

Package ‘cctest’

February 25, 2026

Version 2.3.1

Title Canonical Correlations and Tests of Independence

Description A simple interface for multivariate correlation analysis that unifies various classical statistical procedures including t-tests, tests in univariate and multivariate linear models, parametric and nonparametric tests for correlation, Kruskal-Wallis tests, common approximate versions of Wilcoxon rank-sum and signed rank tests, chi-squared tests of independence, score tests of particular hypotheses in generalized linear models, canonical correlation analysis and linear discriminant analysis.

Author Robert Schlicht [aut, cre]

Maintainer Robert Schlicht <robert.schlicht@tu-dresden.de>

License MIT + file LICENSE | EUPL (>= 1.1)

Imports stats

NeedsCompilation no

Repository CRAN

Date/Publication 2026-02-25 11:20:02 UTC

Contents

cctest 1

Index 8

cctest *Tests of Independence Based on Canonical Correlations*

Description

cctest estimates canonical correlations between two sets of variables, possibly after removing effects of a third set of variables, and performs a classical multivariate test of (conditional) independence based on Pillai’s statistic.

Usage

```
cctest(formula, data=NULL, weights=NULL, ..., tol=1e-07, stats=FALSE)
```

Arguments

formula	An object of the form $Y \sim X \sim A$, where Y represents dependent variables, X represents a second set of dependent variables or explanatory variables not present under the null hypothesis, and A represents explanatory variables that remain under the null hypothesis. Typically A includes at least the constant 1 to specify a model with intercepts; unlike <code>lm</code> , the function never adds this automatically. If <code>stats</code> is <code>FALSE</code> (recommended, see Note), A , X , Y are specified in a simplified notation where all symbols and operators except <code> </code> and <code>:</code> have their regular meaning, with <code> </code> used for joining terms (instead of <code>+</code> ; see <code>cbind</code>) and <code>:</code> for multiplying terms, valid anywhere in the expression. The latter differs from <code>*</code> in that each combination of columns in the arguments contributes a column to the result and multiplication by 0 always yields 0, even for missing values. Single-row arguments in <code>~</code> , <code> </code> , <code>:</code> are expanded to objects with identical rows, and every factor or character variable is represented by its full set of indicator (dummy) variables, with values of character variables <code>sorted</code> by the "radix" method.
data	An optional list (or data frame) in which the variables in <code>formula</code> , <code>weights</code> and <code>...</code> are looked for (prior to lookup in the environment in which the <code>~</code> preceding A was invoked).
weights	An optional object of the form $\sim A_0$ or $w \sim A_0$, where A_0 replaces A for the degrees of freedom computation; the default means both are the same. The weights vector w consists of any nonnegative numbers, equal to 1 by default, that specify how many identical observations each row represents. For rows of data with missing values these numbers are internally replaced with 0, and rows with missing weights are removed.
...	Additional optional arguments in the case that <code>stats</code> is <code>TRUE</code> .
tol	The tolerance in the QR decomposition for detecting linear dependencies of the matrix columns.
stats	A logical value. If <code>TRUE</code> , the expressions in <code>formula</code> , <code>data</code> , <code>weights</code> and <code>...</code> are passed to <code>model.frame</code> and <code>model.matrix</code> for processing; the <code>operators</code> and <code>expansion</code> rules defined for the <code>model</code> part of a <code>formula</code> object here apply to all parts. If <code>FALSE</code> , <code>formula</code> and <code>weights</code> are interpreted according to the simplified notation.

Details

`cctest` unifies various classical statistical procedures that involve the same underlying computations, including t-tests, tests in univariate and multivariate linear models, parametric and nonparametric tests for correlation, Kruskal–Wallis tests, common approximate versions of Wilcoxon rank-sum and signed rank tests, chi-squared tests of independence, score tests of particular hypotheses in generalized linear models, canonical correlation analysis and linear discriminant analysis (see Examples).

