

Package ‘LogisticDx’

July 2, 2014

Type Package

Title Diagnostic tests for logistic regression models

Version 0.1

Date 2014-05-03

Author Chris Dardis

Maintainer Chris Dardis <christopherdardis@gmail.com>

License GPL (>= 2)

Description Diagnostic tests and plots for logistic regression models.

Depends R (>= 2.13.0), multcomp, data.table, pROC, car

Imports gRbase, rms, stats, statmod, graphics, speedglm

LazyLoad yes

Collate 'genLogi.R' 'logiDx.R' 'stukel.R' 'logiGOF.R' 'logiProb.R'
'logiSS.R' 'LogisticDx-internal.R' 'multiPlot.R' 'plotLogiDx.R' 'summaryLogiGOF.R'

Repository CRAN

Repository/R-Forge/Project logisticdx

Repository/R-Forge/Revision 9

Repository/R-Forge/DateTimeStamp 2014-05-14 01:02:09

Date/Publication 2014-05-14 07:28:28

NeedsCompilation no

R topics documented:

LogisticDx-package	2
genLogi	3
logiDx	4
logiGOF	6
logiProb	7
logiSS	9
plotLogiDx	10
stukel	12

Index	14
--------------	-----------

LogisticDx-package	<i>Diagnostic tests for logistic regression models</i>
--------------------	--

Description

Generate sample data frames for logistic regression. Diagnostic tests and plots for logistic regression models

Details

Package:	LogisticDx
Type:	Package
Version:	0.5
Date:	2014-05-01
License:	GPL (>= 2)
LazyLoad:	yes

Author(s)

Christopher Dardis <christopherdardis@gmail.com>

References

Hosmer, D. and Lemeshow, S 2000 *Applied Logistic Regression. Second edition.* John Wiley and Sons, Inc.

genLogi *Generate data for logistic regression*

Description

Generates a data.frame or data.table with a binary outcome, and a logistic model to describe it.

Usage

```
genLogiDf(b = 2L, f = 2L, c = 1L, n = 20L, nlf = 3L,
  pb = 0.5, rc = 0.8, py = 0.5, asFactor = TRUE,
  model = TRUE, timelim = 5, speedglm = FALSE)
```

```
genLogiDt(b = 2L, f = 2L, c = 1L, n = 20L, nlf = 3L,
  pb = 0.5, rc = 0.8, py = 0.5, asFactor = TRUE,
  model = TRUE, timelim = 5, speedglm = FALSE)
```

Arguments

b	<i>binomial predictors</i> , the number of predictors which are binary, i.e. limited to 0 or 1
f	<i>factors</i> , the number of predictors which are factors
c	<i>continuous predictors</i> , the number of predictors which are continuous
n	number of observations in the data frame
nlf	the no. of levels in a factor
pb	<i>probability for binomial predictors</i> : the probability of binomial predictors being = 1 e.g. if pb=0.3, 30% will be 1s, 70% will be 0s
rc	<i>ratio for continuous variables</i> the ratio of levels of continuous variables to the total number of observations <i>n</i> e.g. if rc=0.8 and n=100, it will be in the range 1-80
py	<i>ratio for y</i> the ratio of 1s to total observations for the binomial predictors e.g. if ry=0.5, 50% will be 1s, 50% will be 0s
asFactor	If asFactor=TRUE (the default), predictors given as factors will be converted to factors in the data frame before the model is fit
model	If model=TRUE will also return a model fitted with stats::glm or speedglm::speedglm
timelim	function will timeout after timelim secs. This is present to prevent duplication of rows.
speedglm	If speedglm=TRUE, return a model fitted with speedglm instead of glm

Value

If `model=TRUE`: a list with the following values:

<code>df</code> or <code>dt</code>	A <code>data.frame</code> (for <code>genLogiDf</code>) or <code>data.table</code> (for <code>genLogiDt</code>). Predictors are labelled x_1, x_2, \dots, x_n . Outcome is y . Rows represent to n observations
<code>model</code>	A model fit with <code>stats::glm</code> or <code>speedglm::speedglm</code>

If `model=FALSE` a `data.frame` or `data.table` as above.

Note

`genLogiDt` is faster and more efficient for larger datasets.

Using `asFactor=TRUE` with factors which have a large number of levels (e.g. `nlf > 30`) on large datasets (e.g. $n > 1000$) can cause fitting to be excessively slow.

