

Package ‘binsreg’

July 6, 2023

Type Package

Title Binscatter Estimation and Inference

Date 2023-07-06

Version 1.0

Author Matias D. Cattaneo, Richard K. Crump, Max H. Farrell, Yingjie Feng

Maintainer Yingjie Feng <fengyingjiepku@gmail.com>

Description Provides tools for statistical analysis using the binscatter methods developed by Cattaneo, Crump, Farrell and Feng (2023a) <[arXiv:1902.09608](#)>, Cattaneo, Crump, Farrell and Feng (2023b) <https://nppackages.github.io/references/Cattaneo-Crump-Farrell-Feng_2023_NonlinearBinscatter.pdf> and Cattaneo, Crump, Farrell and Feng (2023c) <[arXiv:1902.09615](#)>. Binscatter provides a flexible way of describing the relationship between two variables based on partitioning/binning of the independent variable of interest. binsreg(), binsqreg() and binsglm() implement binscatter least squares regression, quantile regression and generalized linear regression respectively, with particular focus on constructing binned scatter plots. They also implement robust (pointwise and uniform) inference of regression functions and derivatives thereof. binstest() implements hypothesis testing procedures for parametric functional forms of and nonparametric shape restrictions on the regression function. binspwc() implements hypothesis testing procedures for pairwise group comparison of binscatter estimators. binsregselect() implements data-driven procedures for selecting the number of bins for binscatter estimation. All the commands allow for covariate adjustment, smoothness restrictions and clustering.

Depends R (>= 3.1)

License GPL-2

Encoding UTF-8

Imports ggplot2, sandwich, quantreg, splines, matrixStats

Roxygen list(old_usage = TRUE)

RoxygenNote 7.2.3

R topics documented:

| | |
|---------------------------|----|
| binsreg-package | 2 |
| binsglm | 3 |
| binspwc | 8 |
| binsqreg | 12 |
| binsreg | 18 |
| binsregselect | 24 |
| binstest | 27 |

Description

Binscatter provides a flexible, yet parsimonious way of visualizing and summarizing large data sets and has been a popular methodology in applied microeconomics and other social sciences. The binsreg package provides tools for statistical analysis using the binscatter methods developed in Cattaneo, Crump, Farrell and Feng (2023a) and Cattaneo, Crump, Farrell and Feng (2023b). binsreg implements binscatter least squares regression with robust inference and plots, including curve estimation, pointwise confidence intervals and uniform confidence band. binsqreg implements binscatter quantile regression with robust inference and plots, including curve estimation, pointwise confidence intervals and uniform confidence band. binsglm implements binscatter generalized linear regression with robust inference and plots, including curve estimation, pointwise confidence intervals and uniform confidence band. binstest implements binscatter-based hypothesis testing procedures for parametric specifications of and shape restrictions on the unknown function of interest. binspwc implements hypothesis testing procedures for pairwise group comparison of binscatter estimators. binsregselect implements data-driven number of bins selectors for binscatter implementation using either quantile-spaced or evenly-spaced binning/partitioning. All the commands allow for covariate adjustment, smoothness restrictions, and clustering, among other features.

The companion software article, Cattaneo, Crump, Farrell and Feng (2023c), provides further implementation details and empirical illustration. For related Stata, R and Python packages useful for nonparametric data analysis and statistical inference, visit <https://nppackages.github.io/>.

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepku@gmail.com>.

References

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: **On Binscatter**. Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: **Nonlinear Binscatter Methods**. Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: **Binscatter Regressions**. Working Paper.

binsglm

Data-Driven Binscatter Generalized Linear Regression with Robust Inference Procedures and Plots

Description

`binsglm` implements `binscatter` generalized linear regression with robust inference procedures and plots, following the results in [Cattaneo, Crump, Farrell and Feng \(2023a\)](#) and [Cattaneo, Crump, Farrell and Feng \(2023b\)](#). `Binscatter` provides a flexible way to describe the relationship between two variables, after possibly adjusting for other covariates, based on partitioning/binning of the independent variable of interest. The main purpose of this function is to generate binned scatter plots with curve estimation with robust pointwise confidence intervals and uniform confidence band. If the binning scheme is not set by the user, the companion function `binsregselect` is used to implement `binscatter` in a data-driven way. Hypothesis testing about the function of interest can be conducted via the companion function `binstest`.

