sequences and code them to numerical values; evaluate and select features; build classification model and evaluate its performance.

*sysptm@gmail.com

Here we present an R package (BioSeqClass) to carry out the general workflow for biological sequence based classification. It contains diverse feature coding schemas for RNA, DNA and proteins, supports feature selection, and integrates multiple classification methods.

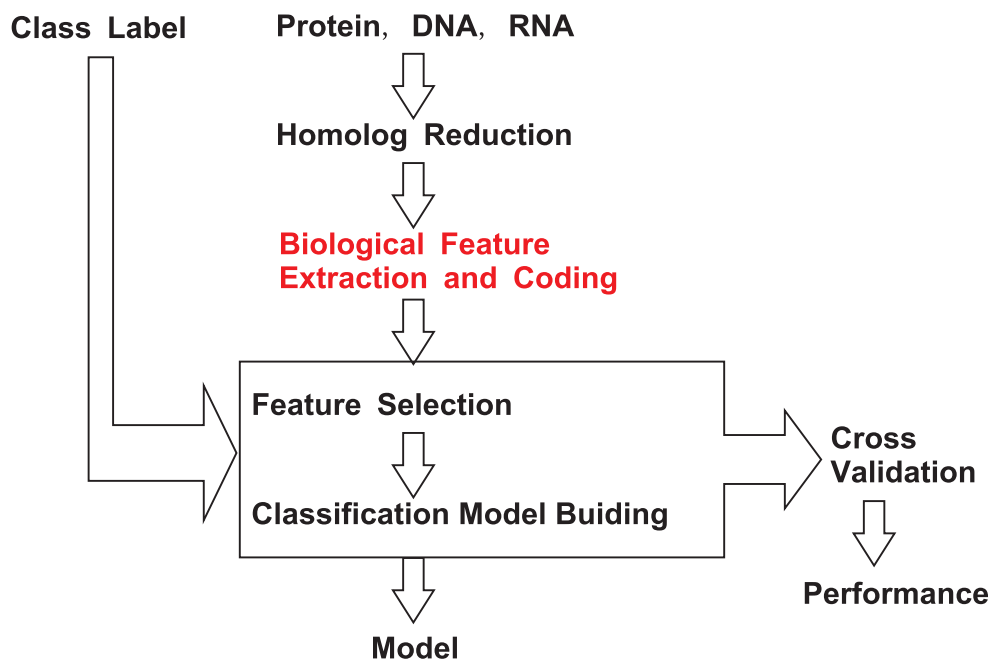


Figure 1: Workflow for Biological Sequence based Classification.

2 Installation

2.1 Requirements

BioSeqClass employs some external programs to extract biological properties and use other R packages to build classification model:

1. BioSeqClass imported R packages are listed in table ???. These packages will be automatically installed when Biocalss is firstly loaded.
2. External programs are used to assist the performance of BioSeqClass (see table ??). Some programs are invoked via their web service, and some ones are needed to be installed at the local computer.

Note: You do not need to install programs listed in table ??, unless you will use the related function in BioSeqClass.

2.2 Installation

The biocLite script is used to install PAnnBuilder from within R:

```
> source("http://bioconductor.org/biocLite.R")
> biocLite("BioSeqClass")
```

Users also can use the installation script "BioSeqClass.R" to download and install BioSeqClass package.

```
> source("http://www.biosino.org/download/BioSeqClass/BioSeqClass.R")
> BioSeqClass()
```

Load package:

```
> library(BioSeqClass)
```

Note: Web Connection is needed to install BioSeqClass and its required packages. All the codes in this vignette were tested in R 2.8.0 and 2.9.0, thus the latest R version is recommended.

3 Function description

3.1 Homolog Reduction

Homologous sequences in training/testing data may lead to overestimation of prediction accuracy. Therefore, the first step for sequence based classification and prediction is homolog reduction based on sequence similarity. Taking computation complexity and similarity restriction into consideration, homolog reductions for full-length sequences and fragment sequences are different. We have designed different functions to deal with them, respectively (see table ??).

