

coin: A Computational Framework for Conditional Inference

Torsten Hothorn¹, Kurt Hornik², Mark van de Wiel³
and Achim Zeileis²

¹Institut für Medizininformatik, Biometrie und Epidemiologie
Friedrich-Alexander-Universität Erlangen-Nürnberg
Waldstraße 6, D-91054 Erlangen, Germany
`Torsten.Hothorn@R-project.org`

²Institut für Statistik und Mathematik, Wirtschaftsuniversität Wien
Augasse 2-6, A-1090 Wien, Austria
`Kurt.Hornik@R-project.org`
`Achim.Zeileis@R-project.org`

³ Department of Mathematics and Computer Science
Eindhoven University of Technology
HG 9.25, P.O. Box 513
5600 MB Eindhoven, The Netherlands
`markvdw@win.tue.nl`

1 Introduction

The `coin` package implements a unified approach for conditional inference procedures commonly known as *permutation tests*. The theoretical basis of design and implementation is the unified framework for permutation tests given by [Strasser and Weber \(1999\)](#). For a very flexible formulation of multivariate linear statistics, [Strasser and Weber \(1999\)](#) derived the conditional expectation and covariance of the conditional (permutation) distribution as well as the multivariate limiting distribution.

Conditional counterparts of a large amount of classical (unconditional) test procedures for continuous, categorical and censored data are part of this framework, for example the Cochran-Mantel-Haenszel test for independence in general contingency tables, linear association tests for ordered categorical data, linear rank tests and multivariate permutation tests.

The conceptual framework of permutation tests by [Strasser and Weber \(1999\)](#) for arbitrary problems is available via the generic `independence_test`. Because

convenience functions for the most prominent problems are available, users will not have to use this extremely flexible procedure. Currently, the conditional variants of the following test procedures are available:

<code>oneway_test</code>	two- and K -sample permutation test
<code>wilcox_test</code>	Wilcoxon-Mann-Whitney rank sum test
<code>normal_test</code>	van der Waerden normal quantile test
<code>median_test</code>	Median test
<code>kruskal_test</code>	Kruskal-Wallis test
<code>ansari_test</code>	Ansari-Bradley test
<code>fligner_test</code>	Fligner-Killeen test
<code>chisq_test</code>	Pearson's χ^2 test
<code>cmh_test</code>	Cochran-Mantel-Haenszel test
<code>lbl_test</code>	linear-by-linear association test
<code>surv_test</code>	two- and K -sample logrank test
<code>maxstat_test</code>	maximally selected statistics
<code>spearman_test</code>	Spearman's test
<code>friedman_test</code>	Friedman test
<code>wilcoxsign_test</code>	Wilcoxon-Signed-Rank test
<code>mh_test</code>	marginal homogeneity test.

Those convenience functions essentially perform a certain transformation of the data, e.g., a rank transformation, and call `independence_test` for the computation of linear statistics, expectation and covariance and the test statistic as well as their null distribution. The exact null distribution can be approximated either by the asymptotic distribution or via conditional Monte-Carlo for all test procedures, the exact null distribution is available for special cases. Moreover, all test procedures allow for the specification of blocks for stratification.

2 Permutation Tests

In the following we assume that we are provided with n observations

$$(\mathbf{Y}_i, \mathbf{X}_i, w_i, b_i), \quad i = 1, \dots, n.$$

The variables \mathbf{Y} and \mathbf{X} from sample spaces \mathcal{Y} and \mathcal{X} may be measured at arbitrary scales and may be multivariate as well. In addition to those measurements, case weights w and a factor b coding blocks may be available. For the sake of simplicity, we assume $w_i = 1$ and $b_i = 0$ for all observations $i = 1, \dots, n$ for the moment.