Usage

```
binsglm(y, x, w = NULL, data = NULL, at = NULL, family = gaussian(),
  deriv = 0, nolink = F, dots = NULL, dotsgrid = 0, dotsgridmean = T,
  line = NULL, linegrid = 20, ci = NULL, cigrid = 0, cigridmean = T,
  cb = NULL, cbgrid = 20, polyreg = NULL, polyreggrid = 20,
  polyregcigrid = 0, by = NULL, bycolors = NULL, bysymbols = NULL,
  bylpatterns = NULL, legendTitle = NULL, legendoff = F, nbins = NULL,
  binspos = "qs", binsmethod = "dpi", nbinsrot = NULL, pselect = NULL,
  sselect = NULL, samebinsby = F, randcut = NULL, nsims = 500,
  simsgrid = 20, simsseed = NULL, vce = "HC1", cluster = NULL,
  asyvar = F, level = 95, noplot = F, dfcheck = c(20, 30),
  masspoints = "on", weights = NULL, subset = NULL, plotxrange = NULL,
  plotyrange = NULL, ...)
```

Arguments

| | |
|---------------------|---|
| <code>y</code> | outcome variable. A vector. |
| <code>x</code> | independent variable of interest. A vector. |
| <code>w</code> | control variables. A matrix, a vector or a formula . |
| <code>data</code> | an optional data frame containing variables in the model. |
| <code>at</code> | value of <code>w</code> at which the estimated function is evaluated. The default is <code>at="mean"</code> , which corresponds to the mean of <code>w</code> . Other options are: <code>at="median"</code> for the median of <code>w</code> , <code>at="zero"</code> for a vector of zeros. <code>at</code> can also be a vector of the same length as the number of columns of <code>w</code> (if <code>w</code> is a matrix) or a data frame containing the same variables as specified in <code>w</code> (when <code>data</code> is specified). Note that when <code>at="mean"</code> or <code>at="median"</code> , all factor variables (if specified) are excluded from the evaluation (set as zero). |
| <code>family</code> | a description of the error distribution and link function to be used in the generalized linear model. (See family for details of family functions.) |
| <code>deriv</code> | derivative order of the regression function for estimation, testing and plotting. The default is <code>deriv=0</code> , which corresponds to the function itself. If <code>nolink=TRUE</code> , <code>deriv</code> cannot be greater than 1. |

| | |
|--------------|--|
| nolink | if true, the function within the inverse link function is reported instead of the conditional mean function for the outcome. |
| dots | a vector or a logical value. If $\text{dots}=\text{c}(p, s)$, a piecewise polynomial of degree p with s smoothness constraints is used for point estimation and plotting as "dots". The default is $\text{dots}=\text{c}(0, 0)$, which corresponds to piecewise constant (canonical binscatter). If $\text{dots}=\text{T}$, the default $\text{dots}=\text{c}(0, 0)$ is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If $\text{dots}=\text{F}$ is specified, the dots are not included in the plot. |
| dotsgrid | number of dots within each bin to be plotted. Given the choice, these dots are point estimates evaluated over an evenly-spaced grid within each bin. The default is <code>dotsgrid=0</code> , and only the point estimates at the mean of x within each bin are presented. |
| dotsgridmean | If true, the dots corresponding to the point estimates evaluated at the mean of x within each bin are presented. By default, they are presented, i.e., <code>dotsgridmean=T</code> . |
| line | a vector or a logical value. If $\text{line}=\text{c}(p, s)$, a piecewise polynomial of degree p with s smoothness constraints is used for plotting as a "line". If $\text{line}=\text{T}$ is specified, $\text{line}=\text{c}(0, 0)$ is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If $\text{line}=\text{F}$ or $\text{line}=\text{NULL}$ is specified, the line is not included in the plot. The default is $\text{line}=\text{NULL}$. |
| linegrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the $\text{line}=\text{c}(p, s)$ option. The default is <code>linegrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for fitting/plotting the line. |
| ci | a vector or a logical value. If $\text{ci}=\text{c}(p, s)$ a piecewise polynomial of degree p with s smoothness constraints is used for constructing confidence intervals. If $\text{ci}=\text{T}$ is specified, $\text{ci}=\text{c}(1, 1)$ is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If $\text{ci}=\text{F}$ or $\text{ci}=\text{NULL}$ is specified, the confidence intervals are not included in the plot. The default is $\text{ci}=\text{NULL}$. |
| cigrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the $\text{ci}=\text{c}(p, s)$ option. The default is <code>cigrid=1</code> , which corresponds to 1 evenly-spaced evaluation point within each bin for confidence interval construction. |
| cigridmean | If true, the confidence intervals corresponding to the point estimates evaluated at the mean of x within each bin are presented. The default is <code>cigridmean=T</code> . |
| cb | a vector or a logical value. If $\text{cb}=\text{c}(p, s)$, a the piecewise polynomial of degree p with s smoothness constraints is used for constructing the confidence band. If the option $\text{cb}=\text{T}$ is specified, $\text{cb}=\text{c}(1, 1)$ is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If $\text{cb}=\text{F}$ or $\text{cb}=\text{NULL}$ is specified, the confidence band is not included in the plot. The default is $\text{cb}=\text{NULL}$. |
| cbgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the $\text{cb}=\text{c}(p, s)$ option. The default is <code>cbgrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for confidence interval construction. |
| polyreg | degree of a global polynomial regression model for plotting. By default, this fit is not included in the plot unless explicitly specified. Recommended specification is <code>polyreg=3</code> , which adds a cubic (global) polynomial fit of the regression function of interest to the binned scatter plot. |

