

Package ‘StepReg’

November 3, 2017

Type Package

Title Stepwise Regression Analysis

Version 1.0.0

Date 2017-10-30

Author Junhui Li, Kun Cheng, Wenxin Liu

Maintainer Junhui Li <junhuili@cau.edu.cn>

Description Stepwise regression analysis for variable selection can be used to get the best candidate final regression model in univariate or multivariate regression analysis with the 'forward' and 'stepwise' steps. Procedure uses Akaike information criterion, the small-sample-size corrected version of Akaike information criterion, Bayesian information criterion, Hannan and Quinn information criterion, the corrected form of Hannan and Quinn information criterion, Schwarz criterion and significance levels as selection criteria, where the significance levels for entry and for stay are set to 0.15 as default. Multicollinearity detection in regression model are performed by checking tolerance value, which is set to 1e-7 as default. Continuous variables nested within class effect are also considered in this package.

License GPL (>= 2)

Imports Rcpp (>= 0.12.13)

LinkingTo Rcpp, RcppEigen

Depends R (>= 2.10)

NeedsCompilation yes

Repository CRAN

Date/Publication 2017-11-03 10:17:07 UTC

R topics documented:

| | |
|-------------------------------|---|
| StepReg-package | 2 |
| bestCandidate_RCcpp | 3 |
| stepwise | 5 |

| | |
|--------------|----------|
| Index | 8 |
|--------------|----------|

Description

Stepwise regression analysis for variable selection can be used to get the best candidate final regression model in univariate or multivariate regression analysis with the 'forward' and 'stepwise' steps. Procedure uses Akaike information criterion, the small-sample-size corrected version of Akaike information criterion, Bayesian information criterion, Hannan and Quinn information criterion, the corrected form of Hannan and Quinn information criterion, Schwarz criterion and significance levels as selection criteria, where the significance levels for entry and for stay are set to 0.15 as default. Multicollinearity detection in regression model are performed by checking tolerance value, which is set to 1e-7 as default. Continuous variables nested within class effect are also considered in this package.

Details

Package: StepReg
Type: Package
Version: 1.0.0
Date: 2017-10-30
License: GPL (>= 2)

Author(s)

Junhui Li, Kun Cheng, Wenxin Liu
Maintainer: Junhui Li <junhuili@cau.edu.cn>

References

- Alsubaihi, A. A., Leeuw, J. D., and Zeileis, A. (2002). Variable selection in multivariable regression using sas/iml. , 07(i12).
- Darlington, R. B. (1968). Multiple regression in psychological research and practice. Psychological Bulletin, 69(3), 161.
- Hannan, E. J., & Quinn, B. G. (1979). The determination of the order of an autoregression. Journal of the Royal Statistical Society, 41(2), 190-195.
- Harold Hotelling. (1992). The Generalization of Student's Ratio. Breakthroughs in Statistics. Springer New York.
- Hurvich, C. M., & Tsai, C. (1989). Regression and time series model selection in small samples. Biometrika, 76(2), 297-307.
- Judge, & George G. (1985). The Theory and practice of econometrics /-2nd ed. The Theory and practice of econometrics /. Wiley.

- Mardia, K. V., Kent, J. T., & Bibby, J. M. (1979). Multivariate analysis. *Mathematical Gazette*, 37(1), 123-131.
- Mckeeon, J. J. (1974). F approximations to the distribution of hotelling's t20. *Biometrika*, 61(2), 381-383.
- Mcquarrie, A. D. R., & Tsai, C. L. (1998). *Regression and Time Series Model Selection. Regression and time series model selection /*. World Scientific.
- Pillai, K. C. S. (2006). Pillai's Trace. *Encyclopedia of Statistical Sciences*. John Wiley & Sons, Inc.
- R.S. Sparks, W. Zucchini, & D. Coutsourides. (1985). On variable selection in multivariate regression. *Communication in Statistics- Theory and Methods*, 14(7), 1569-1587.
- Sawa, T. (1978). Information criteria for discriminating among alternative regression models. *Econometrica*, 46(6), 1273-1291.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6(2), pages. 15-18.

bestCandidate_RCcpp *Obtain one best candidate variable*

Description

Get best candidate variable with forward or backward direction in only one step

Usage

```
bestCandidate_RCcpp(findIn, p, n, sigma, tolerance, Ftrace, criteria, Y, X1, X0, k)
```