We are interested in testing the null hypothesis of independence of \mathbf{Y} and \mathbf{X}

$$H_0 : D(\mathbf{Y}|\mathbf{X}) = D(\mathbf{Y})$$

against arbitrary alternatives. [Strasser and Weber \(1999\)](#) suggest to derive scalar test statistics for testing H_0 from multivariate linear statistics of the

form

$$\mathbf{T} = \text{vec} \left(\sum_{i=1}^n w_i g(\mathbf{X}_i) h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n))^\top \right) \in \mathbb{R}^{pq}. \quad (1)$$

Here, $g : \mathcal{X} \rightarrow \mathbb{R}^p$ is a transformation of the \mathbf{X} measurements and the *influence function* $h : \mathcal{Y} \times \mathcal{Y}^n \rightarrow \mathbb{R}^q$ depends on the responses $(\mathbf{Y}_1, \dots, \mathbf{Y}_n)$ in a permutation symmetric way. We will give specific examples how to choose g and h later on.

The distribution of \mathbf{T} depends on the joint distribution of \mathbf{Y} and \mathbf{X} , which is unknown under almost all practical circumstances. At least under the null hypothesis one can dispose of this dependency by fixing $\mathbf{X}_1, \dots, \mathbf{X}_n$ and conditioning on all possible permutations S of the responses $\mathbf{Y}_1, \dots, \mathbf{Y}_n$. This principle leads to test procedures known as *permutation tests*.

The conditional expectation $\mu \in \mathbb{R}^{pq}$ and covariance $\Sigma \in \mathbb{R}^{pq \times pq}$ of \mathbf{T} under H_0 given all permutations $\sigma \in S$ of the responses are derived by [Strasser and Weber \(1999\)](#):

$$\begin{aligned} \mu &= \mathbb{E}(\mathbf{T}|S) = \text{vec} \left(\left(\sum_{i=1}^n w_i g(\mathbf{X}_i) \right) \mathbb{E}(h|S)^\top \right), \\ \Sigma &= \mathbb{V}(\mathbf{T}|S) \\ &= \frac{\mathbf{w}_\cdot}{\mathbf{w}_\cdot - 1} \mathbb{V}(h|S) \otimes \left(\sum_i w_i g(\mathbf{X}_i) \otimes w_i g(\mathbf{X}_i)^\top \right) \\ &\quad - \frac{1}{\mathbf{w}_\cdot - 1} \mathbb{V}(h|S) \otimes \left(\sum_i w_i g(\mathbf{X}_i) \right) \otimes \left(\sum_i w_i g(\mathbf{X}_i) \right)^\top \end{aligned} \quad (2)$$

where $\mathbf{w}_\cdot = \sum_{i=1}^n w_i$ denotes the sum of the case weights, and \otimes is the Kronecker product. The conditional expectation of the influence function is

$$\mathbb{E}(h|S) = \mathbf{w}_\cdot^{-1} \sum_i w_i h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) \in \mathbb{R}^q$$

with corresponding $q \times q$ covariance matrix

$$\begin{aligned} \mathbb{V}(h|S) &= \mathbf{w}_\cdot^{-1} \sum_i w_i (h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) - \mathbb{E}(h|S)) \\ &\quad (h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) - \mathbb{E}(h|S))^\top. \end{aligned}$$

Having the conditional expectation and covariance at hand we are able to standardize a linear statistic $\mathbf{T} \in \mathbb{R}^{pq}$ of the form (1). Univariate test statistics c mapping an observed linear statistic $\mathbf{t} \in \mathbb{R}^{pq}$ into the real line can be of arbitrary form. An obvious choice is the maximum of the absolute values of the standardized linear statistic

$$c_{\max}(\mathbf{t}, \mu, \Sigma) = \max \left| \frac{\mathbf{t} - \mu}{\text{diag}(\Sigma)^{1/2}} \right|$$

utilizing the conditional expectation μ and covariance matrix Σ . The application of a quadratic form $c_{\text{quad}}(\mathbf{t}, \mu, \Sigma) = (\mathbf{t} - \mu)\Sigma^+(\mathbf{t} - \mu)^\top$ is one alternative, although computationally more expensive because the Moore-Penrose inverse Σ^+ of Σ is involved.

