

Package ‘autoGLM’

July 6, 2016

Type Package

Title Automated fitting, model selection, and documentation functions for GLM modelling of binary response variables.

Version 1.0.0

Author Bo Pieter Johannes Andree

Maintainer <b.p.j.andree@vu.nl>

Description Methods and tools for automation of binary response GLM's. Supports linear models with identity, probit and logit link. The package is written with binomial land-use modelling using large data sets in mind, but the methods and tools stretch out to other research cases. The main function is autoGLM, which wraps around most of the other functions contained in the package. However, some general functions are included as well such as a more elaborated function for package loading, an out of memory version to lapply, memory-threshold garbage collection, sample size reduction routines, variable selection routines, and user friendly methods for logit estimation through GMM.

Depends R (>= 3.2.1)

URL <https://github.com/BPJandree/AutoGLM>

BugReports <https://github.com/BPJandree/AutoGLM>

License MIT

ByteCompile TRUE

LazyData TRUE

RoxygenNote 5.0.1

R topics documented:

accuracy	2
autoGLM	3
corinetable	5
describe	5
drop.na.model	6
easygmmlogit	6
exportWeightsfile	7
fitnew	8
generalizeToSpecific	9
getCall	10
getSamples	11

guessStartVal	12
iapply	13
IC	14
is.speed	15
ITdata	15
logistic	16
MLtoBinomData	16
normalize	17
oftype	17
opt.glm	18
opt.glm.core	19
opt.h	20
opt.ic	20
opt.t	21
pkgTest	22
PtoBin	22
reclassify	23
selectX	24
simulateLogit	25
tgc	26
update.start	27
update.X	27
update.Z	28
warmstart	28
withinInterval	29
Index	30

accuracy	<i>Miscellaneous function used in main routines.</i>
----------	--

Description

Miscellaneous function used in main routines.

Usage

```
accuracy(testdata, fitteddata)
```

Arguments

testdata	a vector of occurrences.
fitteddata	a vector of occurrences.

Value

A numeric indicating the degree of correspondence.

Examples

```
observed <- c(1,1,1,0,0,0,0,0,0,0)
predicted <-c(1,0,1,0,0,0,0,1,1,0)

print(accuracy(observed, predicted))
```

autoGLM	<i>Main function of the package.</i>
---------	--------------------------------------

Description

This function is a wrapper around the optimization and selection routines in the package and can be used for automated calibration of GLM's on semi large datasets.

Usage

```
autoGLM(data, reclassable = "default", class = 1, outputpath = wd(),
  modelname = "autoGLM", tracelevel = 1, actions = c("write", "print",
    "log", "return"), Nval = "default", model = "logit", preselect = "lm",
    method = "opt.ic", crit.t = 1.64, crit.p = 0.1, test = "LR",
    KLIC = "AICc", accuracytolerance = 0.01, confidence.alternative = 0.85,
    use.share = 0.25, maxsampleruns = 50, memorymanagement = TRUE,
    returnall = FALSE, compress = FALSE, JIT = TRUE)
```

Arguments

data	A dataframe with a categorical response variable in the first column, and covariates in subsequent columns. Typically the product of <code>cbind(Y,X)</code> .
reclassable	A table that maps the first column of data into a binary response variable. By default it will be omitted (the binary response variable will be identical to <code>data[,1]</code>). See also corinetable , See also reclassify .
class	The class that should be 1 in the binary response variable, all other classes in the categorical variable will be set to 0. Defaults to 1. See also reclassify .
outputpath	The location on the hard drive where output will be written to. Defaults to <code>wd()</code> .
modelname	The name of the model, will be used when writing a weightsfile. Defaults to "autoGLM". See also exportWeightsfile .
tracelevel	The amount of information to be printed. Passed on to underlying routines. Defaults to 1 for printing, set to 0 for no printing.
actions	Actions to be taken by autoGLM, may include any combination of <code>c("write", "print", "log", "return")</code> , for writing a geoDMS weightsfile, See also exportWeightsfile , printing results, writing a log file, and returning results as a list object.
Nval	Optional categorical variable that should be dropped by the reclassification scheme. See also reclassify .
model	Main model type that should be calibrated, either "lm", "probit", or "logit". See also generalizeToSpecific .
preselect	Optional variable preselection using a first order approximation (linear model) of the logit or probit model, by specifying "lm" (default setting). See also selectX .