- **hr** - It employs **cdhitHR** and **aligndisHR** to filter homolog sequences by sequence similarity. **cdhitHR** is designed to filter full-length protein or gene sequences. **aligndisHR** is designed for aligned sequences with equal length.
- **cdhitHR** - It uses cd-hit program to do homolog reduction ("formatdb" and "blastall" are required for running cd-hit program). CD-HIT is a program for clustering large protein database at high sequence identity threshold (?).
- **aligndisHR** - It uses the number of different residues to do homolog reduction (?). The algorithm proceeds in a stepwise manner by first eliminating sequences that were different from another in exactly 1 position. Elimination proceeds one peptide at a time; Re-evaluate after each peptide is removed. Once no further homologs of distance 1 remain, homologs of distance 2 are eliminated, and so forth until identity between all peptides are less than given cutoff.

3.2 Feature Extraction and Numerical Coding

3.2.1 Biological sequences

RNA, DNA and protein are three kinds of basic biological sequences. RNA and DNA are composed of bases, while proteins are composed of amino acids. Furthermore, amino acids have different physical-chemical properties, which were used to divide amino acids into different groups. These elements and groups are basic objects for feature extraction (see table ??).

3.2.2 Feature Coding

Feature coding means to extract features from sequences and convert them into numerical values. The frequently used features are the basic elements of sequence (bases, amino acids), physical-chemical properties, secondary structures, and so on. There are many methods to convert features to numerical values. The simplest is the composition of element. But more sophisticated conversions are preferable for achieving better distinguishing power. Previous studies have shown that feature coding is the key point for the accuracy of classification and prediction. Here we have summarized various feature coding methods used in published papers and carried out them in BioSeqClass (see table ??). These functions will allow more diverse choices of coding strategies and accelerate the feature coding process. We also provide a function `featureEvaluate` to test the performance of models with different feature coding schemes and different classification algorithms.

3.3 Feature Selection

Features are important for the accuracy of prediction model. However, it does not mean that the more the better. Computation time is usually increased with the increase of number of features. Conflictive features would even reduce the accuracy. Therefore, suitable feature selection is needed for better prediction performance and less computation cost. We provided two functions for feature selection (see table ??).

3.4 Model Building and Performance Evaluation

Besides features, classification method is another factor that influences classification. Different cases may have different preference over classification methods. Multiple classification methods are integrated and available in BioSeqClass (see table ??). To evaluate and compare classification models, performance assessment is done for each model, including precision, sensitivity, specificity, accuracy, and matthews correlation coefficient.

Table 1: Imported R packages.

Existing R Package	Functions used by BioSeqClass
<i>Biostrings</i>	readAAStringSet, writeXStringSet
<i>e1071</i>	svm
<i>ipred</i>	bagging
<i>klaR</i>	svmlight, NaiveBayes
<i>randomForest</i>	randomForest
<i>class</i>	knn
<i>tree</i>	tree
<i>nnet</i>	nnet
<i>rpart</i>	rpart
<i>party</i>	ctree
<i>foreign</i>	write.arff
<i>Biobase</i>	addVigs2WinMenu

4 Examples

To illustrate the use of BioSeqClass, lysine acetylation site prediction is taken as an example.

1. Suppose the original data are protein FASTA sequences and lysine acetylation sites. You can use `getTrain` to extract the flanking peptides of acetylation sites as positive dataset, and filter these peptides based on sequence identity. Lysine without acetylation annotation are regarded as negative dataset, and are filtered like the positive dataset. Considering the computational time, only 20 positive data are used as examples in the following codes.

```
> library(BioSeqClass)
> # Example data in BioSeqClass.
> file=file.path(path.package("BioSeqClass"),"example","acetylation_K.fasta")
> posfile = file.path(path.package("BioSeqClass"),
+   "example", "acetylation_K.site")
> # Only a part of lysine acetylation sites are used for demo.
> posfile1=tempfile()
> write.table(read.table(posfile,sep='\t',header=F)[1:20,], posfile1,
+   sep='\t', quote=F, row.names=F, col.names=F)
> seqList = getTrain(file, posfile1, aa="K", w=7, identity=0.4)

[1] "Positive Site: 18"
[1] "Positive Protein: 9"
[1] "Positive Site After Homolog Reduction: 16"
[1] "Positive Protein After Homolog Reduction: 9"
```

Table 2: Invoked External Programs.

