

Package ‘gene2pathway’

May 10, 2009

Version 1.2.0

Date 2008-12-05

Title Prediction of KEGG pathway membership for individual genes based on InterPro domain signatures

Author Holger Froehlich <h.froehlich@dkfz.de>, contributions by Tim Beissbarth <t.beissbarth@dkfz.de>

Description The package takes a list of genes and predicts to which KEGG pathway each gene maps to. This is done by looking at the InterPro domains of each gene. Each prediction is assigned a confidence score. The package also allows to predict connected component membership of genes within signaling pathways. Separate models for each organism supported by KEGG can be trained.

Maintainer Holger Froehlich <h.froehlich@dkfz.de>

Reference H. Froehlich, M. Fellmann, H. Suelmann, A. Poustka, T. Beissbarth, Predicting Pathway Membership via Domain Signatures, *Bioinformatics*, 24:2137-2142, 2008

Depends R (>= 2.6.0), kernlab (>= 0.9), KEGG.db (>= 2.2.0), biomaRt (>= 1.12.1), KEGGSOAP (>= 1.12.0), RBGL, AnnotationDbi, org.Dm.eg.db

Imports SSOAP, RCurl

LazyLoad Yes

biocViews Microarray, Bioinformatics, Classification, GraphsAndNetworks, Pathways

License GPL (>= 2)

Repository CRAN

Date/Publication 2009-05-10 19:00:37

R topics documented:

classificationModel	2
classificationModelSignalTrans	3
gene2pathway	4
gene2pathway.signaltrans	5
getComponents	7
internal	8
ORF2Entrez	9
retrain	9
retrain.signaltrans	11
test.overrepresentation	12
Index	14

classificationModel

Hierarchical Classification Model

Description

This file contains the hierarchical classification model to predict KEGG pathway branches for genes. The model uses a pruned KEGG hierarchy, where metabolic pathways are not distinguished further, and the KEGG hierarchy for "cellular processes" and "genetic information processing" is pruned at the 2nd level. By default the model uses bagging to improve prediction accuracy. Important: There exists one separate model file for each organism.

Format

List of class "model", where each model has the following entries:

W learned decision hyperplane normal vector

C dictionary of label vectors, which can be predicted individually or which can be used to predict combinations of them

detectors SVM models trained to separate one specific pathway branch from the rest of the hierarchy

used_domains InterPro domains used by the classifier to separate the specific branch from the rest of the hierarchy

alldomains all InterPro domains used to build feature vectors

allpathways hierarchy branches, which can be predicted

treesizes relative size of hierarchy below the corresponding branch

kegg_hierarchy a nested list with information (parent branches, pathway names, pathway IDs, hierarchy level) on all higher hierarchy branches for each pathway

Author(s)

Holger Froehlich

See Also

[classificationModelSignalTrans](#)

classificationModelSignalTrans

Hierarchical Classification Model for Signaling Transduction Pathways and Pathway Components

Description

This file contains the hierarchical classification model to predict KEGG signaling pathways and pathway components for genes. The model contains only pathway components, to which a specified minimum number of genes could be mapped in the training phase (see [retrain.signaltrans](#)). Important: There exists one separate model file for each organism.

Format

List of class "model", where each model has the following entries:

W learned decision hyperplane normal vector

C dictionary of label vectors, which can be predicted individually or which can be used to predict combinations of them

detectors SVM models trained to separate one specific pathway branch from the rest of the hierarchy

used_domains InterPro domains used by the classifier to separate the specific branch from the rest of the hierarchy

alldomains all InterPro domains used to build feature vectors

allpathways hierarchy branches, which can be predicted

treесizes" relative size of hierarchy below the corresponding branch

kegg_hierarchy a nested list with information (parent branches, pathway names, pathway IDs, hierarchy level) on all higher hierarchy branches for each pathway

elemIDs a list of KEGG element IDs mapping to each pathway component - may be used to highlight pathway components with [color.pathway.by.elements](#).

Author(s)

Holger Froehlich

See Also

[classificationModel](#)

gene2pathway

*Pathway membership prediction***Description**

Predicts a gene's membership to a branch in the KEGG hierarchy via the contained InterPro domains.