method	The optimization strategy. Either "opt.ic" to optimize using information criteria, "opt.t" for step-wise elimination of insignificant values (statistically speaking not a sound procedure, but it will provide a parsimonious model that can be usefull as a benchmark), or "opt.h" to optimize by classical hypothesis tests. defaults to "opt.ic". See also opt.ic , opt.t , See also opt.h .
crit.t	The t-value indicating significance when using method "opt.t", defaults to 1.64. opt.t .
crit.p	the p-value used by method "opt.h" in the hypothesis tests. Defaults to 0.05. opt.h .
test	The hypothesis test used by "opt.h". Defaults to "LR" for the Likelihood Ratio test. Other options are "F", for an F test for joint significance of insignificant parameters, or "Chisq" for a wald test against the Chi squared distribution. opt.h .
KLIC	The information criterion used by "opt.ic", either "AIC" or "AICc", defaults to the latter. opt.ic .
accuracytolerance	When aut of sample and within sample accuracy differ more than accuracytolerance, a warning will be issued, which is also logged when specifying "log" in actions. Defaults to 0.01. accuracy .
confidence.alternative	See also getSamples , confidence level used for the alternative of dissimilar samples in the sampling routine. Defaults to .85.
use.share	Share of the data used, See also getSamples . Defaults to .25.
maxsampleruns	See also getSamples , defaults to 50.
memorymanagement	TRUE/FALSE indicating whether garbage collection should be forced regularly when memory usage is high. Defaults to TRUE, recommended setting for large datasets. See also tgc .
returnall	TRUE, FALSE, or "writedisk" indicating whether all the outputted objects for each class should be returned in an array as produced by lapply, or whether only the final output should be returnd as an object. Specifying "writedisk" will write the objects containing results of each class as seperate .RDS files, which you can use to restore the output using readRDS(). iapply . Returning an array of all results can consume large amounts of memory as each object contains copies of the used datasets. When working with countrysize datasets, these array objects can easily require over 64gb of RAM. Specifying returnall = FALSE (default setting), is much more more RAM friendly as it stores results for each class in the same memory adress, overwriting previous results. Seting returnall = FALSE, will still write log files and print diagnostics to screen if specified in actions. returnall="writedisk" is the recommended setting, but it is not default. iapply .
compress,	passed on to iapply. Defaults to no compression of RDS output, which is the recommended setting if computation time is valued of disk space. Keep in mind that when using large datasets, autoGLM objects can be several gigabytes in size. iapply .
JIT,	logical indicating whether just-in-time compilation of internal functions should be used. Mainly for historical reasons.

Examples

```
load(ITdata)
```

```
load(corinetable)
actions <- c("write", "print", "log", "return")
results <- autoGLM(data=ITdata, reclass=corinetable, class=0,
outputpath=wd(), modelname="IT", tracelevel=1, actions=actions,
NAval="default", model="logit", preselect="lm",
method="opt.ic", KLIC="AICc", accuracytolerance=0.01,
confidence.alternative=0.85, use.share=0.25,
maxsampleruns=50, memorymanagement=TRUE)
```

corinetable	<i>A reclassification table for Corine land cover to simplified land cover used under the LUISA platform of the JRC.</i>
-------------	--

Description

A table that maps 51 Corine Land-cover classes into 8 classes and a NoData class.

- 0. Urban.
- 1. Industry.
- 2. Arable.
- 3. Permanent crops.
- 4. Pastures.
- 5. Forest.
- 6. Herbaceous cover.
- 16. Transitional.
- -9999. NoData.

Usage

```
data(corinetable)
```

describe	<i>A simple function to describe a dataset.</i>
----------	---

Description

This function returns a dataframe with variable names, their minimum and maximum values, means and standard deviations.

Usage

```
describe(X)
```

Arguments

X	Matrix or dataframe to be summarized.
---	---------------------------------------

Value

A dataframe with names and statistics.

Examples

```
someVector = 1:10
print(describe(someVector))

df <- data(ITdata)
description <- describe(df)
print(description)
```

drop.na.model	<i>Miscellaneous function used in main routines.</i>
---------------	--

Description

If a model object produced by `lm()` or `glm()` has NA standard errors, `drop.na.model()` drops the problematic variables and refits the model using a warmstart, See also [warmstart](#).

Usage

```
drop.na.model(namodel)
```

Arguments

namodel modelobject.

Value

a model object of the same type as the supplied object.

easygmmlogit	<i>A function for user-friendly estimation of logit models using GMM. Advanced users may use the more elaborate gmm package instead.</i>
--------------	--

Description

This function automizes GMM estimation of the logit model. GMM estimation runs very heavy compared to the IWLS estimator used by the standard routines used in this package, but GMM is a more general approach and works under milder conditions regarding the uniqueness of a solution.

Usage

```
easygmmlogit(Y, X, start = "default", maximizer = "nlminb")
```

Arguments

Y	A binary response variable.
X	A dataframe of multiple exogenous regressors.
start	A vector of starting values to be used by the numerical optimizer. By default, the function calls guessStartVal get a fast initial estimate using an IWLS estimator. If that fails, for example due to bad data input, starting values will be set to zero.
maximizer	The numerical procedure to be used. Either "BFGS" for fast optimization, or "nlnminb" for reliable PORT routines. Defaults to "nlnminb".

Value

A model generated by gmm()

Examples

```
# Load data
data(ITdata)
data(corinetable)
#Grab a sample, gmm runs heavy.
sample <- ITdata[getSamples(ITdata, share =.05),]
# Reclassify
catITdata <- reclassify(sample, reclassstable = corinetable)
# create a binary response dataset.
Y <- MLtoBinomData(catITdata[,1], class =0)
X <- catITdata[,c(2:4)]
GMM <- easygmmlogit (Y, X)
```

exportWeightsfile	<i>A function to export a csv file containing the weights file in a format that can be loaded into geoDMS.</i>
-------------------	--

Description

This function exports a csv file containing the names of the weights, coefficients and country name in a standardized format supported by geoDMS.

Usage

```
exportWeightsfile(model, originaldata, modeldata, coefnamelist, outdir,
  modelname, filename)
```

Arguments

model	The model object from which the estimated coefficients should be extracted.
originaldata	The complete dataset from which the model data has been extracted. Used by the function to determine which variables have a zero weight.
modeldata	The dataset on which the model is estimated. Used by the function to determine which variables have a zero weight.
coefnamelist	A character vector containing the weight names that should be printed in the weights file.

outdir	The directory to which the weightsfile should be exported.
modelname	A string containing the name of the country that should be printed in the weightsfile.
filename	The name of the exported weightsfile.