External Program	Description	Related BioSeqClass Function	Need Installed?	Ref
cd-hit	a program for clustering large protein database at high sequence identity threshold	cdhitHR	Yes	(?)
blastpgp	PSI-BLAST (Position-Specific Iterated BLAST) for capturing the conservation pattern	featurePSSM	Yes	(?)
SVMlight	support vector machine	classifyModelSVMLIGHT	No	(?)
DSSP	a database of secondary structure assignments for protein entries in the Protein Data Bank (PDB)	getDSSP	No	(?)
Proteus2	predict secondary structure	predictPROTEUS	No	(?)
HMMER	predict domains with hmmpfam using models of Pfam database	predictPFAM	No	(?)

Table 3: Summary Table for Homolog Reduction Functions.

Function	Description	Ref
hr	employ cdhitHR and aligndisHR to do homolog reduction	
cdhitHR	invoke cd-hit to cluster sequences	(?)
aligndisHR	calculated identity of aligned sequences	(?)

Table 4: Summary Table for Base and Amino Acid Groups.

Function	Description	Ref
elements	basic elements of biological sequence	
aaClass	amino acids groups depend on their physical-chemical properties: hydrophobicity, normalized Van der Waals volume, polarizability, polarity, and so on	(?)

Table 5: Summary table for Feature Coding Functions.

Type	Function	Feature Coding Scheme	Ref
DNA, RNA, or protein	featureBinary	use 0-1 vector to code each element	(??)
	featureCTD	numeric vector for the composition, transition and distribution of properties	(?)
	featureFragmentComposition	numeric vector for the frequency of k-mer sequence fragment	(??)
	featureGapPairComposition	numeric vector for the frequency of g-spaced element pair	(?)
	featureCKSAAP	integer vector for the number of k-spaced element pair (k cycled from 0 to g)	(?)
protein	featureHydro	hydrophobic effect	(??)
	featureACH	average cumulative hydrophobicity over a sliding window	(??)
	featureAAindex	numeric vector measuring the physicochemical and biochemical properties based on AAindex database	(??)
	featureACI	numeric vector measuring the average cumulative properties in AAindex	(??)
	featureACF	numeric vector measuring the Auto Correlation Function (ACF) of properties in AAindex	(??)
	featurePseudoAAComp	numeric vector for the pseudo amino acid composition proposed by Chou,K.C.	(?)
	featurePSSM	numeric vector for the normalized position-specific score of PSSM generated by PSI-BLAST	(????)
	featureDOMAIN	vector for the number of domain. Domains can be obtained by 'predictPFAM' function.	(?)
DNA or RNA	featureSSC	coding for secondary structure of protein. Secondary structure can be got by 'predictPROTEUS' or 'getDSSP'.	(?)
	featureBDNAVIDEO	Conformational and physicochemical DNA features from B-DNA-VIDEO database	(?)
	featureDIPRODB	conformational and thermodynamic dinucleotide properties from DiProDB database (http://diprodb.fli-leibniz.de)	(?)

Table 6: Summary Table for Feature Selection Functions.

Function	Description	Ref
selectWeka	feature selction by methods in WEKA	(?)
selectFFS	feature forword selction based on the performance of classification model	(?)
classify	build and test model with cross validation, also support feature selection by invoking WEKA	

Table 7: Summary Table for Classification Methods.