| | |
|--------------|--|
| polyreggrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>polyreg=p</code> option. The default is <code>polyreggrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for confidence interval construction. |
| polyregcgrid | number of evaluation points of an evenly-spaced grid within each bin used for constructing confidence intervals based on polynomial regression set by the <code>polyreg=p</code> option. The default is <code>polyregcgrid=0</code> , which corresponds to not plotting confidence intervals for the global polynomial regression approximation. |
| by | a vector containing the group indicator for subgroup analysis; both numeric and string variables are supported. When <code>by</code> is specified, <code>binsreg</code> implements estimation and inference for each subgroup separately, but produces a common binned scatter plot. By default, the binning structure is selected for each subgroup separately, but see the option <code>samebinsby</code> below for imposing a common binning structure across subgroups. |
| bycolors | an ordered list of colors for plotting each subgroup series defined by the option <code>by</code> . |
| bysymbols | an ordered list of symbols for plotting each subgroup series defined by the option <code>by</code> . |
| bylpatterns | an ordered list of line patterns for plotting each subgroup series defined by the option <code>by</code> . |
| legendTitle | String, title of legend. |
| legendoff | If true, no legend is added. |
| nbins | number of bins for partitioning/binning of x . If <code>nbins=T</code> or <code>nbins=NULL</code> (default) is specified, the number of bins is selected via the companion command binsregselect in a data-driven, optimal way whenever possible. If a vector with more than one number is specified, the number of bins is selected within this vector via the companion command binsregselect . |
| binspos | position of binning knots. The default is <code>binspos="qs"</code> , which corresponds to quantile-spaced binning (canonical <code>binscatter</code>). The other options are <code>"es"</code> for evenly-spaced binning, or a vector for manual specification of the positions of inner knots (which must be within the range of x). |
| binsmethod | method for data-driven selection of the number of bins. The default is <code>binsmethod="dpi"</code> , which corresponds to the IMSE-optimal direct plug-in rule. The other option is: <code>"rot"</code> for rule of thumb implementation. |
| nbinsrot | initial number of bins value used to construct the DPI number of bins selector. If not specified, the data-driven ROT selector is used instead. |
| pselect | vector of numbers within which the degree of polynomial p for point estimation is selected. Piecewise polynomials of the selected optimal degree p are used to construct dots or line if <code>dots=T</code> or <code>line=T</code> is specified, whereas piecewise polynomials of degree $p+1$ are used to construct confidence intervals or confidence band if <code>ci=T</code> or <code>cb=T</code> is specified. <i>Note:</i> To implement the degree or smoothness selection, in addition to <code>pselect</code> or <code>sselect</code> , <code>nbins=#</code> must be specified. |
| sselect | vector of numbers within which the number of smoothness constraints s for point estimation is selected. Piecewise polynomials with the selected optimal s smoothness constraints are used to construct dots or line if <code>dots=T</code> or <code>line=T</code> is specified, whereas piecewise polynomials with $s+1$ constraints are used to construct confidence intervals or confidence band if <code>ci=T</code> or <code>cb=T</code> is specified. If not specified, for each value p supplied in the option <code>pselect</code> , only the piecewise polynomial with the maximum smoothness is considered, i.e., $s=p$. |