Specifically, for the matrices with ranks K and L obtained from X and Y by subtracting from each column its orthogonal projection on the column space of A , the function computes factorizations

$\tilde{X}U$ and $\tilde{Y}V$ with \tilde{X} and \tilde{Y} having K and L columns, respectively, such that both $\tilde{X}^\top \tilde{X} = rI$ and $\tilde{Y}^\top \tilde{Y} = rI$, and $\tilde{X}^\top \tilde{Y} = rD$ is zero except for nonnegative elements on its main diagonal arranged in decreasing order. The scaling factor r , which should be nonzero, is the dimension of the orthogonal complement of the column space of A_0 .

The function realizes this variant of the singular value decomposition by first computing preliminary QR factorizations of the stated form (taking $r = 1$) without the requirement on D , and then, in a second step, modifying these based on a standard singular value decomposition of that matrix. The main work is done in a rotated coordinate system where the column space of A aligns with the coordinate axes. The basic approach and the rank detection algorithm are inspired by the implementations in [cancor](#) and in [lm](#), respectively.

The main diagonal elements of D , or singular values, are the estimated *canonical correlations* (Hotelling, 1936) of the variables represented by X and Y if these follow a linear model $(X \ Y) = A(\alpha \ \beta) + (\delta \ \epsilon)$ with known A , unknown $(\alpha \ \beta)$ and error terms $(\delta \ \epsilon)$ that have uncorrelated rows with expectation zero and an identical unknown covariance matrix. In the most common case, where A is given as a constant 1, these are the sample canonical correlations (i.e., based on simple centering) most often presented in the literature for full column ranks K and L . They are always decreasing and between 0 and 1.

In the case of the linear model with independent normally distributed rows and $A_0 = A$, the ranks K and L equal, with probability 1, the ranks of the covariance matrices of the rows of X and Y , respectively, or r , whichever is smaller. Under the hypothesis of independence of X and Y , given those ranks, the joint distribution of the J squared singular values, where J is the smaller of the two ranks, is then known and in the case $r \geq K + L$ has a probability density (Hsu, 1939, Anderson, 2003, Anderson, 2007) given by

$$\rho(t_1, \dots, t_J) \propto \prod_{j=1}^J t_j^{(|K-L|-1)/2} (1 - t_j)^{(r-K-L-1)/2} \prod_{j'>j} (t_j - t_{j'}),$$

$1 \geq t_1 \geq \dots \geq t_J \geq 0$. For $J = 1$ this reduces to the well-known case of a single beta distributed R^2 or equivalently an F distributed $\frac{R^2/(KL)}{(1-R^2)/(r-KL)}$, with the divisors in the numerator and denominator representing the degrees of freedom, or twice the parameters of the beta distribution.

Pillai's statistic is the sum of squares of the canonical correlations, which equals, even without the diagonal requirement on D , the squared Frobenius norm of that matrix (or trace of $D^\top D$). Replacing the distribution of that statistic divided by J (i.e., of the mean of squares) with beta or gamma distributions with first or shape parameter $KL/2$ and expectation $KL/(rJ)$ leads to the F and chi-squared approximations that the p-values returned by `cctest` are based on.

The F or beta approximation (Pillai, 1954, p. 99, p. 44) is usually used with $A_0 = A$ and then is exact if $J = 1$. The chi-squared approximation represents Rao's (1948) score test (with a test statistic that is r times Pillai's statistic) in the model obtained after removing (or conditioning on) the orthogonal projections on the column space of A_0 provided that is a subset of the column space of A ; see Mardia and Kent (1991) for the case with independent identically distributed rows.

Value

A list with class `htest` containing the following components:

<code>x, y</code>	matrices \tilde{X} and \tilde{Y} of new transformed variables
<code>xinv, yinv</code>	matrices U and V representing the inverse coordinate transformations

<code>estimate</code>	vector of canonical correlations, i.e., the diagonal elements of D , possibly (only if $K = L = 1$) with a name indicating the direction of the correlation
<code>statistic</code>	vector of p-values based on Pillai's statistic and classical F (beta) and chi-squared (gamma) approximations
<code>df.residual</code>	the number r
<code>method</code>	the name of the function
<code>data.name</code>	a character string representation of formula (possibly shortened)

Note

The handling of the weights differs from that in `lm` unless the nonzero weights are scaled so as to have a mean of 1. Also, to facilitate predictions for rows with zero weights (see Examples), the square roots of the weights, used internally for scaling the data, are always computed as nonzero numbers, even for zero weights, where they are so small that their square is still numerically zero and hence without effect on the correlation analysis.