Examples

```
set.seed(1)
genLogiDf()
genLogiDt(b=0, c=2, n=100, rc=0.7, model=FALSE)
```

logiDx

Diagnostics for logistic regression

Description

Returns standard diagnostic measures for a logistic regression model by covariate pattern

Usage

```
logiDx(x, round = FALSE, roundTo = 3)
```

Arguments

<code>x</code>	A model of class <code>glm</code>
<code>round</code>	If <code>round=TRUE</code> , digits will be rounded to <code>roundTo</code> decimal places
<code>roundTo</code>	No. decimal places to which to round digits

Value

A data table. There is one row per covariate pattern with at least one observation. These are sorted by dBhat (see below).

The initial columns give all combinations of the predictor variables with at least one observation.

Subsequent columns are labelled as follows:

obs	Number of observations with this covariate pattern
prob	Probability of this covariate pattern
yhat	Number of observations of $y = 1$, <i>predicted</i> by the model
y	<i>Actual</i> number of observations of $y = 1$ from the data
lev	<i>Leverage</i> , the diagonal of the hat matrix used to generate the model; a measure of influence of this covariate pattern
devR	<i>Deviance residual</i> , calculated by covariate pattern; a measure of influence of this covariate pattern
PeR	<i>Pearson residual</i> , calculated by covariate pattern; a measure of influence of this covariate pattern. Given by:

$$\sqrt{\text{obs}} \sqrt{\frac{\text{prob}}{(1 - \text{prob})}}$$

sPeR	<i>Standardized Pearson residual</i> calculated by covariate pattern; a measure of influence of this covariate pattern. Given by:
------	---

$$\frac{\text{PeR}}{\sqrt{(1 - \text{lev})}}$$

dBhat	<i>Change in Bhat</i> , the standardized difference between the original maximum likelihood estimates B and that estimates with this covariate pattern excluded
dXsq	<i>Change in chi-square</i> , decrease in the value of Pearson chi-square statistic with this covariate pattern excluded. Given by:

$$sPeR^2$$

dDev	<i>Change in deviance D</i> with this covariate pattern excluded. Given by:
------	---

$$\frac{\text{dev}^2}{(1 - \text{lev})}$$

Note

Values for the statistics are calculated by *covariate pattern*. Different values may be obtained if calculated for each individual observation (i.e. row in data frame).

Generally, the values calculated by covariate pattern are preferred, particularly where no. observations are > 5 .

See Also[plotLogiDx](#)**Examples**

```
d1 <- genLogiDt(model=FALSE)
f1 <- stats::glm(y ~ I(x5^2)*x1 -1, family=binomial("logit"), data=d1)
logiDx(f1)
```

logiGOF

*Goodness of fit tests for a logistic regression model***Description**

Gives 15 commonly employed measures of goodness of fit for a logistic regression model

Usage

```
logiGOF(x, g = 10)
```

Arguments

x	A model of class glm
g	No. groups (quantiles) into which to split observations for Hosmer-Lemeshow and modified Hosmer-Lemeshow tests.

Value

A list of class logiGOF with the following items:

chiPearCov	Pearsons chi-square, calculated by <i>covariate group</i> , with <i>p</i> value and interpretation
chiPearIndiv	Pearsons chi-square, calculated by <i>individual observation</i> , with <i>p</i> value and interpretation
chiPearTab	Pearsons chi-square, calculated by <i>table of covariate patterns by outcome</i> , with <i>p</i> value and interpretation
OsRo	Osius & Rojek test of the logistic link, with <i>p</i> value and interpretation
chiDevCov	Deviance chi-square, calculated by <i>covariate group</i> , with <i>p</i> value and interpretation
chiDevIndiv	Deviance chi-square, calculated by <i>individual observation</i> , with <i>p</i> value and interpretation
chiDevTab	Deviance chi-square, calculated by <i>table of covariate patterns by outcome</i> , with <i>p</i> value and interpretation
covPatTab	Matrix of covariance patterns, used to calculate above chi-square tests of Pearson residuals and deviance

HosLem	Hosmer & Lemeshow goodness of fit test, with g quantile groups, with p value and interpretation
modHosLem	modified Hosmer & Lemeshow goodness of fit test, with g quantile groups, with p value and interpretation
CesHou	le Cessie, van Houwelingen, Copas & Hosmer unweighted sum of squares test for global goodness of fit, with p value and interpretation
Stuk	Stukels test of the appropriateness of the logistic link, with p value and interpretation
PR2	Pearsons R^2 , correlation of observed outcome with predicted
ssR2	Linear regression-like sum-of-squares R^2 , using covariate patterns
l1R2	Log-likelihood based R^2 , calculated by covariate group
ROC	Area under the Receiver Operating Curve, with 95% CI by method of DeLong

Note

A summary method is available

Warning: Will fail if cannot generate a hat matrix for the model using logiDx

Author(s)

Modified Hosmer & Lemeshow goodness of fit test: adapted from existing work by Yongmei Ni

See Also

[logiDx](#)

Examples

```
set.seed(1)
m1 <- genLogiDf(n=100)$model
logiGOF(m1)
```

logiProb	<i>Logit, odds ratio and probability for coefficients of a logistic regression</i>
----------	--

Description

Generate logit, log odds ratio, odds ratio and probability for coefficients in a logistic regression. Can be generalized to all combinations of coefficients.