The definition of one- and two-sided p -values used for the computations in the `coin` package is

$$\begin{aligned} P(c(\mathbf{T}, \mu, \Sigma) &\leq c(\mathbf{t}, \mu, \Sigma)) && \text{(less)} \\ P(c(\mathbf{T}, \mu, \Sigma) &\geq c(\mathbf{t}, \mu, \Sigma)) && \text{(greater)} \\ P(|c(\mathbf{T}, \mu, \Sigma)| &\leq |c(\mathbf{t}, \mu, \Sigma)|) && \text{(two-sided).} \end{aligned}$$

Note that for quadratic forms only two-sided p -values are available and that in the one-sided case maximum type test statistics are replaced by

$$\min \left(\frac{\mathbf{t} - \mu}{\text{diag}(\Sigma)^{1/2}} \right) \quad \text{(less)} \quad \text{and} \quad \max \left(\frac{\mathbf{t} - \mu}{\text{diag}(\Sigma)^{1/2}} \right) \quad \text{(greater)}.$$

The conditional distribution and thus the p -value of the statistics $c(\mathbf{t}, \mu, \Sigma)$ can be computed in several different ways. For some special forms of the linear statistic, the exact distribution of the test statistic is trackable. For two-sample problems, the shift-algorithm by [Streitberg and Röhmel \(1986\)](#) and [Streitberg and Röhmel \(1987\)](#) and the split-up algorithm by [van de Wiel \(2001\)](#) are implemented as part of the package. Conditional Monte-Carlo procedures can be used to approximate the exact distribution. [Strasser and Weber \(1999\)](#) proved (Theorem 2.3) that the conditional distribution of linear statistics \mathbf{T} with conditional expectation μ and covariance Σ tends to a multivariate normal distribution with parameters μ and Σ as $n, s \rightarrow \infty$. Thus, the asymptotic conditional distribution of test statistics of the form c_{max} is normal and can be computed directly in the univariate case ($pq = 1$) or approximated by means of quasi-randomized Monte-Carlo procedures in the multivariate setting ([Genz, 1992](#)). For quadratic forms c_{quad} which follow a χ^2 distribution with degrees of freedom given by the rank of Σ (Theorem 6.20, [Rasch, 1995](#)), exact probabilities can be computed efficiently.

3 Illustrations and Applications

The main workhorse `independence_test` essentially allows for the specification of \mathbf{Y}, \mathbf{X} and b through a formula interface of the form `y ~ x | b`, weights can be defined by a formula with one variable on the right hand side only. Four additional arguments are available for the specification of the transformation g (`xtrans`), the influence function h (`ytrans`), the form of the test statistic c (`teststat`) and the null distribution (`distribution`).

Independent K -Sample Problems. When we want to compare the distribution of an univariate qualitative response \mathbf{Y} in K groups given by a factor \mathbf{X}

at K levels, the transformation g is the dummy matrix coding the groups and h is either the identity transformation or a some form of rank transformation.

For example, the Kruskal-Wallis test may be computed as follows (example taken from [Hollander and Wolfe, 1999](#), Table 6.3, page 200):

```
> library(coin)
> YOY <- data.frame(length = c(46, 28, 46, 37, 32,
+   41, 42, 45, 38, 44, 42, 60, 32, 42, 45, 58, 27,
+   51, 42, 52, 38, 33, 26, 25, 28, 28, 26, 27, 27,
+   27, 31, 30, 27, 29, 30, 25, 25, 24, 27, 30),
+   site = factor(c(rep("I", 10), rep("II", 10),
+     rep("III", 10), rep("IV", 10))))
> it <- independence_test(length ~ site, data = YOY,
+   ytrafo = function(data) trafo(data, numeric_trafo = rank),
+   teststat = "quadtype")
> it
```

Asymptotic General Independence Test

```
data: length by groups I, II, III, IV
T = 22.8524, df = 3, p-value = 4.335e-05
```

The linear statistic \mathbf{T} is the sum of the ranks in each group and can be extracted via

```
> statistic(it, "linear")
```

```
      [,1]
I       278
II      307
III     119
IV     116
```

Note that `statistic(..., "linear")` currently returns the linear statistic in matrix form, i.e.

$$\sum_{i=1}^n w_i g(\mathbf{X}_i) h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n))^{\top} \in \mathbb{R}^{p \times q}.$$