Function	Description	Depended R package
<code>classifyModelLIBSVM</code>	support vector machine by LIBSVM	<i>e1071</i>
<code>classifyModelSVMLIGHT</code>	support vector machine by SVM-light	<i>klaR</i>
<code>classifyModelNB</code>	naive bayes	<i>klaR</i>
<code>classifyModelRF</code>	random forest	<i>randomForest</i>
<code>classifyModelKNN</code>	k-nearest neighbor	<i>class</i>
<code>classifyModelTree</code>	tree model	<i>tree</i>
<code>classifyModelNNET</code>	neural net algorithm	<i>VR</i>
<code>classifyModelRPART</code>	recursive partitioning trees	<i>rpart</i>
<code>classifyModelCTREE</code>	conditional inference trees	<i>party</i>
<code>classifyModelCTREELIBSVM</code>	combine conditional inference trees and support vector machine	<i>party</i> , <i>e1071</i>
<code>classifyModelBAG</code>	bagging method	<i>ipred</i>

```
[1] "Negative Site: 215"
[1] "Negative Protein: 9"
[1] "Negative Site After Homolog Reduction and Random Selection: 16"
[1] "Negative Protein After Homolog Reduction and Random Selection: 6"
```

2. If the original data are non-redundant positive/negative peptides. We directly read the data into R and assign class labels for them.

```
> tmpDir=file.path(path.package('BioSeqClass'), 'example')
> tmpFile1=file.path(tmpDir, 'acetylation_K.pos40.pep')
> tmpFile2=file.path(tmpDir, 'acetylation_K.neg40.pep')
> posSeq=as.matrix(read.csv(tmpFile1,header=F,sep='\t',row.names=1))[,1]
> negSeq=as.matrix(read.csv(tmpFile2,header=F,sep='\t',row.names=1))[,1]
> seq=c(posSeq,negSeq)
> classLable=c(rep("+1",length(posSeq)),rep("-1",length(negSeq)) )
> length(seq)

[1] 1200
```

3. Once you have positive/negative datasets, you can code them to numeric vectors by functions listed in table ?? . Function `featureBinary` and `featureGapPairComposition` are taken as examples of different coding methods, which use binary 0-1 coding and the composition of gapped amino acid pair, respectively. Other functions can be used in the same way.


```
> # Use 0-1 binary coding.
> feature1 = featureBinary(seq,elements("aminoacid"))
> dim(feature1)
```

```
[1] 1200 300
```

```
> # Use the composition of paired amino acids.
> feature2 = featureGapPairComposition(seq,0,elements("aminoacid"))
> dim(feature2)
```

```
[1] 1200 400
```

4. `classify` is used to build classification model under cross validation. It also supports feature selection by invoking WEKA. Models built with selected features usually can obtain higher accuracy. In the following codes, two models are built by `classify`. The 1st classification model 'LIBSVM_CV5' is built by support vector machine with linear kernel and get an accuracy of 0.56 under 5-fold cross validation. The 2nd classification model 'FS_LIBSVM_CV5' is also built by support vector machine with linear kernel, but a feature selection method called "CfsSubsetEval" is used before building model. Thus the 2nd model using feature selection achieves an higher accuracy of 0.62 than the 1st model using all features.

```
> data = data.frame(feature1,classLable)
> # Use support vector machine and 5 fold cross validation to do classification.
> LIBSVM_CV5 = classify(data, classifyMethod='libsvm',
+                       cv=5, svm.kernel='linear',svm.scale=F)
> LIBSVM_CV5[["totalPerformance"]]
```

tp	tn	fp	fn	prc	sn
337.0000000	337.0000000	263.0000000	263.0000000	0.5599397	0.5616667
sp	acc	mcc	pc		
0.5616667	0.5616667	0.1243355	0.3918323		

```
> # Features selection is done by envoking "CfsSubsetEval" method in WEKA.
> FS_LIBSVM_CV5 = classify(data, classifyMethod='libsvm',
+                          cv=5, evaluator='CfsSubsetEval', search='BestFirst',
+                          svm.kernel='linear', svm.scale=F)
> FS_LIBSVM_CV5[["totalPerformance"]] ## Accuracy is increased by feature select
```