Values are calculated for a change in the value of the coefficient for the predictor from 0 to 1. (For continuous predictors changes of more than one unit may have more practical significance).

Usage

```
logiProb(x, usePrim = FALSE, all = FALSE)
```

Arguments

x	A logistic regression model of class <code>glm</code>
all	If <code>all=FALSE</code> (the default) return values for all coefficients in model, considered together. If <code>all=TRUE</code> return values for all <i>combinations</i> of coefficients in model.
usePrim	If <code>usePrim=FALSE</code> (the default) use <code>utils::combn</code> to generate combinations. If <code>usePrim=TRUE</code> use <code>gRbase::combnPrim</code> instead (faster).

Value

If `all=TRUE`, a `data.table` giving, for each *combination* of coefficients:

coef	Combination of coefficients
logit	The logit for a given combination of coefficients
lnOR	Natural log of Odds Ratio
OR	Odds Ratio
p	probability

This is sorted by OR (low to high).

If `all=FALSE`, a `data.frame` giving the above values for all predictors.

Note

To use `gRbase::combnPrim` the following dependencies may be necessary. Install as follows:

```
source("http://bioconductor.org/biocLite.R")
biocLite("graph")
biocLite("BiocGenerics")
biocLite("RBGL")
```

Examples

```
set.seed(1)
f1 <- genLogiDf(n=50)$model
logiProb(f1)
d1 <- genLogiDt(n=50, model=FALSE)
logiProb(glm(y ~ x1 + x3 -1, data=d1, family=binomial()), all=TRUE)
```

 logiSS

Sample size for given coefficient and events per covariate for model

Description

Gives sample size necessary to demonstrate that coefficient in model for given predictor is equal to its given value (rather than equal to zero) for a given level of power and significance.

Also number of events (smaller of outcome $y = 0$ and outcome $y = 1$) per predictor.

Uses different methods depending on whether model has one binomial, one continuous or multiple predictors.

Usage

```
logiSS(x, alpha = 0.05, beta = 0.8, coeff = "x1")
```

Arguments

x	A logistic regression model of class glm
alpha	significance level α for null-hypothesis significance test
beta	power β for null-hypothesis significance test
coeff	Name of predictor (coefficient) in model to be tested

Value

A list of:

res	Result: Sample size required to show coefficient for predictor is as given in the model rather than 0
epc	Events per covariate; should be >10 to make meaningful statements about coefficients obtained

Examples

```
set.seed(1)
### one coefficient, which is binomial
f1 <- genLogiDf(b=1, c=0, n=50)$model
logiSS(f1)
###
### one coefficient, which is continuous
f1 <- genLogiDf(f=0, b=0, c=1, n=50)$model
logiSS(f1, coeff="x1")
###
### binomial coefficient
f1 <- genLogiDf(f=0, b=1, c=1, n=50)$model
logiSS(f1, coeff="x1")
```

```
###
### continuous coefficient
f1 <- genLogiDf(f=0, b=1, c=1, n=50)$model
logiSS(f1, coeff="x2")
```

plotLogiDx

Diagnostic plots for a logistic regression

Description

Common diagnostic plots for a logistic regression model

Usage

```
plotLogiDx(x, noPerPage = 6,
  cols = c("deepskyblue", "dodgerblue"), cex = 2,
  pch = 21, inches = 0.25, identify = FALSE,
  extras = FALSE, width = NULL, height = NULL)
```

Arguments

x	A logistic regression model of class <code>glm</code>
noPerPage	Number of plots per page (for initial plots). Will be used as <i>guidance</i> and optimised for ease of display
cols	Colours. As used by <code>graphics::points</code>
cex	Cex Character exp ansion. See <code>?graphics::plot.default</code>
pch	P lotting ch aracter. See <code>?graphics::points</code>
inches	Width of circles for bubble plot. See <code>?graphics::symbols</code>
identify	If TRUE will give option to identify individual points on a number of the plots produced. The number which appears next to the point corresponds to the relevant row as given by <code>logiGOF</code>
extras	If TRUE produces additional plots, detailed below
width	Width of screen(display device) in pixels
height	Height of screen(display device) in pixels