Usage

```
gene2pathway(geneIDs=NULL, flyBase=FALSE, gene2Domains=NULL, organism="hsa", useKEGG=FALSE)
```

Arguments

geneIDs	a character vector of Entrez gene IDs or FlyBase identifiers (not necessary, if the argument gene2Domains is provided)
flyBase	Are FlyBase identifiers provided? Default: No
gene2Domains	By default associations between genes and InterPro domains are retrieved via biomaRt from Ensembl. Alternatively, the user can provide its own mapping of genes to InterPro domains in form of a list here (see details).
organism	KEGG letter code describing an organism. Please refer to <URL: http://www.genome.jp/kegg-bin/create_kegg_menu > for a complete list of organisms (and their letter codes) supported by KEGG.
useKEGG	Should KEGG information instead of a prediction be used when possible?
KEGG.package	If useKEGG=TRUE: Instead of retrieving information directly from KEGG, one can use the KEGG.db package instead, which is significantly faster. However, the KEGG.db package only supports a fraction of organisms so far. Please refer to the manual pages of the KEGG.db package for further information. Default: use KEGG.db package, if useKEGG=TRUE

Details

A hierarchical classification model based on SVMs and a ranking perceptron is used. This model is usually additionally bagged to improve prediction quality. The model is stored in the package data directory and is recommended to be retrained from time to time.

The current version of the KEGG hierarchy is always retrieved directly from KEGG via FTP. By default associations between genes and InterPro domains are retrieved automatically via biomaRt from Ensembl. Please refer to <URL:<http://www.ebi.ac.uk/ensembl/>> for a list of organisms supported by Ensembl. Alternatively to using Ensembl and biomaRt, the user can provide its own mapping of genes to InterPro domains in form of a list. This especially allows for using organisms, which are supported by KEGG, but not by Ensembl so far. The list has the form genes -> InterPro domains, and each list entry is named by a gene identifier of the corresponding gene. If useKEGG=TRUE, Entrez gene IDs or FlyBase identifiers have to be used. Otherwise, arbitrary identifiers are allowed.

Value

gene2Path	mapping of gene IDs to corresponding KEGG pathway IDs
gene2Pathname	mapping of gene IDs to corresponding KEGG pathway names
byKEGG	inticates by TRUE/FALSE for each gene whether the mapping information was obtained directly from KEGG or whether it was predicted
scores	confidence scores for the prediction (0, if no prediction was performed): see notes for details

Note

By default a bagged model prediction is used, i.e. each of the individual sub-models is giving a vote for a specific output. The final output is determined by the majority of the votes for each hierarchy branch separately. The corresponding fraction voting for a specific branch may be interpreted as its probability. In the ideal case all individual branch probabilities should always be close to 1, if the gene maps to that part of the KEGG hierarchy, and close to 0 otherwise. A cumulative measure of confidence is thus the average over all probabilities > 0.5 and one minus the average over all probabilities < 0.5 . We combine both measure by taking the average of both and report it as a reliability score.

If the user decides to retrain a model WITHOUT using bagging, then the reliability score is simply the margin between the highest and the second highest ranked solution. This margin should be larger 2 for good confidence.

Author(s)

Holger Froehlich

See Also

[retrain, classificationModel](#)

Examples

```
## Not run:  
gene2pathway("10327")  
## End(Not run)
```

gene2pathway.signaltrans

Pathway membership prediction

Description

Predicts a gene's membership to a KEGG signaling pathway and/or pathway component via the contained InterPro domains.

Usage

```
gene2pathway.signaltrans(geneIDs=NULL, flyBase=FALSE, gene2Domains=NULL, organism='')
```

Arguments

geneIDs	a character vector of Entrez gene IDs or FlyBase identifiers (not necessary, if the argument gene2Domains is provided)
flyBase	Are FlyBase identifiers provided? Default: No
gene2Domains	By default associations between genes and InterPro domains are retrieved via biomaRt from Ensembl. Alternatively, the user can provide its own mapping of genes to InterPro domains in form of a list here (see details).
organism	KEGG letter code describing an organism. Please refer to <URL: http://www.genome.jp/kegg-bin/create_kegg_menu > for a complete list of organisms (and their letter codes) supported by KEGG.
useKEGG	Should KEGG information instead of a prediction be used when possible?