Arguments

| | |
|-----------|--|
| findIn | Logical value, if FALSE then add independent variable to regression model, otherwise remove independent variable from regression model |
| p | The number of independent variable entered in regression |
| n | The sample size |
| sigma | Pure error variance from full regressoin model for Bayesian information criterion(BIC) |
| tolerance | Tolerance value for multicollinearity |
| Ftrace | Statistic of multivariate regression including Wilks' lambda, Pillai trace and Hotelling-lawley trace |
| criteria | Information criterion including AIC, AICc, BIC, SBC, HQ, HQc and SL |
| Y | Data set for dependent variable |
| X1 | Data set for independent variables not in regression model |
| X0 | Data set for independent variables entered in regression model |
| k | Forces the first k effects entered in regression model, and the selection methods are performed on the other effects in the data set |

Details

This function can compute probability value or information criteria statistics with multivariate and univariate regression using least square method

Value

| | |
|------|---|
| PIc | P value or Information Criteria statistic value |
| seq | Pointer for independent variable enter or eliminate |
| SSE | Maximum or minimum of SSE |
| RkCh | Rank changed or not |

Author(s)

Junhui Li

References

- Alsubaihi, A. A., Leeuw, J. D., and Zeileis, A. (2002). Variable selection in multivariable regression using sas/iml. , 07(i12).
- Darlington, R. B. (1968). Multiple regression in psychological research and practice. *Psychological Bulletin*, 69(3), 161.
- Hannan, E. J., & Quinn, B. G. (1979). The determination of the order of an autoregression. *Journal of the Royal Statistical Society*, 41(2), 190-195.
- Harold Hotelling. (1992). *The Generalization of Student's Ratio. Breakthroughs in Statistics.* Springer New York.
- Hurvich, C. M., & Tsai, C. (1989). Regression and time series model selection in small samples. *Biometrika*, 76(2), 297-307.
- Judge, & GeorgeG. (1985). *The Theory and practice of econometrics /-2nd ed. The Theory and practice of econometrics /.* Wiley.
- Mardia, K. V., Kent, J. T., & Bibby, J. M. (1979). Multivariate analysis. *Mathematical Gazette*, 37(1), 123-131.
- Mckeeon, J. J. (1974). F approximations to the distribution of hotelling's t20. *Biometrika*, 61(2), 381-383.
- Mcquarrie, A. D. R., & Tsai, C. L. (1998). *Regression and Time Series Model Selection. Regression and time series model selection /.* World Scientific.
- Pillai, K. C. S. (2006). Pillai's Trace. *Encyclopedia of Statistical Sciences.* John Wiley & Sons, Inc.
- R.S. Sparks, W. Zucchini, & D. Coutsourides. (1985). On variable selection in multivariate regression. *Communication in Statistics- Theory and Methods*, 14(7), 1569-1587.
- Sawa, T. (1978). Information criteria for discriminating among alternative regression models. *Econometrica*, 46(6), 1273-1291.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6(2), pags. 15-18.

Description

Stepwise regression analysis for variable selection can be used to get the best candidate final regression model in univariate or multivariate regression analysis with the 'forward' and 'stepwise' steps. Procedure uses Akaike information criterion, the small-sample-size corrected version of Akaike information criterion, Bayesian information criterion, Hannan and Quinn information criterion, the corrected form of Hannan and Quinn information criterion, Schwarz criterion and significance levels as selection criteria, where the significance levels for entry and for stay are set to 0.15 as default. Multicollinearity detection in regression model are performed by checking tolerance value, which is set to 1e-7 as default. Continuous variables nested within class effect are also considered in this package.

Usage

```
stepwise(data, y, notX, include, Class, selection, select, sle, sls, tolerance,
Trace, Choose)
```

Arguments

| | |
|-----------|--|
| data | Data set including dependent and independent variables to be analyzed |
| y | Numeric or character vector for dependent variables |
| notX | Numeric or character vector for independent variables removed from stepwise regression analysis |
| include | Forces the effects vector listed in the data to be included in all models. The selection methods are performed on the other effects in the data set |
| Class | Class effect variable |
| selection | Model selection method including "forward" and "stepwise", forward selection starts with no effects in the model and adds effects, while stepwise regression is similar to the forward method except that effects already in the model do not necessarily stay there |
| select | specifies the criterion that uses to determine the order in which effects enter and/or leave at each step of the specified selection method including Akaike Information Criterion(AIC), the Corrected form of Akaike Information Criterion(AICc), Bayesian Information Criterion(BIC), Schwarz criterion(SBC), Hannan and Quinn Information Criterion(HQ), Significant Levels(SL) and so on |
| sle | Specifies the significance level for entry |
| sls | Specifies the significance level for staying in the model |
| tolerance | Tolerance value for multicollinearity, default is 1e-7 |
| Trace | Statistic for multivariate regression analysis, including Wilks' lamda ("Wilks"), Pillai Trace ("Pillai") and Hotelling-Lawley's Trace ("Hotelling") |

