

Package ‘YatchewTest’

April 26, 2024

Title Yatchew (1997), De Chaisemartin & D'Haultfoeuille (2024)
Linearity Test

Version 1.0.2

Maintainer Diego Ciccia <diego.ciccia@sciencespo.fr>

Description Test of linearity originally proposed by Yatchew (1997) <[doi:10.1016/S0165-1765\(97\)00218-8](https://doi.org/10.1016/S0165-1765(97)00218-8)> and improved by de Chaisemartin & D'Haultfoeuille (2024) <[doi:10.2139/ssrn.4284811](https://doi.org/10.2139/ssrn.4284811)> to be robust under heteroskedasticity.

License MIT + file LICENSE

Imports Rcpp, ggplot2

LinkingTo Rcpp

Author Diego Ciccia [aut, cre],
Felix Knau [aut],
Doulo Sow [aut],
Clément de Chaisemartin [aut],
Xavier D'Haultfoeuille [aut]

Encoding UTF-8

RoxygenNote 7.2.3

Suggests testthat (>= 3.0.0)

Config/testthat/edition 3

NeedsCompilation yes

Repository CRAN

Date/Publication 2024-04-26 10:30:03 UTC

R topics documented:

yatchew_test	2
yatchew_test.data.frame	2

Index	5
--------------	----------

yatchew_test	<i>Main function</i>
--------------	----------------------

Description

Test of Linearity of a Conditional Expectation Function (Yatchew, 1997; de Chaisemartin and D'Haultfoeuille, 2024)

Usage

```
yatchew_test(data, ...)
```

Arguments

data	A data object.
...	Undocumented.

Value

Method dispatch depending on the data object class.

yatchew_test.data.frame	<i>General yatchew_test method for unclassified dataframes</i>
-------------------------	--

Description

General yatchew_test method for unclassified dataframes

Usage

```
## S3 method for class 'data.frame'
yatchew_test(data, Y, D, het_robust = FALSE, path_plot = FALSE, ...)
```

Arguments

data	(data.frame) A dataframe.
Y	(char) Dependent variable.
D	(char) Independent variable.
het_robust	(logical) If FALSE, the test is performed under the assumption of homoskedasticity (Yatchew, 1997). If TRUE, the test is performed using the heteroskedasticity-robust test statistic proposed by de Chaisemartin and D'Haultfoeuille (2024).
path_plot	(logical) if TRUE and D has length 2, the assigned object will include a plot of the sequence of (D_{1i}, D_{2i}) s that minimizes the euclidean distance between each pair of consecutive observations (see Overview for further details).
...	Undocumented.

Value

A list with test results.

Overview

This program implements the linearity test proposed by Yatchew (1997) and its heteroskedasticity-robust version proposed by de Chaisemartin and D’Haultfoeuille (2024). In this overview, we sketch the intuition behind the two tests, as to motivate the use of the package and its options. Please refer to Yatchew (1997) and Section 3 of de Chaisemartin and D’Haultfoeuille (2024) for further details.

Yatchew (1997) proposes a useful extension of the test with multiple independent variables. The program implements this extension when the `D` argument has length > 1 . It should be noted that the power and consistency of the test in the multivariate case are not backed by proven theoretical results. We implemented this extension to allow for testing and exploratory research. Future theoretical exploration of the multivariate test will depend on the demand and usage of the package.

Univariate Yatchew Test:

Let Y and D be two random variables. Let $m(D) = E[Y|D]$. The null hypothesis of the test is that $m(D) = \alpha_0 + \alpha_1 D$ for two real numbers α_0 and α_1 . This means that, under the null, $m(\cdot)$ is linear in D . The outcome variable can be decomposed as $Y = m(D) + \varepsilon$, with $E[\varepsilon|D] = 0$ and $\Delta Y = \Delta \varepsilon$ for $\Delta D \rightarrow 0$. In a dataset with N i.i.d. realisations of (Y, D) , one can test this hypothesis as follows:

1. sort the dataset by D ;
2. denote the corresponding observations by $(Y_{(i)}, D_{(i)})$, with $i \in \{1, \dots, N\}$;
3. approximate $\hat{\sigma}_{\text{diff}}^2$, i.e. the variance of the first differenced residuals $\varepsilon_{(i)} - \varepsilon_{(i-1)}$, by the variance of $Y_{(i)} - Y_{(i-1)}$;
4. compute $\hat{\sigma}_{\text{lin}}^2$, i.e. the variance of the residuals from an OLS regression of Y on D .