The conditional expectation and covariance are available from

```
> expectation(it)

[1] 205 205 205 205

> covariance(it)
```

```

      [,1]      [,2]      [,3]      [,4]
[1,] 1019.0385 -339.6795 -339.6795 -339.6795
[2,] -339.6795 1019.0385 -339.6795 -339.6795
[3,] -339.6795 -339.6795 1019.0385 -339.6795
[4,] -339.6795 -339.6795 -339.6795 1019.0385

```

and the standardized linear statistic $(\mathbf{T} - \mu)\text{diag}(\Sigma)^{-1/2}$ is

```
> statistic(it, "standardized")
```

```

      [,1]
I      2.286797
II     3.195250
III   -2.694035
IV    -2.788013

```

Since a quadratic form of the test statistic was requested via `teststat = "quadtype"`, the test statistic is

```
> statistic(it)

[1] 22.85242
```

By default, the asymptotic distribution of the test statistic is computed, the p -value is

```
> pvalue(it)

[1] 4.334659e-05
```

Life is much simpler with convenience functions very similar to those available in package `stats` for a long time. The exact null distribution of the Kruskal-Wallis test can be approximated by 9999 Monte-Carlo replications via

```
> kw <- kruskal_test(length ~ site, data = YOY, distribution = approximate(B = 9999))
> kw
```

Approximative Kruskal-Wallis Test

```
data: length by groups I, II, III, IV
T = 22.8524, p-value < 2.2e-16
```

with p -value (and 99% confidence interval) of

```
> pvalue(kw)

[1] 0
99 percent confidence interval:
 0.0000000000 0.0005297444
```

Of course it is possible to choose a c_{\max} type test statistic instead of a quadratic form.

Independence in Contingency Tables. Independence in general two- or three-dimensional contingency tables can be tested by the Cochran-Mantel-Haenszel test. Here, both g and h are dummy matrices (example data from [Agresti, 2002](#), Table 7.8, page 288):

```
> data(jobsatisfaction, package = "coin")
> it <- cmh_test(jobsatisfaction)
> it
```

Asymptotic Generalised Cochran-Mantel-Haenszel Test

```
data: Job.Satisfaction by
      groups <5000, 5000-15000, 15000-25000, >25000
      stratified by Gender
T = 10.2001, df = 9, p-value = 0.3345
```

The standardized contingency table allowing for an inspection of the deviation from the null hypothesis of independence of income and jobsatisfaction (stratified by gender) is

```
> statistic(it, "standardized")
```

	Very Dissatisfied	A Little Dissatisfied
<5000	1.3112789	0.69201053
5000-15000	0.6481783	0.83462550
15000-25000	-1.0958361	-1.50130926
>25000	-1.0377629	-0.08983052

	Moderately Satisfied	Very Satisfied
<5000	-0.2478705	-0.9293458
5000-15000	0.5175755	-1.6257547
15000-25000	0.2361231	1.4614123
>25000	-0.5946119	1.2031648

Ordered Alternatives. Of course, both job satisfaction and income are ordered variables. When \mathbf{Y} is measured at J levels and \mathbf{X} at K levels, \mathbf{Y} and \mathbf{X} are associated with score vectors $\xi \in \mathbb{R}^J$ and $\gamma \in \mathbb{R}^K$, respectively. The linear statistic is now a linear combination of the linear statistic \mathbf{T} of the form

$$\mathbf{MT} = \text{vec} \left(\sum_{i=1}^n w_i \gamma^\top g(\mathbf{X}_i) (\xi^\top h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)))^\top \right) \in \mathbb{R} \text{ with } \mathbf{M} = \xi \otimes \gamma.$$

By default, scores are $\xi = 1, \dots, J$ and $\gamma = 1, \dots, K$.

```
> lbl_test(jobsatisfaction)
```

Asymptotic Linear-by-Linear Association Test

```
data: Job.Satisfaction (ordered) by
```

```

      groups <5000 < 5000-15000 < 15000-25000 < >25000
      stratified by Gender
T = 6.6235, df = 1, p-value = 0.01006

```

The scores ξ and γ can be specified to the linear-by-linear association test via a list those names correspond to the variable names

```

> lbl_test(jobsatisfaction, scores = list(Job.Satisfaction = c(1,
+   3, 4, 5), Income = c(3, 10, 20, 35)))