Value

The following are plotted, for each covariate group:

p_X_lev	Probability of $y = 1$ for this group by leverage (diagonal of hat matrix, a measure of influence)
p_X_dXsq	Probability as above by $dXsq$ change in Pearson chi-square statistic with deletion of this group
p_X_dBhat	Probability by $dBhat$ change in Bhat; the difference in the maximum likelihood estimators Beta for model coefficients with all subjects included vs those with this group, standardized by the estimated covariance matrix of Beta
p_X_dDev	Probability by $dDev$, the change in deviance when this group is excluded
bubbleplot	Probability by $dXsq$, with area of circle proportional to $dBhat$
lev_X_dXsq	Leverage by $dXsq$, the change in the Pearson chi-square statistic when this group is excluded
lev_X_dBhat	Leverage by $dBhat$, the difference in the maximum likelihood estimators Beta for model coefficients with all subjects included vs those when this group is excluded. This is standardized by the estimated covariance matrix of Beta
lev_X_dDev	Leverage by $dDev$, the change in deviance when this group is excluded
ROC	Receiver Operator Curve

Additional plots are given when extras=TRUE:

influenceplot	See <code>?car::influencePlot</code>
sr_X_hat	<i>Studentized residual</i> by hat values. Studentized residual = residual / estimate of standard deviation of residual
slp	Spread-level plot. See from <code>?car::spreadLevelPlot</code>
qqPlot	quantile-quantile plot vs Normal for residuals. See <code>?stats::qqplot</code>
iip	Influence-index plot. Gives Cooks distance, studentized residual and hat values for each observation
pairs	Pairs plot for the measures of influence $dBhat$, $dXsq$ and $dDev$. See <code>?graphics::pairs</code>
crPlots	Component + residual plots. See <code>?car::crPlots</code>
avPlots	Added-variable plots. See <code>?car::avPlots</code>
mmps	Marginal model plots. These require that the <code>data.frame</code> used to fit the model be present in the current environment. See <code>?car::mmps</code>

Note

Different colors can be found with e.g.
`grDevices::colours()[grep("blue",grDevices::colours())]`

Examples

```
set.seed(1)
### generate up to 8x covariate patterns
mod1 <- genLogiDf(b=3, f=0, c=0, n=50)$model
plotLogiDx(mod1, cex=8, noPerPage=1)
plotLogiDx(mod1, cex=3, noPerPage=6, extras=TRUE)
df1 <- genLogiDf(b=0, f=0, c=2, n=50, model=FALSE)
g1 <- glm(y ~ ., family=binomial("logit"), data=df1)
plotLogiDx(g1)
```

 stukel

Stukels test of the logistic link

Description

Calculates Stukels test for a logistic regression model.

This determines the appropriateness of the logistic link.

Two new covariates, $z1$ and $z2$ are generated such that

$$z1 = 0.5 * \text{logit}^2 * I(pi \geq 0.5), z2 = -0.5 * \text{logit}^2 * I(pi \leq 0.5)$$

where $I(arg) = 1$ where arg is true and $= 0$ where false.

Tests if significant change occurs in the model with the addition of these coefficients.

Usage

```
stukel(object,
       alternative = c("both", "alpha1", "alpha2"))
```

Arguments

object	An object of class glm
alternative	add both $z1$ and $z2$ to model or just one of the above

Value

A list with the following elements:

statistic	value of statistic
parameter	degrees of freedom
p.value	if < 0.05 suggests should accept null hypothesis that logistic model is incorrect for the data
alternative	alternative
method	method
data.name	name of object
allstat	statistics for all tests
allpar	degrees of freedom
allpval	all p values

Author(s)

Brett Presnell

References

[University of Florida](#)

Examples

```
set.seed(1)
m1 <- genLogiDf(b=3, f=0, c=0, n=50)$model
stukel(m1)
```

Index

*Topic **datagen**

genLogi, [3](#)

*Topic **hplot**

plotLogiDx, [10](#)

*Topic **htest**

logiGOF, [6](#)

logiSS, [9](#)

stukel, [12](#)

*Topic **package**

LogisticDx-package, [2](#)

genLogi, [3](#)

genLogiDf (genLogi), [3](#)

genLogiDt (genLogi), [3](#)

logiDx, [4](#), [7](#)

logiGOF, [6](#), [10](#)

logiProb, [7](#)

logiSS, [9](#)

LogisticDx (LogisticDx-package), [2](#)

logisticDx (LogisticDx-package), [2](#)

LogisticDx-package, [2](#)

plotLogiDx, [6](#), [10](#)

stukel, [12](#)