tp	tn	fp	fn	prc	sn
258.0000000	488.0000000	112.0000000	342.0000000	0.6924855	0.4300000
sp	acc	mcc	pc		
0.8133333	0.6216667	0.2621407	0.3651539		

```

> # Selected features:
> colnames(data)[FS_LIBSVM_CV5$features[[1]]]

[1] "BIN.1_A" "BIN.2_R" "BIN.3_K" "BIN.4_K" "BIN.5_K" "BIN.5_N"
[7] "BIN.6_N" "BIN.6_G" "BIN.6_T" "BIN.6_P" "BIN.7_P" "BIN.9_W"
[13] "BIN.9_H" "BIN.9_Y" "BIN.9_L" "BIN.10_Q" "BIN.11_W" "BIN.11_V"
[19] "BIN.12_K" "BIN.12_F" "BIN.13_K" "BIN.13_I" "BIN.13_M" "BIN.14_A"
[25] "BIN.15_K" "BIN.15_E" "BIN.15_A" "BIN.15_T"

```

5. Different feature coding methods usually might result in different prediction performance. `featureEvaluate` can be used to test multiple feature coding methods. Figure ?? shows the 3D plot of prediction accuracy varied with feature coding functions and parameters. It can be generated by employing `featureEvaluate` as follows (Note: It may be time consuming!):

```

> fileName = tempfile()
> # Note: It may be time consuming.
> testFeatureSet = featureEvaluate(seq, classLable, fileName, cv=5,
+   ele.type='aminoacid', featureMethod=c('Binary','GapPairComposition'),
+   classifyMethod='libsvm',
+   group=c('aaH', 'aaV', 'aaZ', 'aaP', 'aaF', 'aaS', 'aaE'), g=0,
+   hydro.methods=c('kpm', 'SARAH1'), hydro.indexs=c('hydroE', 'hydroF', 'hydroO')
> summary = read.csv(fileName,sep="\t",header=T)
> # Plot the result of 'featureEvaluate'
> require("scatterplot3d")
> tmp1 = summary[, "Feature_Function"]
> tmp2 = as.factor(sapply(as.vector(summary[, 'Feature_Parameter']),
+   function(x){unlist(strsplit(x,split='; '))[1]}))
> testResult = data.frame(as.integer(tmp2), as.integer(tmp1), summary[, "acc"])
> s3d=scatterplot3d(testResult,
+   color=c('red','blue')[testResult[,2]], pch=19, xlab='Parameter',
+   ylab='Feature Coding',
+   zlab='Accuracy', lab=c(9,3,7),
+   x.ticklabs=gsub('class: ','',sort(unique(tmp2))),
+   type='h',ylim=c(0,3),y.margin.add=2.5,
+   y.ticklabs=c('',gsub('feature','',sort(unique(tmp1))),''))

```

6. Features from multiple functions can be combined and re-selected to increase the prediction accuracy. In the following code chunk, the first three feature sets from 'testFeatureSet' are combined together ('testFeatureSet' is generated in the aforementioned codes by `featureEvaluate`). Then feature selection functions (`classify` and `selectFFS`) can be employed to select features. (`classify` has been illustrated in the aforementioned code chunk. Thus `selectFFS` is used here to do feature forward selection to