```

Asymptotic Linear-by-Linear Association Test

```

data:  Job.Satisfaction (ordered) by
      groups <5000 < 5000-15000 < 15000-25000 < >25000
      stratified by Gender
T = 6.1563, df = 1, p-value = 0.01309

```

Incomplete Randomised Blocks. [Rayner and Best \(2001\)](#), Chapter 7, discuss the application of Durbin's test to data from sensoric experiments, where incomplete block designs are common. As an example, data from taste-testing on ten dried eggs where mean scores for off-flavour from seven judges are given and one wants to assess whether there is any difference in the scores between the ten egg samples. The sittings are a block variable which can be added to the formula via '|'.

```

> egg_data <- data.frame(scores = c(9.7, 8.7, 5.4,
+   5, 9.6, 8.8, 5.6, 3.6, 9, 7.3, 3.8, 4.3, 9.3,
+   8.7, 6.8, 3.8, 10, 7.5, 4.2, 2.8, 9.6, 5.1, 4.6,
+   3.6, 9.8, 7.4, 4.4, 3.8, 9.4, 6.3, 5.1, 2, 9.4,
+   9.3, 8.2, 3.3, 8.7, 9, 6, 3.3, 9.7, 6.7, 6.6,
+   2.8, 9.3, 8.1, 3.7, 2.6, 9.8, 7.3, 5.4, 4, 9,
+   8.3, 4.8, 3.8, 9.3, 8.3, 6.3, 3.8), sitting = factor(rep(c(1:15),
+   rep(4, 15))), product = factor(c(1, 2, 4, 5,
+   2, 3, 6, 10, 2, 4, 6, 7, 1, 3, 5, 7, 1, 4, 8,
+   10, 2, 7, 8, 9, 2, 5, 8, 10, 5, 7, 9, 10, 1,
+   2, 3, 9, 4, 5, 6, 9, 1, 6, 7, 10, 3, 4, 9, 10,
+   1, 6, 8, 9, 3, 4, 7, 8, 3, 5, 6, 8)))
> independence_test(scores ~ product | sitting, data = egg_data,
+   teststat = "quadtype", ytrafo = function(data) trafo(data,
+   numeric_trafo = rank, block = egg_data$sitting))

```

Asymptotic General Independence Test

```

data:  scores by
      groups 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
      stratified by sitting
T = 39.12, df = 9, p-value = 1.096e-05

```


and the Monte-Carlo p -value can be computed via

```
> pvalue(independence_test(scores ~ product | sitting,
+   data = egg_data, teststat = "quadtype", ytrafo = function(data) trafo(data,
+   numeric_trafo = rank, block = egg_data$sitting),
+   distribution = approximate(B = 19999)))

[1] 0
99 percent confidence interval:
 0.000000000 0.000264894
```

If we assume that the products are ordered, the Page test is appropriate and can be computed as follows

```
> independence_test(scores ~ product | sitting, data = egg_data,
+   scores = list(product = 1:10), ytrafo = function(data) trafo(data,
+   numeric_trafo = rank, block = egg_data$sitting))
```

Asymptotic General Independence Test

```
data: scores by
      groups 1 < 2 < 3 < 4 < 5 < 6 < 7 < 8 < 9 < 10
      stratified by sitting
T = 6.2166, p-value = 5.081e-10
```

Multiple Tests. One may be interested in testing multiple hypotheses simultaneously, either by using a linear combination of the linear statistic **KT**, or by specifying multivariate variables **Y** and / or **X**. For example, all pair comparisons may be implemented via

```
> if (require(multcomp)) {
+   it <- independence_test(length ~ site, data = YOY,
+   xtrafo = function(data) trafo(data, factor_trafo = function(x) model.matrix(~x -
+   1) %*% t(contrMat(table(x), "Tukey"))),
+   teststat = "max", distribution = approximate(B = 9999))
+   print(pvalue(it))
+   print(pvalue(it, method = "single-step"))
+ }
```

```
[1] 0.00010001
99 percent confidence interval:
 5.013042e-07 7.428484e-04
```

```
      [,1]
II-I    0.64726473
III-I   0.03680368
IV-I    0.02110211
III-II  0.00010001
```

```
IV-II 0.00010001
IV-III 0.99799980
```