Details

A hierarchical classification model based on SVMs and a ranking perceptron is used. This model is usually additionally bagged to improve prediction quality. The model is stored in the package data directory and is recommended to be retrained from time to time.

The current version of the KEGG hierarchy is always retrieved directly from KEGG via FTP. By default associations between genes and InterPro domains are retrieved automatically via biomaRt from Ensembl. Please refer to <URL:<http://www.ebi.ac.uk/ensembl/>> for a list of organisms supported by Ensembl. Alternatively to using Ensembl and biomaRt, the user can provide its own mapping of genes to InterPro domains in form of a list. This especially allows for using organisms, which are supported by KEGG, but not by Ensembl so far. The list has the form genes -> InterPro domains, and each list entry is named by a gene identifier of the corresponding gene. If useKEGG=TRUE, Entrez gene IDs or FlyBase identifiers have to be used. Otherwise, arbitrary identifiers can be allowed.

Value

gene2Path	mapping of gene IDs to corresponding KEGG pathway IDs
gene2Pathname	mapping of gene IDs to corresponding KEGG pathway names
byKEGG	inticates by TRUE/FALSE for each gene whether the mapping information was obtained directly from KEGG or whether it was predicted
scores	confidence scores for the prediction (0, if no prediction was performed): see notes for details
elemIDs	KEGG elements mapping to the corresponding predicted pathway components, if there are any, otherwise NULL. May be used to highlight pathway components with <code>color.pathway.by.elements</code> .

Note

By default a bagged model prediction is used, i.e. each of the individual sub-models is giving a vote for a specific output. The final output is determined by the majority of the votes for each hierarchy branch separately. The corresponding fraction voting for a specific branch may be interpreted as its probability. In the ideal case all individual branch probabilities should always be close to 1, if the gene maps to that part of the KEGG hierarchy, and close to 0 otherwise. A cumulative measure of confidence is thus the average over all probabilities > 0.5 and one minus the average over all probabilities < 0.5 . We combine both measure by taking the average of both and report it as a reliability score.

If the user decides to retrain a model WITHOUT using bagging, then the reliability score is simply the margin between the highest and the second highest ranked solution. This margin should be larger 2 for good confidence.

Author(s)

Holger Froehlich

See Also

[retrain.signaltrans](#), [classificationModelSignalTrans](#)

Examples

```
## Not run:
  gene2pathway.signaltrans("1443")
## End(Not run)
```

getComponents

KEGG pathway information

Description

1. get connected pathway components; 2. get all elements of a given pathway; 3. color certain elements in a pathway.

Usage

```
getComponents(pathway.id)

get.elements.by.pathway(pathway.id)

color.pathway.by.elements(pathway.id, elements)
```

Arguments

```
pathway.id  KEGG pathway ID, e.g. "path:hsa04012"
elements    KEGG element IDs: character vector of numbers
```

Details

All functions use the KEGG SOAP service.

Value

getComponents: a list with the entries

geneIDs Entrez gene IDs mapping to each pathway component

elemIDs KEGG element IDs mapping to each pathway component

get.elements.by.pathway: list, see <URL http://www.genome.jp/kegg/soap/doc/keggapi_manual.html>
for details

color.pathway.by.elements: an URL of a colored gif file, see <URL http://www.genome.jp/kegg/soap/doc/keggapi_manual.html>
for details

Author(s)

Holger Froehlich

Examples

```
## Not run:  
  comp = getComponents("path:hsa04020") # get all connected components  
  color.pathway.by.elements("path:hsa04020", comp$elemIDs[[1]]) # mark first component  
## End(Not run)
```

internal

internal functions

Description

internal functions: do not call these functions directly.

Usage

various

Arguments

various

Value

various

Author(s)

Holger Froehlich

ORF2Entrez	<i>Conversion locus tag -> Entrez ID and Entrez ID -> locus tag according to KEGG</i>
------------	---

Description

Conversion locus tag -> Entrez ID and Entrez ID -> locus tag according to KEGG (see <URL:ftp://ftp.genome.jp/pub/kegg/ge

Usage

```
ORF2Entrez (ORFIDs, organism="dme")
Entrez2ORF (entrezIDs, organism="dme")
```