Value

Any messages that may be printed by file(), writeLines() or close().

Examples

```
logitpars <- c(-0.5, 0.4, -0.3, 0.2, -0.1, 0, 0, 0, 0, 0, 0, 0)
someData <- simulateLogit(nobs=2500, nxregs=12, pars=logitpars)
useData <- someData[,c(1:6)]
someModel <- glm(formula(useData), family=binomial(link='logit'), data=useData)

exportWeightsfile(model = someModel, originaldata = someData, modeldata = useData,
  coefnamelist = colnames(someData[-1,]), outdir = "C:\\Users\\",
  modelname = "myfirstglm", filename = "important_result.csv")

ITdata <- data(ITdata)
ITsample <- getSamples(data = ITdata, share = 0.05, confidence.alternative=0.90, max.iter =100)
sampledIT = ITdata[ITsample,]
reclIT <- reclassify(sampledIT)
urbanIT <- MLtoBinomData(reclIT, class=0)
trainY = urbanIT[,1]
trainX = urbanIT[,-1]
bestX <- selectX(trainY, trainX)
bestlogit <- glm(formula(cbind(trainY,trainX)), family=binomial(link='logit'), data=cbind(trainY,trainX))

exportWeightsfile(model = bestlogit, originaldata = ITdata, modeldata = cbind(trainY,trainX),
  coefnamelist = colnames(trainX), outdir = "C:\\Users\\",
  modelname = "IT", filename = "Urban_weights.csv")
```

fitnew

Miscellaneous function used in main routines.

Description

Miscellaneous function used in main routines.

Usage

```
fitnew(newspec, model, newZ, newstart)
```

Arguments

newspec	formula object
model	Either "lm", "logit", or "probit".
newZ	dataframe.
newstart	vector of starting values.

Value

model object of requested type

generalizeToSpecific	<i>A function to select from an entire set of exogenous regressors, those regressors that minimize Kullback-Leibler divergence in the supplied model. The underlying function uses a "warm start" algorithm initialized by guessStartval, in a fashion similar to the method of sieves. If you plan to use a Lasso estimator on large data, you might consider using this function to determine the zero parameters with option method = "lm", and continue with a Ridge estimator for computational reasons.</i>
----------------------	---

Description

This function is a wrapper to the functions `opt.ic()`, `opt.t()` and `opt.h()`.

Usage

```
generalizeToSpecific(model = "lm", Y, X, method = "opt.ic", KLIC = "AICc",
  crit.t = 1.64, crit.p = 0.05, test = "LR", tracelevel = 1,
  memorymanagement = TRUE)
```

Arguments

model	Either "lm" for the linear probability model, "logit" for the logistic probability model, or "probit", for the probit model. The logit and probit models are solved using Iterated Weighted Least Squares, and optimization of the logit model is significantly faster than the probit model. Defaults to "lm".
Y	A binary response variable.
X	A dataframe of multiple exogenous regressors.
method	The optimization strategy. Either "opt.ic" to optimize using information criteria, "opt.t" for step-wise elimination of insignificant values (statistically speaking not a sound procedure, but it will provide a parsimonious model that can be usefull as a benchmark), or "opt.h" to optimize by classical hypothesis tests. defaults to "opt.ic".
KLIC	the information criterion used by "opt.ic", either "AIC" or "AICc", defaults to the latter.
crit.t	The t-value indicating significance when using method "opt.t", defaults to 1.64.
crit.p	the p-value used by method "opt.h" in the hypothesis tests. Defaults to 0.05.
test	The hypothesis test used by "opt.h". Defaults to "LR" for the Likelihood Ratio test. Other options are "F", for an F test for joint significance of insignificant parameters, or "Chisq" for a wald test against the Chi squared distribution. Recommended setting is either "LR" as it is less dependent on correct estimation of the standard errors. Keep in mind that "Chisq" is an asymptotic test, anf "F" is more appropriate for small sample tests. However "Chisq" holds under milder conditions and should be used if no small sample theory is available for the model.

tracelevel	the amount of information to be printed. Passed on to underlying routines. Defaults to 1 for printing, set to 0 for no printing.
memorymanagement	TRUE/FALSE indicating whether garbage collection should be forec regularly when memory usage is high. Defaults to TRUE, recommended setting for large datasets.
share	between 0-1, specifying the amount of data that should be passed on to the optimization strategies. Defaults to 0.75, to improve speed. Uses getSamples() to maintain first and second moments of the data.

Value

Either a dataframe of exogenous variables, or a vector containing the collumn names indicating the optimal variables extracted from the supplied dataset.

Examples

```
data(ITdata)
data(corinetable)
#Grab a sample (optional).
sample <- ITdata[getSamples(ITdata, share =.05),]
# Reclassify
catITdata <- reclassify(sample, reclassstable = corinetable)
# create a binary response dataset.
Y <- MLtoBinomData(catITdata[,1], class =1)
X <- catITdata[,-1]
bestm <- generalizeToSpecific(model ="lm", Y, X)
notabadm <- generalizeToSpecific(model ="lm", Y, X, "opt.t")
```

getCall	<i>Returns the name of a call as string.</i>
---------	--

Description

When called inside a function, returns the parent function's name.

Usage

```
getCall(level = -1)
```

Arguments

level	0 returns "getCall", -1 returns the name of the function in which getCall(-1) is called etc.
-------	--

Value

string

getSamples

A function to extract small samples that maintain important characteristics of the population sample.