The simplified formula notation is intended to provide a simpler, more consistent behavior than the legacy stats procedure based on `terms.formula`, `model.frame` and `model.matrix`. Inconsistent or hard-to-predict behavior can result in `model.matrix`, in particular, from the special interpretation of common symbols, the identification of variables by deparsed expressions, the locale-dependent conversion of character variables to factors and the attempts at reducing linear dependencies subject to options("contrasts"). The stats procedure may be removed in the future.

For the classical rank tests shown in the Examples, data must be passed to `cctest` with subsetting and removal of missing values done beforehand.

Author(s)

Robert Schlicht

References

- Hotelling, H. (1936). Relations between two sets of variates. *Biometrika* 28, 321–377. doi:10.1093/biomet/28.34.321, doi:10.2307/2333955
- Hsu, P.L. (1939). On the distribution of roots of certain determinantal equations. *Annals of Eugenics* 9, 250–258. doi:10.1111/j.14691809.1939.tb02212.x
- Rao, C.R. (1948). Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation. *Mathematical Proceedings of the Cambridge Philosophical Society* 44, 50–57. doi:10.1017/S0305004100023987
- Pillai, K.C.S. (1954). *On some distribution problems in multivariate analysis* (Institute of Statistics mimeo series 88). North Carolina State University, Dept. of Statistics.
- Mardia, K.V., Kent, J.T. (1991). Rao score tests for goodness of fit and independence. *Biometrika* 78, 355–363. doi:10.1093/biomet/78.2.355
- Anderson, T.W. (2003). *An introduction to multivariate statistical analysis*, 3rd edition, Ch. 12–13. Wiley.
- Anderson, T.W. (2007). Multiple discoveries: distribution of roots of determinantal equations. *Journal of Statistical Planning and Inference* 137, 3240–3248. doi:10.1016/j.jspi.2007.03.008

See Also

Functions [cancor](#), [anova.mlm](#) in package `stats` and implementations of canonical correlation analysis in other packages such as `CCP` (tests only), `MVar`, `candisc` (both including tests based on Wilks' statistic), `yacca`, `CCA`, `acca`, `whitening`.