| | |
|------------|---|
| samebinsby | if true, a common partitioning/binning structure across all subgroups specified by the option by is forced. The knots positions are selected according to the option binspos and using the full sample. If nbins is not specified, then the number of bins is selected via the companion command binsregselect and using the full sample. |
| randcut | upper bound on a uniformly distributed variable used to draw a subsample for bins/degree/smoothness selection. Observations for which <code>runif()<=#</code> are used. # must be between 0 and 1. By default, <code>max(5000, 0.01n)</code> observations are used if the samples size <code>n>5000</code> . |
| nsims | number of random draws for constructing confidence bands. The default is <code>nsims=500</code> , which corresponds to 500 draws from a standard Gaussian random vector of size $[(p+1)*J - (J-1)*s]$. Setting at least <code>nsims=2000</code> is recommended to obtain the final results. |
| simsgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the supremum operation needed to construct confidence bands. The default is <code>simsgrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for approximating the supremum operator. Setting at least <code>simsgrid=50</code> is recommended to obtain the final results. |
| simsseed | seed for simulation. |
| vce | Procedure to compute the variance-covariance matrix estimator. Options are <ul style="list-style-type: none"> • "const" homoskedastic variance estimator. • "HC0" heteroskedasticity-robust plug-in residuals variance estimator without weights. • "HC1" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc1</code> weights. Default. • "HC2" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc2</code> weights. • "HC3" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc3</code> weights. |
| cluster | cluster ID. Used for compute cluster-robust standard errors. |
| asyvar | if true, the standard error of the nonparametric component is computed and the uncertainty related to control variables is omitted. Default is <code>asyvar=FALSE</code> , that is, the uncertainty related to control variables is taken into account. |
| level | nominal confidence level for confidence interval and confidence band estimation. Default is <code>level=95</code> . |
| noplot | if true, no plot produced. |
| dfcheck | adjustments for minimum effective sample size checks, which take into account number of unique values of <code>x</code> (i.e., number of mass points), number of clusters, and degrees of freedom of the different stat models considered. The default is <code>dfcheck=c(20, 30)</code> . See Cattaneo, Crump, Farrell and Feng (2023c) for more details. |
| masspoints | how mass points in <code>x</code> are handled. Available options: <ul style="list-style-type: none"> • "on" all mass point and degrees of freedom checks are implemented. Default. • "noadjust" mass point checks and the corresponding effective sample size adjustments are omitted. • "nolocalcheck" within-bin mass point and degrees of freedom checks are omitted. |

| | |
|------------|---|
| | <ul style="list-style-type: none"> • "off" "noadjust" and "nolocalcheck" are set simultaneously. • "veryfew" forces the function to proceed as if x has only a few number of mass points (i.e., distinct values). In other words, forces the function to proceed as if the mass point and degrees of freedom checks were failed. |
| weights | an optional vector of weights to be used in the fitting process. Should be NULL or a numeric vector. For more details, see lm . |
| subset | optional rule specifying a subset of observations to be used. |
| plotxrange | a vector. <code>plotxrange=c(min,max)</code> specifies a range of the x-axis for binscatter plot. Observations outside the range are dropped in the plot. |
| plotyrange | a vector. <code>plotyrange=c(min,max)</code> specifies a range of the y-axis for binscatter plot. Observations outside the range are dropped in the plot. |
| ... | optional arguments used by glm . |

Value

| | |
|--------------|--|
| bins_plot | A ggplot object for binscatter plot. |
| data.plot | <p>A list containing data for plotting. Each item is a sublist of data frames for each group. Each sublist may contain the following data frames:</p> <ul style="list-style-type: none"> • <code>data.dots</code> Data for dots. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.line</code> Data for line. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.ci</code> Data for CI. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>ci.l</code> and <code>ci.r</code>, left and right boundaries of each confidence intervals. • <code>data.cb</code> Data for CB. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>cb.l</code> and <code>cb.r</code>, left and right boundaries of the confidence band. • <code>data.poly</code> Data for polynomial regression. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.polyci</code> Data for confidence intervals based on polynomial regression. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>polyci.l</code> and <code>polyci.r</code>, left and right boundaries of each confidence intervals. • <code>data.bin</code> Data for the binning structure. It contains: <code>bin.id</code>, ID for each bin; <code>left.endpoint</code> and <code>right.endpoint</code>, left and right endpoints of each bin. |
| imse.var.rot | Variance constant in IMSE, ROT selection. |
| imse.bsq.rot | Bias constant in IMSE, ROT selection. |
| imse.var.dpi | Variance constant in IMSE, DPI selection. |
| imse.bsq.dpi | Bias constant in IMSE, DPI selection. |
| cval.by | A vector of critical values for constructing confidence band for each group. |
| opt | A list containing options passed to the function, as well as <code>N.by</code> (total sample size for each group), <code>Ndist.by</code> (number of distinct values in x for each group), <code>Nclust.by</code> (number of clusters for each group), and <code>nbins.by</code> (number of bins |

for each group), and byvals (number of distinct values in by). The degree and smoothness of polynomials for dots, line, confidence intervals and confidence band for each group are saved in dots, line, ci, and cb.