Description

This function returns a sample extracted from the supplied population data, that has a similar distribution to the supplied population dataset. The function is called by guessStartval() to estimate initial values for numerical optimization procedures, but can also be used directly to reduce the sample size such that computationally intensive models can be estimated on a representative sample of an entire dataset. The function makes use of var.test() to compute an F test for the ratio of sample/population variance, and t.test() to compare their means.

Usage

```
getSamples(data, share = 0.25, confidence.alternative = 0.9,
           max.iter = 50, tracelevel = 1, memorymanagement = TRUE)
```

Arguments

data	The population data from which a sample needs to be taken.
share	The size of the sample in terms of the share of the population data. Defaults to .25.
confidence.alternative	The confidence level used in the F and t-tests defined as the probability level at which the alternative is accepted. For confidence.alternative = .9, we need less evidence to accept the alternative hypothesis that the samples are unequal than at confidence.alternative = .95, hence .90 is stricter than .95.
max.iter	The maximum number of draws to be taken. The program breaks either when a suitable sample is found or when max.iter is reached.
tracelevel	Similar to a verbose statement. Should information be printed during execution? defaults to 1 for printing. set to 0 for no printing.
memorymanagement	TRUE/FALSE indicating whether garbage collection should be forced using gc(). Defaults to TRUE. Recommended setting for large datasets.

Value

A sample of the population dataset that has significantly similar means and variances, or a message indicating that no suitable dat

Examples

```
getSamples (data = ITdata, share = 0.025, confidence.alternative=0.90, max.iter =100)
```

guessStartVal	<i>A function to efficiently obtain starting values for numerical optimization procedures. Used to initialize the "warm start" optimization routines in generalizeToSpecific. The function itself uses a "warm start" algorithm over a growing dataset similar to a sieves estimator for an unbounded parameter space.</i>
---------------	--

Description

This function is called by generalizeToSpecific(), but may also be called by users directly to obtain an initial guess of starting values to be passed on to easygmmlogit().

Usage

```
guessStartVal(Y, X, model, tracelevel = 1, memorymanagement = TRUE)
```

Arguments

Y	A binary response variable.
X	A dataset containing multiple exogenous regressors.
model	The model for which starting values should be estimated. Either "logit" or "probit" for the logit or probit model, or "gmm_nlminb" for a logit model estimated with gmm using PORT routines (reliable) or "gmm_bfgs" using the BFGS algorithm (fast, but still very slow compared to option "logit").
tracelevel	Whether information should be printed during execution. Defaults to 1 for printing, set to 0 for no printing.
memorymanagement	TRUE/FALSE indicating whether garbage collection should be forced regularly when memory usage is high. Defaults to TRUE, recommended setting for large datasets.

Value

A vector of coefficients that can be passed on to numerical optimization algorithms.

Examples

```
data(ITdata)
data(corinetable)
#Grab a sample (optional).
sample <- ITdata[getSamples(ITdata, share =.05),]
# Reclassify
catITdata <- reclassify(sample, reclassstable = corinetable)
# create a binary response dataset.
Y <- MLtoBinomData(catITdata[,1], class =0)
X <- catITdata[,c(2:4)]
initial <- guessStartVal(Y, X, model="logit")
```

iapply	<i>lapply without memory buildup. Iterative application of a function over a vector of arguments, returning only the last result as an object or writing all output to disk.</i>
--------	--

Description

Similar to well-know lapply, but returns only the last result as an object. Very usefull when working with big datasets. If the object to be returned is large in memory, say it contains copies of a large dataset, iapply iteratively applies a function but returns only the last output of the function application. This is useful if the function prints diagnostics to screen or saves results in a log file, but you would like to evaluate the last object to understand the output better. Additionally, the function allows you to write all the output to disk using a Serialization Interface for Single Objects. This allows you to restore any output to an object, possibly with a different name.

Usage

```
iapply(X, FUN, writedisk = FALSE, outdir = "default", name = "default",
       compress = FALSE)
```

Arguments

X	a vector of input variables similar to lapply.
FUN	a function to be applied iteratively over the input arguments. similar to FUN in lapply.
writedisk	TRUE/FALSE indicating whether output of application of FUN to elements of X should be written to disk as serialized representations in RDS files before it gets overwritten in memory. Will also write the final output to disk.
outdir	output directory where .RDS files should be written.
name	When writedisk = TRUE, output is saved with this name followed by the iteration. By default, uses getCall(), thus RDS files are "iapply" followed by the number of th element. However storing iapply in another object, will change the output name.
compress	a logical specifying whether outputted RDS files use "gzip" compress, or one of "gzip", "bzip2" or "xz" to indicate the type of compress to be used. Defaults to no compress. You can later restore the objects using readRDS().

Value

the output of FUN(X[X[length(X)]], plus any prints, messages, warnings, errors that lapply(X, FUN) would produce. If FUN writes results to a disk, these files will be created too.

Examples

```
f <- function(x){
  print (x)
  return(x)
}
x=1:10
# this will print all elements of x, finaloutput will only contain the last element of x.
```

```

finaloutput <- iapply (x, f)
print(finaloutput)

# this will print all elements of x, create 10 RDS files that you can use to restore f(x),
# finaloutput will contain the last element of x.
finaloutput <- iapply (x, f, writedisk = TRUE, outdir = "C:\\Users\\")
print(finaloutput)

# this will do the same, but stores names as "binomials1", "binomials2" etc.
binomials <- iapply
finaloutput <- binomials (x, f, writedisk = TRUE, outdir = "C:\\Users\\")
print(finaloutput)

```

IC

A simple function to extract information criteria from a model object.