Examples

```
## Artificial observations in 5-by-5 meter subplots in a forest for
## comparing cctest analyses with equivalent 'stats' methods:
dat <- within(data.frame(row.names=1:150), {
  u <- function() replicate(150, z<<-(z*69069+2^-32)%%1); z<-0
  plot <- factor(u() < .5, , c("a","b")) # plot a or b
  x <- as.integer(30*u()) + c(1,82)[plot] # x position on grid
  y <- as.integer(30*u()) + c(1,62)[plot] # y position on grid
  ori <- factor(u()%/.25, c("E", "N", "S", "W")) # orientation of slope
  elev <- 40*u() + c(605,610)[plot] # elevation (in meters)
  h <- 115 - .15*elev + 2*log(1/u()-1) # tree height (in meters)
  h5 <- h + log(1/u()-1) # tree height 5 years earlier
  h10 <- h5 + log(1/u()-1) # tree height 10 years earlier
  c15 <- as.integer(h10 + log(1/u()-1) > 20) # 0-1 coded, 15 years earlier
  sap1 <- as.integer(log(1/u())^.8*elev/40) # number of saplings
  rm(u, z)
})
## Not run:
dat

## t-tests:
cctest(h~plot~1, dat)
t.test(h~plot, dat, var.equal=TRUE)
summary(lm(h~plot, dat))
cctest(h~20~1~0, dat)
t.test(dat$h, mu=20)
t.test(h~1, dat, mu=20)
cctest(h~h5~1~0, dat)
t.test(dat$h, dat$h5, paired=TRUE)
t.test(Pair(h,h5)~1, dat)

## Test for correlation:
cctest(h~elev~1, dat)
cor.test(~h+elev, dat)

## One-way analysis of variance:
cctest(h~ori~1, dat)
anova(lm(h~ori, dat))
oneway.test(h~ori, dat, var.equal=TRUE)

## F-tests in linear models:
cctest(h~ori~1|elev, dat)
anova(lm(h~1+elev, dat), lm(h~ori+elev, dat))
cctest(h~h5~(h5-h10):(1|x|x^2)~0, dat)
summary(lm(h~h5~0+I(h5-h10)+I(h5-h10):(x+I(x^2)), dat))
```

```

## Test in multivariate linear model based on Pillai's statistic:
cctest(h|h5|h10~x|y~1|elev, dat)
  anova(lm(cbind(h,h5,h10)~elev, dat), lm(cbind(h,h5,h10)~elev+x+y, dat))

## Test based on Spearman's rank correlation coefficient:
cctest(rank(h)~rank(elev)~1, dat)
  cor.test(~h+elev, dat, method="spearman", exact=FALSE)

## Kruskal-Wallis and Wilcoxon rank-sum tests:
cctest(rank(h)~ori~1, dat)
  kruskal.test(h~ori, dat)
cctest(rank(h)~plot~1, dat)
  wilcox.test(h~plot, dat, exact=FALSE, correct=FALSE)

## Wilcoxon signed rank test:
cctest(rank(abs(h-h5))~sign(h-h5)~0, subset(dat, h-h5 != 0))
#dat|> within(d<-h-h5)|> subset(d|0)|> with(rank(abs(d))~sign(d)~0)|> cctest()
  wilcox.test(h-h5 ~ 1, dat, exact=FALSE, correct=FALSE)

## Chi-squared test of independence:
cctest(ori~plot~1, dat, ~0)
cctest(ori~plot~1, as.data.frame(xtabs(~ori+plot,dat)), Freq~0)
  summary(xtabs(~ori+plot, dat, drop.unused.levels=TRUE))
  chisq.test(dat$ori, dat$plot, correct=FALSE)

## Score test in logistic regression (logit model, ...~1 only):
cctest(c15~x|y~1, dat, ~0)
  anova(glm(c15~1, binomial, dat, epsilon=1e-20, maxit=200),
    glm(c15~1+x+y, binomial, dat), test="Rao")

## Score test in multinomial logit model (...~1 only):
cctest(ori~x|y~1, dat, ~0)
  with(expand.grid(stringsAsFactors=FALSE,i=row.names(dat),j=levels(dat$ori)),
    anova(glm(ori==j ~ j+x+y, poisson, dat[i,], epsilon=1e-20, maxit=200),
      glm(ori==j ~ j*(x+y), poisson, dat[i,]), test="Rao"))

## Absolute values of (partial) correlation coefficients:
cctest(h~elev~1, dat)$est
  cor(dat$h, dat$elev)
cctest(h~elev~1|x|y, dat)$est
  cov2cor(estVar(lm(cbind(h,elev)~1+x+y, dat)))
cctest(h~x|y|elev~1, dat)$est^2
  summary(lm(h~1+x+y+elev, dat))$r.squared

## Canonical correlations:
cctest(h|h5|h10~x|y~1, dat)$est
  cancortest(dat[c("x","y")],dat[c("h","h5","h10")])$cor

## Linear discriminant analysis:
with(cctest(h|h5|h10~ori~1, dat, ~ori), y / sqrt(1-estimate^2)[col(y)])[1:7,]
  #predict(MASS::lda(ori~h+h5+h10,dat))$x[1:7,]

## Correspondence analysis:

```

```

cctest(ori~plot~1, as.data.frame(xtabs(~ori+plot,dat)), Freq~0)[1:2]
#MASS::corresp(~plot+ori, dat)

## Prediction in multivariate linear model:
with(cctest(h|h5|h10~1|x|y~0, dat, plot=="a"~0),
     x %% diag(estimate,ncol(x),ncol(y)) %% yinv)[1:7,]
predict(lm(cbind(h,h5,h10)~1+x+y, dat, subset=plot=="a"), dat)[1:7,]

## Other constructions:
cctest(ave(h,plot,FUN=rank)~ori~plot, aggregate(h~ori+plot,dat,mean))
friedman.test(h~ori|plot, aggregate(h~ori+plot,dat,mean))
cctest(ori:cut(elev,4)-cut(elev,4):ori~1~0, dat)
mcnemar.test(dat$ori, cut(dat$elev,4), correct=FALSE)
cctest(ave(abs(rank(h,ties.method="f")-mean(rank(h))),h)~plot~1, dat)
ansari.test(h~plot, dat, exact=FALSE)
cctest(ave((rank(h,ties.method="f")-mean(rank(h)))^2,h)~plot~1, dat)
mood.test(h~plot, dat)
cctest(qnorm(rank(-abs(h-ave(h,ori,FUN=median)))/(length(h)+1)/2)~ori~1, dat)
fligner.test(h~ori, dat)
cctest(h~diag(plot=="a",length(plot))~plot, dat)
var.test(h~plot, dat, alternative="greater")
cctest(ori:sum(Freq)/Freq~1~0, as.data.frame(xtabs(~ori,dat)),
      (if(all(Freq)) Freq^2/sum(Freq)/c(.2,.3,.4,.1) else stop())~0)
chisq.test(xtabs(~ori,dat), p=c(.2,.3,.4,.1))
with(cctest({h|h5|h10;0;0} ~ {h|h5|h10;diag(3)} ~ {1|0*x;0;0;0},
          c(dat,`~`=rbind, ~ {h|h5|h10;diag(3)} | {1|0*x;0;0;0}),
      list(estimate/(s<-sqrt((1-estimate^2)*df.residual)), t(xinv*s)))
prcomp(~h+h5+h10, dat)

## Handling of additional arguments and edge cases:
cctest(1:150~ori=="E"|ori=="W"~1, c(dat,`~`=:`~,`|`=~`|`))
anova(lm(1:150~ori=="E"|ori=="W", dat))
cctest(h~h5~h10~h5~1|x|y, dat, (ori=="E")^NA~1|x|y)
cctest(h~h10~1+x+y, dat, offset=h5, subset=ori=="E", stats=TRUE)
anova(lm(h~h5~x+y, dat, ori=="E"), lm(h~h5~x+y+I(h10-h5), dat, ori=="E"))
cctest(h~x~1, dat, weights=sapl/mean(sapl[sapl!=0])~1)
anova(lm(h~1, dat, weights=sapl), lm(h~1+x, dat, weights=sapl))
cctest(h*c15/c15~elev~1, dat[1:6,])[1:2]
cctest(I(h*c15/c15)~elev~1, dat[1:6,], stats=TRUE, na.action=na.exclude)[1:2]
scale(with(dat[1:6,], cbind(elev,h)*c15/c15))
cctest(c15~h~1, dat, tol=0.999*sqrt(1-cctest(h~1~0,dat)$est^2))
summary(lm(c15~h, dat, tol=0.999*sqrt(1-cctest(h~1~0,dat)$est^2)))
cctest(c15~h~1, dat, tol=1.001*sqrt(1-cctest(h~1~0,dat)$est^2))
summary(lm(c15~h, dat, tol=1.001*sqrt(1-cctest(h~1~0,dat)$est^2)))
cctest(NULL~NULL~NULL)
cctest(0~0~0)
anova(lm(0~0), lm(0~0+0))
cctest(h^0~1~0, dat)
cctest(1~1~0, dat, stats=TRUE)
anova(lm(h^0~0, dat), lm(h^0~0+1, dat))
## End(Not run)

```

Index

- * **htest**
 - cctest, 1
- * **multivariate**
 - cctest, 1
- ~, 2
- anova.mlm, 5
- cancor, 3, 5
- cbind, 2
- cctest, 1
- expansion, 2
- lm, 2–4
- model.frame, 4
- model.matrix, 4
- operators, 2
- sorted, 2
- terms.formula, 4