select a subset with maximum prediction accuracy (Note: It may be time consuming!). The process of feature selection and the increasing accuracy are shown in Figure ??.

```
> feature.index = 1:3
> tmp <- testFeatureSet[[1]]$data
> tmp1 <- testFeatureSet[[feature.index[1]]]$model
> colnames(tmp) <- paste(
+   tmp1["Feature_Function"],
+   tmp1["Feature_Parameter"],
+   colnames(tmp),sep=" ; ")
> data <- tmp[,-ncol(tmp)]
> for(i in 2:length(feature.index) ){
+   tmp <- testFeatureSet[[feature.index[i]]]$data
+   tmp1 <- testFeatureSet[[feature.index[i]]]$model
+   colnames(tmp) <- paste(
+     tmp1["Feature_Function"],
+     tmp1["Feature_Parameter"],
+     colnames(tmp),sep=" ; ")
+   data <- data.frame(data, tmp[,-ncol(tmp)] )
+ }
> name <- colnames(data)
> data <- data.frame(data, tmp[,ncol(tmp)] ) ## Combined features
> # Use 'selectFFS' to do feature forward selection.
> # Note: It may be time consuming.
> combineFeatureResult = selectFFS(data,accCutoff=0.005,
+   classifyMethod="knn",cv=5) ## It is time consuming.
> tmp = sapply(combineFeatureResult,function(x){
+   c(length(x$features),x$performance["acc"])}))
> plot(tmp[1,],tmp[2,],xlab="Feature Number",ylab="Accuracy",
+   , pch=19)
> lines(tmp[1,],tmp[2,])
```

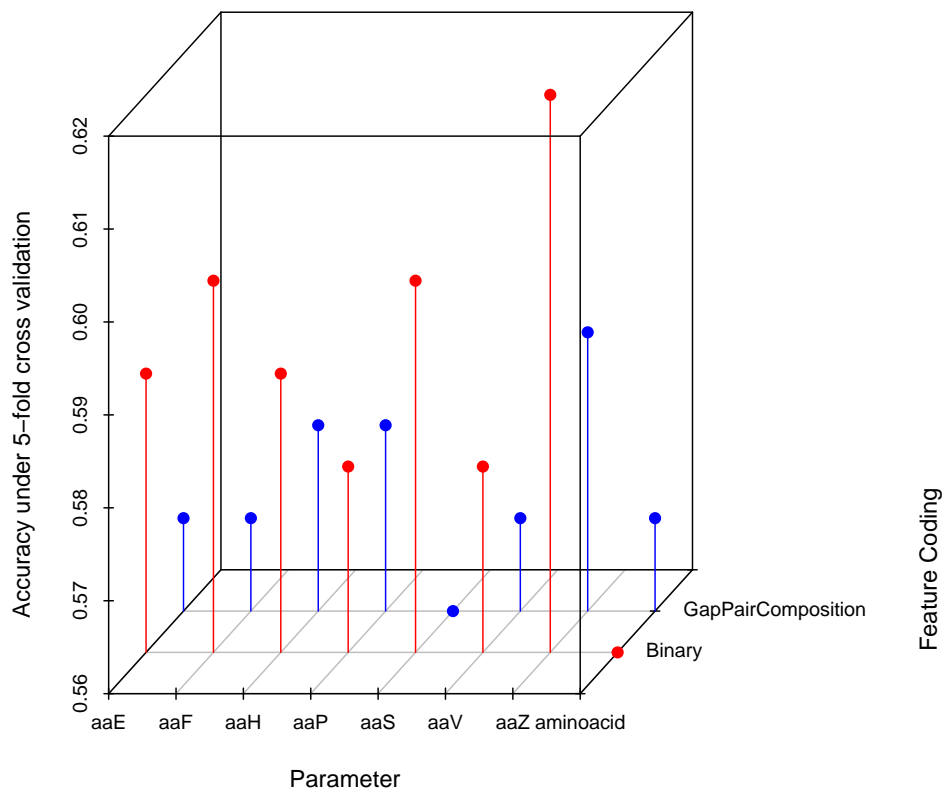


Figure 2: Result of `featureEvaluate`.