Description

This function returns a Kullback-Leibler Information Criterion.

Usage

```
IC(model, KLIC = "AICc", sigma.is.estimated = FALSE)
```

Arguments

<code>model</code>	A model object that is supported by AIC(). For example produced by <code>lm()</code> , <code>glm()</code> , or <code>speedglm()</code> .
<code>KLIC</code>	either "AIC" for the AIC, or "AICc" for the corrected criterion. defaults to "AICc".
<code>sigma.is.estimated</code>	boolean indicating whether variance has been estimated as part of the parameters (Likelihood), or not (LS), If FALSE, number of parameters is increased by 1. Defaults to FALSE. Should not impact the results when models of the same type are compared, but the option has been added to allow users to change the number of counted parameters when comparing ML objects wit hLS objects.

Value

numeric value.

Examples

```

somepars <- c(-0.5, 0.4, -0.3, 0.2, -0.1, 0, 0, 0)
someData <- simulateLogit(nobs=2500, nxregs=8, pars=somepars)
someModel <- lm(formula(someData), data = someData)
IC(someModel)
# this uses k+1 parameters to make it comparable to an MLE object
IC(someModel, KLIC="AIC", sigma.is.estimated = FALSE)
# this calculates with k parameters
IC(someModel, KLIC="AIC", sigma.is.estimated = TRUE)

```

is.speed	<i>Miscellaneous function used in main routines.</i>
----------	--

Description

Miscellaneous function used in main routines.

Usage

```
is.speed(model)
```

Arguments

model a modelobject.

Value

Boolean indicating whether supplied object is of class speedglm.

Examples

```
y=c(1,0,1,0,1,0)
x=c(0,0,1,1,1,0)
is.speed(lm(y~x))
```

ITdata	<i>A small sample dataset produced with getSamples with 157550 observations of 100x100 meter resolution Corine land cover data for Italy with various exogenous variables.</i>
--------	--

Description

To inspect the data, describe(data(ITdata))

Usage

```
data(ITdata)
```

Format

a data frame with 157550 rows and 27 columns

logistic	<i>The logistic function.</i>
----------	-------------------------------

Description

Called by various routines in the package. May also be used to predict fitted conditional probabilities by supplying a set of variables and corresponding coefficients estimated from a logit model.

Usage

```
logistic(theta, data)
```

Arguments

theta	A vector of coefficients.
data	A dataframe of multiple exogenous regressors.

Value

A vector of values produced by a logistic formula under specified parameters and data.

Examples

```
predicted <- logistic(coef(logitmodel), logitmodel$model[, -1])
```

MLtoBinomData	<i>Miscellaneous function used in main routines.</i>
---------------	--

Description

Miscellaneous function used in main routines.

Usage

```
MLtoBinomData(data, class)
```

Arguments

data	factorial or numeric data of a categorial nature.
class	the category that should be reclassified into 1, all other values return as 0.

Value

a binary dataset.

Examples

```
y = c(1,2,3,4,1,2,3,4,1,2,3,4)
```

normalize	<i>A simple function to normalize a dataset of variables to the 0-1 interval.</i>
-----------	---

Description

This function returns a dataframe with normalized variables.

Usage

```
normalize(X)
```

Arguments

X Matrix or dataframe to be normalized.

Value

A dataframe with normalized values ranging from 0 to 1.

Examples

```
someVector = 1:10
normalize(someVector)
df <- data(ITdata)
normalizedIT <- normalize(df)
describe(normalizedIT)
```

oftype	<i>Miscellaneous function used in main routines.</i>
--------	--

Description

Miscellaneous function used in main routines.

Usage

```
oftype(bestfit)
```

Arguments

bestfit a model object.

Value

either "lm", "logit", or "probit".

Examples

```
y=c(1,0,1,0,1,0)
x=c(0,0,1,1,1,0)
oftype(lm(y~x))
```

opt.glm

Core function used by autoGLM, for details and usage see autoGLM.

Description

Core function used by autoGLM, for details and usage see autoGLM.

Usage

```
opt.glm(data, reclassstable, class, outputpath, modelname, tracelevel = 1,
  actions = c("write", "print", "log", "return"), NAval = -9999,
  model = "logit", preselect = "lm", method = "opt.ic", crit.t = 1.64,
  crit.p = 0.1, test = "LR", KLIC = "AICc", accuracytolerance = 0.01,
  confidence.alternative = 0.9, use.share = 0.25, maxsampleruns = 100,
  memorymanagement = TRUE)
```

Arguments

data	See also autoGLM .
reclassstable	autoGLM .
class	autoGLM .
outputpath	autoGLM .
modelname	autoGLM .
tracelevel	autoGLM .
actions	autoGLM .
NAval	autoGLM .
model	autoGLM .
preselect	autoGLM .
method	autoGLM .
crit.t	autoGLM .
crit.p	autoGLM .
test	autoGLM .
KLIC	autoGLM .
accuracytolerance	autoGLM .
confidence.alternative	autoGLM .
use.share	autoGLM .
maxsampleruns	autoGLM .
memorymanagement	autoGLM .

Value

a model object.

opt.glm.core*Core function used by opt.glm, for details and usage see autoGLM.*

Description

Core function used by opt.glm, for details and usage see autoGLM.