Arguments

ORFIDs	character vector of locus tags / ORF identifiers
organism	KEGG letter code describing an organism. Please refer to <URL:http://www.genome.jp/kegg-bin/create_kegg_menu> for a complete list of organisms (and their letter codes) supported by KEGG.
entrezIDs	character vector of Entrez gene IDs

Value

character vector

Examples

```
entrez=ORF2Entrez ("Dmel_CG4942", organism="dme")
Entrez2ORF (entrez, organism="dme") # should be "Dmel_CG4942"
```

retrain	<i>Retrain classification model</i>
---------	-------------------------------------

Description

Retrains the hierarchical classification model. This way new information from InterPro and KEGG databases can be incorporated to give better predictions. Retraining should be done on a regular basis from time to time.

Usage

```
retrain(minnmap=30, level1Only="Metabolism", level2Only="Genetic Information Proces
```

Arguments

<code>minnmap</code>	prune hierarchy branches with < minnmap mapping genes
<code>level1Only</code>	for these hierarchy branches only the first level is used
<code>level2Only</code>	for these hierarchy branches only the first and the second levels are used
<code>organism</code>	KEGG letter code describing an organism. Please refer to <URL:http://www.genome.jp/kegg-bin/create_kegg_menu> for a complete list of organisms (and their letter codes) supported by KEGG.
<code>gene2Domains</code>	By default associations between genes and InterPro domains are retrieved via biomaRt from Ensembl. Alternatively, the user can provide its own mapping of genes to InterPro domains in form of a list here (see details).
<code>KEGG.package</code>	Instead of retrieving information directly from KEGG, one can use the KEGG.db package instead, which is significantly faster. However, the KEGG.db package only supports a fraction of organisms so far. Please refer to the manual pages of the KEGG.db package for further information. Default: Don't use KEGG.db package
<code>remove.duplicates</code>	remove genes having the same InterPro domains prior training. Default: Don't do this
<code>use.bagging</code>	use bagging
<code>nbag</code>	number of models to average over

Details

A hierarchical classification model based on SVMs and a ranking perceptron algorithm is trained. This model is usually additionally bagged to improve prediction quality. The method produces a "classificationModel_[organism].rda" (e.g. "classificationModel_hsa.rda") file, which should be stored in the package data directory. Once a new model has been trained, the complete package should be reloaded.

The current version of the KEGG hierarchy is always retrieved directly from KEGG via FTP. By default associations between genes and InterPro domains are retrieved automatically via biomaRt from Ensembl. Please refer to <URL:http://www.ebi.ac.uk/ensembl/> for a list of organisms supported by Ensembl. Alternatively to using Ensembl and biomaRt, the user can provide its own mapping of genes to InterPro domains in form of a list. This especially allows for using organisms, which are supported by KEGG, but not by Ensembl so far. The list has the form genes -> InterPro domains, and each list entry is named by the Entrez gene ID of the corresponding gene. This is, because KEGG uses Entrez gene IDs for the mapping genes -> KEGG pathways.

Value

The model structure. See `classificationModel` for details.

Author(s)

Holger Froehlich

See Also

[gene2pathway](#), [classificationModel](#)

Examples

```
## Not run:
      retrain(KEGG.package=TRUE, organism="dme") # retrain classification model for drosophila
## End(Not run)
```

```
retrain.signaltrans
```

Retrain classification model for signaling pathways

Description

Retrains the hierarchical classification model for signaling pathway components. This way new information from InterPro and KEGG databases can be incorporated to give better predictions. Retraining should be done on a regular basis from time to time.

Usage

```
retrain.signaltrans(minnmap=10, organism="hsa", gene2Domains=NULL, remove.duplicates=FALSE)
```

Arguments

<code>minnmap</code>	prune hierarchy branches with < minnmap mapping genes
<code>organism</code>	KEGG letter code describing an organism. Please refer to <URL: http://www.genome.jp/kegg-bin/create_kegg_menu > for a complete list of organisms (and their letter codes) supported by KEGG.
<code>gene2Domains</code>	By default associations between genes and InterPro domains are retrieved via biomaRt from Ensembl. Alternatively, the user can provide its own mapping of genes to InterPro domains in form of a list here (see details).
<code>remove.duplicates</code>	remove genes having the same InterPro domains prior training
<code>use.bagging</code>	use bagging
<code>nbag</code>	number of models to average over

Details

A hierarchical classification model based on SVMs and a ranking perceptron algorithm is trained. This model is usually additionally bagged to improve prediction quality. The method produces a "classificationModelSignalTrans_[organism].rda" (e.g. "classificationModelSignalTrans_hsa.rda") file, which should be stored in the package data directory. Once a new model has been trained, the complete package should be reloaded.