5 Session Information

This vignette was generated using the following package versions:

```
R version 3.3.0 RC (2016-04-26 r70550)
Platform: x86_64-apple-darwin13.4.0 (64-bit)
Running under: OS X 10.9.5 (Mavericks)
```

locale:

```
[1] C/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
```

attached base packages:

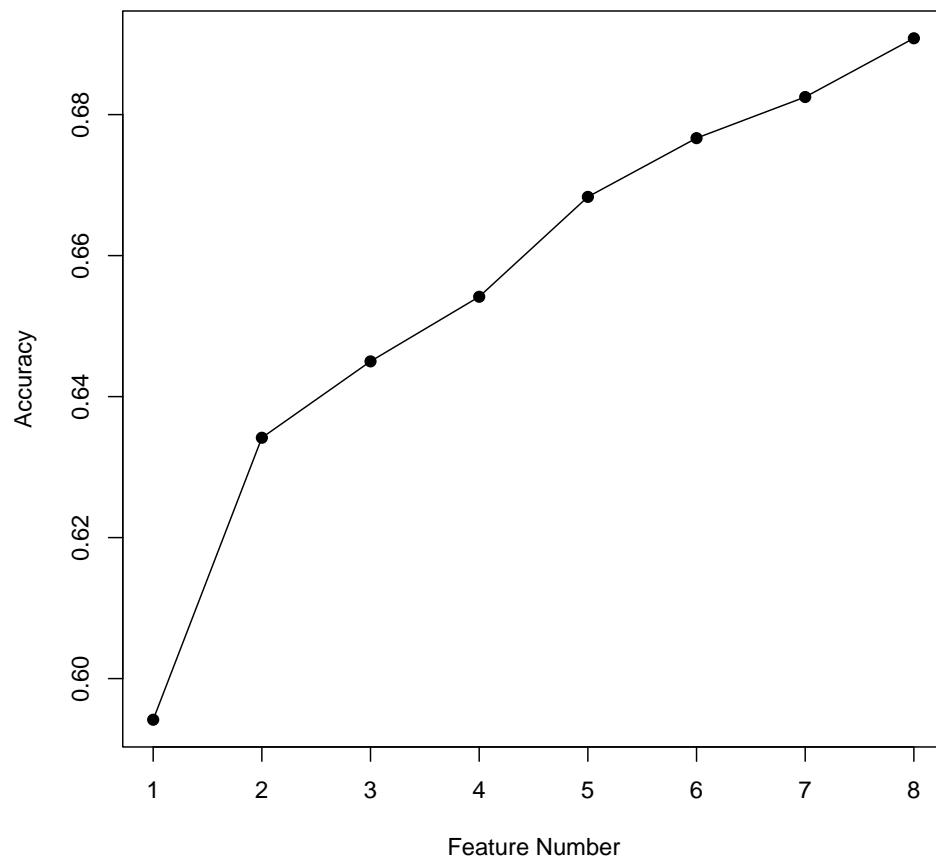


Figure 3: Result of `selectFFS`.

```
[1] stats      graphics  grDevices  utils      datasets  methods   base
```

other attached packages:

```
[1] BioSeqClass_1.30.0  scatterplot3d_0.3-36
```

loaded via a namespace (and not attached):

```
[1] Rcpp_0.12.4.5      XVector_0.12.0      zlibbioc_1.18.0
[4] IRanges_2.6.0      BiocGenerics_0.18.0 splines_3.3.0
[7] MASS_7.3-45        prodlim_1.5.7        lattice_0.20-33
[10] multcomp_1.4-4      tools_3.3.0          nnet_7.3-12
[13] parallel_3.3.0      grid_3.3.0           ipred_0.9-5
[16] Biobase_2.32.0      TH.data_1.0-7        e1071_1.6-7
[19] modeltools_0.2-21   class_7.3-14         randomForest_4.6-12
[22] survival_2.39-2     Matrix_1.2-6         lava_1.4.3
[25] party_1.0-25        S4Vectors_0.10.0     codetools_0.2-14
[28] rpart_4.1-10        strucchange_1.5-1    coin_1.1-2
[31] sandwich_2.3-4      klaR_0.6-12          Biostrings_2.40.0
[34] combinat_0.0-8      stats4_3.3.0         tree_1.0-37
[37] mvtnorm_1.0-5       foreign_0.8-66       zoo_1.7-12
```