Usage

```
opt.glm.core(data, reclassstable, class, outputpath, modelname, tracelevel = 1,  
  actions = c("write", "print", "log", "return"), NAval = -9999,  
  model = "logit", preselect = "lm", method = "opt.ic", crit.t = 1.64,  
  crit.p = 0.05, test = "LR", KLIC = "AICc", accuracytolerance = 0.01,  
  confidence.alternative = 0.9, use.share = 0.25, maxsampleruns = 50,  
  memorymanagement = TRUE)
```

Arguments

data	See also autoGLM .
reclassstable	autoGLM .
class	autoGLM .
outputpath	autoGLM .
modelname	autoGLM .
tracelevel	autoGLM .
actions	autoGLM .
NAval	autoGLM .
model	autoGLM .
preselect	autoGLM .
method	autoGLM .
crit.t	autoGLM .
crit.p	autoGLM .
test	autoGLM .
KLIC	autoGLM .
accuracytolerance	autoGLM .
confidence.alternative	autoGLM .
use.share	autoGLM .
maxsampleruns	autoGLM .
memorymanagement	autoGLM .

Value

a model object.

 opt.h

Optimization routine based on hypothesis testing.

Description

This function searches

Usage

```
opt.h(model, Y, X, returntype, tracelevel, crit.p, test, memorymanagement)
```

Arguments

model	The model to be optimized. Supports "lm" for the linear probability model, "logit" for the logistic probability model, and "probit" for the probit model.
Y	The binary response variable.
X	A dataframe with columns of exogenous regressors.
returntype	"model", "data", or "colnames"
tracelevel	level of printing.
crit.p	p value used in hypothesis tests.
test	type of test, either "LR", "F", or "Chisq".
memorymanagement	logical, indicating whether memory should be more actively managed.

Value

"model", "data", or "colnames", to be specified in returntype.

 opt.ic

Optimization routine based on information criteria.

Description

This function searches

Usage

```
opt.ic(model, Y, X, KLIC, returntype, tracelevel, memorymanagement)
```

Arguments

model	The model to be optimized. Supports "lm" for the linear probability model, "logit" for the logistic probability model, and "probit" for the probit model.
Y	The binary response variable.
X	A dataframe with columns of exogenous regressors.
KLIC	information criterion to be used, "AIC" or "AICc", See also IC
returntype	"model", "data", or "colnames"
tracelevel	level of printing.
memorymanagement	logical, indicating whether memory should be more actively managed.

Value

"model", "data", or "colnames", to be specified in returntype.

Examples

```
opt.ic(model, Y, X, KLIC, returntype, tracelevel, memorymanagement)
```

opt.t	<i>Optimization routine based on step-wise elimination using t-values.</i>
-------	--

Description

This function searches

Usage

```
opt.t(model, Y, X, returntype, tracelevel, crit.t, memorymanagement)
```

Arguments

model	The model to be optimized. Supports "lm" for the linear probability model, "logit" for the logistic probability model, and "probit" for the probit model.
Y	The binary response variable.
X	A dataframe with columns of exogenous regressors.
returntype	"model", "data", or "colnames"
tracelevel	level of printing.
crit.t	t-value used for significance test.
memorymanagement	logical, indicating whether memory should be more actively managed.

Value

"model", "data", or "colnames", to be specified in returntype.

Examples

```
opt.t(model, Y, X, returntype, tracelevel, crit.t, memorymanagement)
```

pkgTest

An improved function to load packages.

Description

This function checks if a package is installed before loading it, and pulls it from the nearest CRAN server if the package is missing.

Usage

```
pkgTest(package, silent = FALSE)
```

Arguments

package	Package to be loaded. Accepts a vector of packages to load multiple packages.
silent	Should startup messages printed by the packages be suppressed? Will not suppress warnings or errors. Defaults to FALSE.

Value

The messages printed by the loaded package, or a message if the package is not found.

Examples

```
verbose_packages <- c("gmm", "foreign", "sp")
silent_packages <- c("data.table", "compiler", "speedglm")

pkgTest(verbose_packages)
pkgTest(silent_packages)
```

PtoBin

Miscellaneous function used in main routines.

Description

Miscellaneous function used in main routines.

Usage

```
PtoBin(Pmap)
```

Arguments

Pmap	a vector of probabilities.
------	----------------------------

Value

A binary vector indicating occurrences.

reclassify

*A function for reclassification of factorial data.***Description**

This function works similar to ArcMaps raster reclassify function, only fast and without starting hell on earth. If you need more speed contact the creator of the package. I could implement a parallel job.

Usage

```
reclassify(LUdata, reclasstable = "default", JIT = TRUE, dropNA = TRUE,
  NAvail = "default")
```

Arguments

LUdata	A dataset (vector, matrix or dataframe) that contains data to be reclassified, e.g., CORINE land cover data. If an object with multiple columns is supplied, the function will attempt to reclassify the first column.
reclasstable	the path to a csv file containing the reclass table. By default, the function loads "CORINE_LUISA_codes", supplied with the package. This table maps Corine land-cover into the Classification scheme used in the LUISA framework of the Joint European Research centre. To inspect the table: <code>data(CORINE_LUISA_codes)</code> .
JIT	a boolean indicating whether Just in Time compilation should be used. Can improve speed in very large datasets. Defaults to TRUE.
dropNA	TRUE/FALSE indicating whether values that are recoded into nodata should be dropped. By default TRUE.
NAvail	the value in the reclasstable that corresponds to nodata. By default the function assumes that the classes have positive values. it assumes that the lowest negative value corresponds to the class that should be dropped (e.g., if your reclass table maps 1,2,3 into 1,2,-9999, the function will return a dataset with classes 1,2.). If your classification scheme maps into negative values and the lowest value does not correspond to the class that should be omitted in the output, specify NAvail if dropNA = TRUE, or the output will miss one class.