Choose Chooses from the list of models at the steps of the selection process the model that yields the best value of the specified criterion. If the optimal value of the specified criterion occurs for models at more than one step, then the model with the smallest number of parameters is chosen. If you do not specify the Choose option, then the model selected is the model at the final step in the selection process

Details

Multivariate regression and univariate regression can be detected by parameter 'y', where numbers of elements in 'y' is more than 1, then multivariate regression is carried out otherwise univariate regression

Author(s)

Junhui Li

References

- Alsubaihi, A. A., Leeuw, J. D., and Zeileis, A. (2002). Variable selection in multivariable regression using sas/iml. , 07(i12).
- Darlington, R. B. (1968). Multiple regression in psychological research and practice. *Psychological Bulletin*, 69(3), 161.
- Hannan, E. J., & Quinn, B. G. (1979). The determination of the order of an autoregression. *Journal of the Royal Statistical Society*, 41(2), 190-195.
- Harold Hotelling. (1992). *The Generalization of Student's Ratio. Breakthroughs in Statistics.* Springer New York.
- Hurvich, C. M., & Tsai, C. (1989). Regression and time series model selection in small samples. *Biometrika*, 76(2), 297-307.
- Judge, & George G. (1985). *The Theory and practice of econometrics /-2nd ed. The Theory and practice of econometrics /.* Wiley.
- Mardia, K. V., Kent, J. T., & Bibby, J. M. (1979). Multivariate analysis. *Mathematical Gazette*, 37(1), 123-131.
- Mckeeon, J. J. (1974). F approximations to the distribution of hotelling's t^2 . *Biometrika*, 61(2), 381-383.
- Mcquarrie, A. D. R., & Tsai, C. L. (1998). *Regression and Time Series Model Selection. Regression and time series model selection /.* World Scientific.
- Pillai, K. C. S. (2006). Pillai's Trace. *Encyclopedia of Statistical Sciences.* John Wiley & Sons, Inc.
- R.S. Sparks, W. Zucchini, & D. Coutsourides. (1985). On variable selection in multivariate regression. *Communication in Statistics- Theory and Methods*, 14(7), 1569-1587.
- Sawa, T. (1978). Information criteria for discriminating among alternative regression models. *Econometrica*, 46(6), 1273-1291.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6(2), pages. 15-18.

Examples

```

set.seed(4)
dfY <- data.frame(matrix(c(rnorm(20,0,2),c(rep(1,10),rep(2,10))),rnorm(20,2,3)),20,3))
colnames(dfY) <- paste("Y",1:3,sep="")
dfX <- data.frame(matrix(c(rnorm(100,0,2),rnorm(100,2,1)),20,10))
colnames(dfX) <- paste("X",1:10,sep="")
dfyx <- cbind(dfY,dfX)

#for univariate regression
y <- c("Y1")
notX <- c("Y3")
#for multivariate regression you can use this
ym <- c("Y1","Y3")
notXm <- NULL
#* with continuous variable nested in class effect
ClassY2 <- c("Y2")
#* without continuous variable nested in class effect
Class0 <- NULL
# without forced effect in regression model
include0 <- NULL
# force the 'Y2' into the regression model
includeY2 <- c("Y2")
selection <- 'stepwise'
tolerance <- 1e-7
Trace <- "Pillai"
sle <- 0.15
sls <- 0.15

#univariate regression for 'SBC' select and 'AIC' choose
#without forced effect and continuous variable nested in class effect
stepwise(dfyx, y, notX, include0, Class0, selection, "SBC", sle, sls, tolerance, Trace, 'AIC')

#univariate regression for 'AICc' select and 'HQc' choose
#with forced effect and continuous variable nested in class effect
stepwise(dfyx, y, notX, includeY2, ClassY2, selection, 'AICc', sle, sls, tolerance, Trace, 'HQc')

#multivariate regression for 'HQ' select and 'BIC' choose
#with forced effect and continuous variable nested in class effect
stepwise(dfyx, ym, notXm, includeY2, ClassY2, selection, 'HQ', sle, sls, tolerance, Trace, 'BIC')

```

Index

*Topic **package**

StepReg-package, [2](#)

*Topic **stepwise regression**

bestCandidate_RCcpp, [3](#)

stepwise, [5](#)

bestCandidate_RCcpp, [3](#)

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

StepReg-package, [2](#)

stepwise, [5](#)