When either g or h are multivariate, single-step adjusted p -values based on maximum-type statistics are computed as described in [Westfall and Young \(1993\)](#), algorithm 2.5 (page 47) and equation (2.8), page 50. Note that for the example shown above only the *minimum* p -value is adjusted appropriately because the subset pivotality condition is violated, i.e., the distribution of the test statistics under the complete null-hypothesis of no treatment effect of **site** is the basis of all adjustments instead of the corresponding partial null-hypothesis.

Another important application are simultaneous tests for many response variables. This problem frequently occurs in microarray expression studies and we shall have a look at an artificial example: 100 variables (from a normal distribution) are to be tested in a one-way classification with $n = 40$ observations. Only the first variable shows a difference and we are interested in both a global test and the adjusted p -values. Here, the example is formulated within the Biobase framework:

```
> if (require(Biobase)) {
+   p <- 100
+   pd <- new("phenoData", pData = data.frame(group = gl(2,
+     20)), varLabels = list(group = c("1", "2")))
+   exprs <- matrix(rnorm(p * 40), nrow = p)
+   exprs[1, 1:20] <- exprs[1, 1:20] + 1.5
+   ex <- new("exprSet", exprs = exprs, phenoData = pd)
+   it <- independence_test(group ~ ., data = ex,
+     distribution = approximate(B = 1000))
+   print(pvalue(it))
+   print(which(pvalue(it, method = "step-down") <
+     0.05))
+ }

[1] 0.002
99 percent confidence interval:
 0.0001035410 0.0092401306

[1] 1
```

4 Quality Assurance

The test procedures implemented in package **coin** are continuously checked against results obtained by the corresponding implementations in package **stats** (where available). In addition, the test statistics and exact, approximative and asymptotic p -values for data examples given in the **StatXact-6** user manual ([Mehta and Patel, 2003](#)) are compared with the results reported in the **StatXact-6** manual. For details on the test procedures we refer to the R-transcript files in directory **coin/tests**.

References

- Alan Agresti. *Categorical Data Analysis*. John Wiley & Sons, Hoboken, New Jersey, 2nd edition, 2002. [7](#)
- Alan Genz. Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1:141–149, 1992. [4](#)
- Myles Hollander and Douglas A. Wolfe. *Nonparametric statistical inference*. John Wiley & Sons, New York, 2nd edition, 1999. [5](#)
- Cyrus R. Mehta and Nitin R. Patel. *StatXact-6: Statistical Software for Exact Nonparametric Inference*. Cytel Software Cooperation, Cambridge, USA, 2003. [10](#)
- Dieter Rasch. *Mathematische Statistik*. Johann Ambrosius Barth Verlag, Heidelberg, Leipzig, 1995. [4](#)
- J. C. W. Rayner and D. J. Best. *A contingency table approach to nonparametric testing*. Chapman & Hall, New York, 2001. [8](#)
- Helmut Strasser and Christian Weber. On the asymptotic theory of permutation statistics. *Mathematical Methods of Statistics*, 8:220–250, 1999. [1](#), [2](#), [3](#), [4](#)
- Bernd Streitberg and Joachim Röhmel. Exact distributions for permutations and rank tests: An introduction to some recently published algorithms. *Statistical Software Newsletter*, 12(1):10–17, 1986. ISSN 1609-3631. [4](#)
- Bernd Streitberg and Joachim Röhmel. Exakte Verteilungen für Rang- und Randomisierungstests im allgemeinen c -Stichprobenfall. *EDV in Medizin und Biologie*, 18(1):12–19, 1987. [4](#)
- Mark A. van de Wiel. The split-up algorithm: a fast symbolic method for computing p-values of rank statistics. *Computational Statistics*, 16:519–538, 2001. [4](#)
- Peter H. Westfall and S. Stanley Young. *Resampling-based Multiple Testing*. John Wiley & Sons, New York, 1993. [10](#)