Value

A vector of reclassified values.

Examples

```
trainSample <- data(ITdata)
reclasstable <- data(corinetable)
# reclassify a vector of land-use data:
landcover <- reclassify(LUdata=trainSample[,1], reclasstable, JIT=TRUE)

#reclassify a multi-collumn dataset in which the first collumn represents factorial data:
reclass_IT <- reclassify(LUdata=trainSample, reclasstable, JIT=TRUE)

# the corine_to_LUISA codes are loaded by default, so if you wish to reclassify corine land cover to LUISA dat.
reclass_IT <- reclassify(LUdata=trainSample)
```

```
# the LUISA codes are a simplification and drop certain land cover classes. By default, the categories reclass
# If you wish to keep all categories:
reclass_IT <- reclassify(LUdata=trainSample, reclassstable, JIT=TRUE, dropNA=FALSE)

# or if the value that represents NoData is a positive integer, say 9999, and you wish to drop it:
drop9999table <-reclassstable
drop9999table[drop9999table==9999]<-9999
reclass_IT <- reclassify(LUdata=trainSample, reclassstable=drop9999table, JIT=TRUE, dropNA=TRUE, Nval = 9999)
```

selectX

A function to select from an entire set of exogenous regressors, those regressors that minimize Kullback-Leibler divergence in the supplied model. The underlying function uses a "warm start" algorithm initialized by guessStartval, in a fashion similar to the method of sieves. If you plan to use a Lasso estimator on large data, you might consider using this function to determine the zero parameters with option method = "lm", and continue with a Ridge estimator for computational reasons.

Description

This function is a wrapper to the functions bestlinearX(), bestlogitX() and bestprobitX(), with an additional option to call getSamples for improved speed. Take into account that sampling itself takes time, such that total computational burden is a trade-off between the load of the getSample function and the model optimization itself.

Usage

```
selectX(Y, X, model = "lm", returntype = "data", method = "opt.ic",
        KLIC = "AICc", crit.t = 1.64, crit.p = 0.05, test = "LR",
        share = 0.75, confidence.alternative = 0.85, max.iter = 50,
        tracelevel = 1, memorymanagement = TRUE)
```

Arguments

Y	A binary response variable.
X	A dataframe of multiple exogenous regressors.
model	Either "lm" for the linear probability model, "logit" for the logistic probability model, or "probit", for the probit model. The logit and probit models are solved using Iterated Weighted Least Squares, and optimization of the logit model is significantly faster than the probit model. Defaults to "lm".
returntype	Either "data" to return a dataset, or colnames" to only return the column names of the variables that are used in the optimal model. "data" by default.
method	The optimization strategy. Either "opt.ic" to optimize using information criteria, "opt.t" for step-wise elimination of insignificant values (statistically speaking not a sound procedure, but it will provide a parsimonious model that can be usefull as a benchmark), or "opt.h" to optimize by classical hypothesis tests. defaults to "opt.ic".
KLIC	the information criterion used by "opt.ic", either "AIC" or "AICc", defaults to the latter.

<code>crit.t</code>	The t-value indicating significance when using method "opt.t", defaults to 1.64.
<code>crit.p</code>	the p-value used by method "opt.h" in the hypothesis tests. Defaults to 0.05.
<code>test</code>	The hypothesis test used by "opt.h". Defaults to "LR" for the Likelihood Ratio test. Other options are "F", for an F test for joint significance of insignificant parameters, or "Chisq" for a wald test against the Chi squared distribution. Recommended setting is either "LR" as it is less dependent on correct estimation of the standard errors. Keep in mind that "Chisq" is an asymptotic test, and "F" is more appropriate for small sample tests. However "Chisq" holds under milder conditions and should be used if no small sample theory is available for the model.
<code>share</code>	between 0-1, specifying the amount of data that should be passed on to the optimization strategies. Defaults to 0.75, to improve speed. Uses <code>getSamples()</code> to maintain first and second moments of the data.
<code>confidence.alternative</code>	passed on to <code>getSample</code> . Defaults to .85.
<code>max.iter</code>	passed on to <code>getSample</code> . Defaults to 50.
<code>tracelevel</code>	the amount of information to be printed. Passed on to underlying routines. Defaults to 1 for printing, set to 0 for no printing.
<code>memorymanagement</code>	TRUE/FALSE indicating whether garbage collection should be forced regularly when memory usage is high. Defaults to TRUE, recommended setting for large datasets.

Value

Either a dataframe of exogenous variables, or a vector containing the column names indicating the optimal variables extracted from the supplied dataset.

Examples

```
# load data
data(ITdata)
data(corinetable)
#Grab a sample (optional).
sample <- ITdata[getSamples(ITdata, share =.05),]
# Reclassify
catITdata <- reclassify(sample, reclassstable = corinetable)
# create a binary response dataset.
Y <- MLtoBinomData(catITdata[,1], class =1)
X <- catITdata[,-1]
selectX(Y, X, model ="lm", returntype = "colnames", method = "opt.t")
bestX <- selectX(Y, X)
```

`simulateLogit`

A function to generate random data from a logit model.