Heuristically, the validity of step (3) derives from the fact that $Y_{(i)} - Y_{(i-1)} = m(D_{(i)}) - m(D_{(i-1)}) + \varepsilon_{(i)} - \varepsilon_{(i-1)}$ and the first difference term is close to zero for $D_{(i)} \approx D_{(i-1)}$. Sorting at step (1) ensures that consecutive $D_{(i)}$ s are as close as possible, and when the sample size goes to infinity the distance between consecutive observations goes to zero. Then, Yatchew (1997) shows that under homoskedasticity and regularity conditions

$$T := \sqrt{G} \left(\frac{\hat{\sigma}_{\text{lin}}^2}{\hat{\sigma}_{\text{diff}}^2} - 1 \right) \xrightarrow{d} \mathcal{N}(0, 1)$$

Then, one can reject the linearity of $m(\cdot)$ with significance level α if $T > \Phi(1 - \alpha)$.

If the homoskedasticity assumption fails, this test leads to overrejection. De Chaisemartin and D’Haultfoeuille (2024) propose a heteroskedasticity-robust version of the test statistic above. This version of the Yatchew (1997) test can be implemented by running the command with the option `het_robust = TRUE`.

Multivariate Yatchew Test:

Let \mathbf{D} be a vector of K random variables. Let $g(\mathbf{D}) = E[Y|\mathbf{D}]$. Denote with $\|\cdot, \cdot\|$ the Euclidean distance between two vectors. The null hypothesis of the multivariate test is $g(\mathbf{D}) = \alpha_0 + A'\mathbf{D}$, with $A = (\alpha_1, \dots, \alpha_K)$, for $K + 1$ real numbers $\alpha_0, \alpha_1, \dots, \alpha_K$. This means that, under the null, $g(\cdot)$ is linear in \mathbf{D} . Following the same logic as the univariate case, in a dataset with N i.i.d. realisations of (Y, \mathbf{D}) we can approximate the first difference $\Delta \varepsilon$ by ΔY valuing $g(\cdot)$ between

consecutive observations. The program runs a nearest neighbor algorithm to find the sequence of observations such that the Euclidean distance between consecutive positions is minimized. The algorithm has been programmed in C++ and it has been integrated in R thanks to the Rcpp library. The program follows a very simple nearest neighbor approach:

1. collect all the Euclidean distances between all the possible unique pairs of rows in \mathbf{D} in the matrix M , where $M_{n,m} = \|\mathbf{D}_n, \mathbf{D}_m\|$ with $n, m \in \{1, \dots, N\}$;
2. setup the queue to $Q = \{1, \dots, N\}$, the (empty) path vector $I = \{\}$ and the starting index $i = 1$;
3. remove i from Q and find the column index j of M such that $M_{i,j} = \min_{c \in Q} M_{i,c}$;
4. append j to I and start again from step 3 with $i = j$ until Q is empty.

To improve efficiency, the program collects only the $N(N - 1)/2$ Euclidean distances corresponding to the lower triangle of matrix M and chooses j such that $M_{i,j} = \min_{c \in Q} 1\{c < i\}M_{i,c} + 1\{c > i\}M_{c,i}$. The output of the algorithm, i.e. the vector I , is a sequence of row numbers such that the distance between the corresponding rows \mathbf{D}_i s is minimized. The program also uses two refinements suggested in Appendix A of Yatchew (1997):

- The entries in \mathbf{D} are normalized in $[0, 1]$;
- The algorithm is applied to sub-cubes, i.e. partitions of the $[0, 1]^K$ space, and the full path is obtained by joining the extrema of the subpaths.

By convention, the program computes $(2^{\lceil \log_{10} N \rceil})^K$ subcubes, where each univariate partition is defined by grouping observations in $2^{\lceil \log_{10} N \rceil}$ quantile bins. If $K = 2$, the user can visualize in a ggplot graph the exact path across the normalized \mathbf{D}_i s by running the command with the option `path_plot = TRUE`.

Once the dataset is sorted by I , the program resumes from step (2) of the univariate case.

Contacts

If you wish to inquire about the functionalities of this package or to report bugs/suggestions, feel free to post your question in the Issues section of the [yatchew_test GitHub repository](#).

References

de Chaisemartin, C., d'Haultfoeuille, X. (2024). Two-way Fixed Effects and Difference-in-Difference Estimators in Heterogeneous Adoption Designs.

Yatchew, A. (1997). An elementary estimator of the partial linear model.

Examples

```
df <- as.data.frame(matrix(NA, nrow = 1E3, ncol = 0))
df$x <- rnorm(1E3)
df$b <- runif(1E3)
df$y <- 2 + df$b * df$x
yatchew_test(data = df, Y = "y", D = "x")
```

Index

yatchew_test, [2](#)

yatchew_test.data.frame, [2](#)