The current version of the KEGG hierarchy is always retrieved directly from KEGG via FTP. Labels for the training set are obtained via the function [getComponents](#), which uses the KEGG SOAP

service. By default associations between genes and InterPro domains are retrieved automatically via biomaRt from Ensembl. Please refer to <URL: <http://www.ebi.ac.uk/ensembl/>> for a list of organisms supported by Ensembl. Alternatively to using Ensembl and biomaRt, the user can provide its own mapping of genes to InterPro domains in form of a list. This especially allows for using organisms, which are supported by KEGG, but not by Ensembl so far. The list has the form genes -> InterPro domains, and each list entry is named by the Entrez gene ID of the corresponding gene. This is, because KEGG uses Entrez gene IDs for the mapping genes -> KEGG pathways.

Value

The model structure. See `classificationModelSignalTrans` for details.

Author(s)

Holger Froehlich

See Also

`gene2pathway.signaltrans`, `classificationModelSignalTrans`

Examples

```
## Not run:
      retrain.signaltrans() # retrain classification model for signal transduction pathway
## End (Not run)
```

```
test.overrepresentation
```

Test statistical overrepresentation of KEGG pathways in a list of genes

Description

Test the statistical overrepresentation of KEGG pathways in a group of genes using Fisher's exact test. The analysis can either be based on all KEGG pathways predicted by `gene2pathway/gene2pathway.signaltrans` or on original KEGG annotation only.

Usage

```
test.overrepresentation(genesOfInterest, predpath, KEGGonly=FALSE, cutoff=0.1, min.
```

Arguments

<code>genesOfInterest</code>	a character vector of gene identifiers (see <code>gene2pathway</code> , <code>gene2pathway.signaltrans</code>) for a gene list of interest
<code>predpath</code>	predictions of <code>gene2pathway</code> or <code>gene2pathway.signaltrans</code>
<code>KEGGonly</code>	use KEGG annotation only

<code>cutoff</code>	p-value significance cutoff
<code>min.conf</code>	filter predictions such that only those with a confidence score > min.conf are considered
<code>adj.method</code>	multiple testing correction method. Default: Benjamini-Yekutieli

Value

Table with two columns: KEGG pathway and adjusted p-value (adjustment according to Benjamini-Yekutieli)

Index

*Topic **datasets**

classificationModel, 2
classificationModelSignalTrans,
3

*Topic **file**

gene2pathway, 4
gene2pathway.signaltrans, 5
getComponents, 7
internal, 8
ORF2Entrez, 9
retrain, 9
retrain.signaltrans, 11
test.overrepresentation, 12

buildTrainingSet (*internal*), 8

classificationModel, 2, 3, 5, 10, 11
classificationModelSignalTrans,
2, 3, 7, 12

code_test (*internal*), 8
code_train (*internal*), 8
color.pathway.by.elements, 3, 6
color.pathway.by.elements
(*getComponents*), 7

Entrez2ORF (*ORF2Entrez*), 9
Entrez2ORF.*internal* (*internal*), 8

gene2pathway, 4, 11, 12
gene2pathway.signaltrans, 5, 12
get.element.relations.by.pathway
(*internal*), 8
get.elements.by.pathway
(*getComponents*), 7
getComponents, 7, 11
getComponents.*internal*
(*internal*), 8
getInterProDomains (*internal*), 8
getKEGGHierarchy (*internal*), 8

internal, 8

KEGG2Entrez (*internal*), 8

loss (*internal*), 8

modelKEGG (*classificationModel*), 2
modelSignalTrans
(*classificationModelSignalTrans*),
3

ORF2Entrez, 9

predict.gene2pathway (*internal*), 8

retrain, 5, 9

retrain.signaltrans, 3, 7, 11

struct_predict (*internal*), 8
struct_train (*internal*), 8
svmlearn (*internal*), 8
svmpredict (*internal*), 8

test.overrepresentation, 12