Description

This function generates random data simulated from a logit model based on the parameters supplied. The function is a crude implementation suited for testing purposes, not designed with Monte Carlo applications in mind. The returned response variable is either binary or truncated probabilities between 0-1.

Usage

```
simulateLogit(nobs = 2500, nxregs = 10, pars = "default", Ytype = "bin")
```

Arguments

nobs	number of observations that the returned dataset should include. Defaults to 2500.
nxregs	number of standardized random regressors that the returned dataset should include. The regressors are used in the data generating process according to the supplied parameters. Defaults to 10.
pars	A vector of parameters to be used in the data generating process. The length of the vector must correspond to the number of random regressors. Regressors may have a 0 coefficient. Defaults to 0 for all variables.
Ytype	Either "bin" to return a binary response variable, or "truncP" to return truncated probabilities between 0 and 1. Defaults to "bin".

Value

A simulated dataset from a logistic probability model with nobs observations, a binary response variable, and nxregs regressors ranging from 0 to 100.

Examples

```
logitpars <- c(-0.5, 0.4, -0.3, 0.2, -0.1, 0, 0, 0, 0, 0, 0, 0)
logit_data <- simulateLogit(nobs=2500, nxregs=12, pars=logitpars, Ytype="bin")

random_binaryResponse_data <- simulateLogit(nobs=500, nxregs=6, pars="default", Ytype="bin")
```

tgc

A function for memory management.

Description

A call to tgc() forces a garbage collection if memory usage is above a certain threshold, be default whgen approximately half of the available mememory is in use by R. Calling tgc(0) performs a standard garbage collection and is therefore the same as gc() with the additional option to delete the output of gc() from memory, such that the call is entirely silent. The function makes a request to the OS to retun memory usage statistics, if this call fails garbage collection will be forced even if memory usage is below the supplied threshold. The function has only been developed for Windows machines, on other OS, the function defaults to tgc(0).

Usage

```
tgc(maxmemshare = 0.5, verbose = FALSE, emptythrash = TRUE)
```

Arguments

maxmemshare	the threshold share of memory usage that triggers garbage collection.
verbose	should diagnostics be printed? FALSE by default.
emptythrash	should the object that contains the output from the garbage collection, including messages, be deleted from memory? TRUE by default. If TRUE, tgc() is silent.

Value

Prints diagnostics of the garbage collection, or silent when emptytrash is TRUE.

Examples

```
tgc()

tgc(maxmemshare=0, verbose = FALSE, emptytrash = TRUE)

tgc(0, TRUE, FALSE)
```

update.start	<i>Miscellaneous function used in main routines.</i>
--------------	--

Description

Miscellaneous function used in main routines.

Usage

```
## S3 method for class 'start'
update(bestfit)
```

Arguments

bestfit modelobject.

Value

vector of starting values

update.X	<i>Miscellaneous function used in main routines.</i>
----------	--

Description

Miscellaneous function used in main routines.

Usage

```
## S3 method for class 'X'
update(bestfit, newX)
```

Arguments

bestfit model object
newX dataframe.

Value

dataframe

update.Z	<i>Miscellaneous function used in main routines.</i>
----------	--

Description

Miscellaneous function used in main routines.

Usage

```
## S3 method for class 'Z'  
update(Y, newX)
```

Arguments

Y	vector
newX	dataframe with corresponding rows.

Value

dataframe

warmstart	<i>Miscellaneous function used in main routines.</i>
-----------	--

Description

Miscellaneous function used in main routines.

Usage

```
warmstart(oldstart, warmth = 1)
```

Arguments

oldstart	vector of starting values.
warmth	numeric between 0-1, to be multiplied by the supplied vector oldstart.

Value

vector of starting values

withinInterval	<i>A function that determines whether a point lies in a specified interval.</i>
----------------	---

Description

A function that determines whether a point lies in a specified interval.

Usage

```
withinInterval(x, interval)
```

Arguments

x	numeric.
interval	vector of length 2 specifying the outer values of the interval.

Value

TRUE/FALSE

Examples

```
withinInterval(5,c(1,10))
```

Index

- *Topic **datasets**
 - corinetable, 5
 - ITdata, 15
- *Topic **data**
 - simulateLogit, 25
- *Topic **generate**
 - simulateLogit, 25
- *Topic **loading**
 - pkgTest, 22
- *Topic **logit**
 - simulateLogit, 25
- *Topic **package**
 - pkgTest, 22
- *Topic **random**
 - simulateLogit, 25
- accuracy, 2, 4
- autoGLM, 3, 18, 19
- corinetable, 3, 5
- describe, 5
- drop.na.model, 6
- easygmmlogit, 6
- exportWeightsfile, 3, 7
- fitnew, 8
- generalizeToSpecific, 3, 9
- getCall, 10
- getSamples, 4, 11, 15
- guessStartVal, 12
- iapply, 4, 13
- IC, 14, 20
- is.speed, 15
- ITdata, 15
- logistic, 16
- MLtoBinomData, 16
- normalize, 17
- oftype, 17
- opt.glm, 18
- opt.glm.core, 19
- opt.h, 4, 20
- opt.ic, 4, 20
- opt.t, 4, 21
- pkgTest, 22
- PtoBin, 22
- reclassify, 3, 23
- selectX, 3, 24
- simulateLogit, 25
- tgc, 4, 26
- update.start, 27
- update.X, 27
- update.Z, 28
- warmstart, 6, 28
- withinInterval, 29