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepk@gmail.com>.

References

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: [On Binscatter](#). Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: [Nonlinear Binscatter Methods](#). Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: [Binscatter Regressions](#). Working Paper.

See Also

[binsregselect](#), [binstest](#).

Examples

```
x <- runif(500); d <- 1*(runif(500)<=x)
## Binned scatterplot
binsglm(d, x, family=binomial())
```

binspwc

Data-Driven Pairwise Group Comparison using Binscatter Methods

Description

binspwc implements hypothesis testing procedures for pairwise group comparison of binscatter estimators, following the results in [Cattaneo, Crump, Farrell and Feng \(2023a\)](#) and [Cattaneo, Crump, Farrell and Feng \(2023b\)](#). If the binning scheme is not set by the user, the companion function [binsregselect](#) is used to implement binscatter in a data-driven way. Binned scatter plots based on different methods can be constructed using the companion functions [binsreg](#), [binsqreg](#) or [binsglm](#). Hypothesis testing for parametric functional forms of and shape restrictions on the regression function of interest can be conducted via the companion function [binstest](#).

Usage

```
binspwc(y, x, w = NULL, data = NULL, estmethod = "reg",
  family = gaussian(), quantile = NULL, deriv = 0, at = NULL,
  nolink = F, by = NULL, pwc = NULL, testtype = "two-sided",
  lp = Inf, bins = NULL, bynbins = NULL, binspos = "qs",
  pselect = NULL, sselect = NULL, binsmethod = "dpi", nbinsrot = NULL,
  samebinsby = FALSE, randcut = NULL, nsims = 500, simsgrid = 20,
```



```

simsseed = NULL, vce = NULL, cluster = NULL, asyvar = F,
dfcheck = c(20, 30), masspoints = "on", weights = NULL,
subset = NULL, numdist = NULL, numclust = NULL, estmethodopt = NULL,
...)

```

Arguments

| | |
|-----------|--|
| y | outcome variable. A vector. |
| x | independent variable of interest. A vector. |
| w | control variables. A matrix, a vector or a formula . |
| data | an optional data frame containing variables used in the model. |
| estmethod | estimation method. The default is <code>estmethod="reg"</code> for tests based on binscatter least squares regression. Other options are <code>"qreg"</code> for quantile regression and <code>"glm"</code> for generalized linear regression. If <code>estmethod="glm"</code> , the option family must be specified. |
| family | a description of the error distribution and link function to be used in the generalized linear model when <code>estmethod="glm"</code> . (See family for details of family functions.) |
| quantile | the quantile to be estimated. A number strictly between 0 and 1. |
| deriv | derivative order of the regression function for estimation, testing and plotting. The default is <code>deriv=0</code> , which corresponds to the function itself. |
| at | value of w at which the estimated function is evaluated. The default is <code>at="mean"</code> , which corresponds to the mean of w. Other options are: <code>at="median"</code> for the median of w, <code>at="zero"</code> for a vector of zeros. <code>at</code> can also be a vector of the same length as the number of columns of w (if w is a matrix) or a data frame containing the same variables as specified in w (when data is specified). Note that when <code>at="mean"</code> or <code>at="median"</code> , all factor variables (if specified) are excluded from the evaluation (set as zero). |
| nolink | if true, the function within the inverse link function is reported instead of the conditional mean function for the outcome. |
| by | a vector containing the group indicator for subgroup analysis; both numeric and string variables are supported. When <code>by</code> is specified, <code>binsreg</code> implements estimation and inference for each subgroup separately, but produces a common binned scatter plot. By default, the binning structure is selected for each subgroup separately, but see the option <code>samebinsby</code> below for imposing a common binning structure across subgroups. |
| pwc | a vector or a logical value. If <code>pwc=c(p,s)</code> , a piecewise polynomial of degree p with s smoothness constraints is used for testing the difference between groups. If <code>pwc=T</code> or <code>pwc=NULL</code> (default) is specified, <code>pwc=c(1,1)</code> is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). |
| testtype | type of pairwise comparison test. The default is <code>testtype="two-sided"</code> , which corresponds to a two-sided test of the form $H_0: \mu_1(x) = \mu_2(x)$. Other options are: <code>testtype="left"</code> for the one-sided test form $H_0: \mu_1(x) \leq \mu_2(x)$ and <code>testtype="right"</code> for the one-sided test of the form $H_0: \mu_1(x) \geq \mu_2(x)$. |
| lp | an Lp metric used for (two-sided) parametric model specification testing and/or shape restriction testing. The default is <code>lp=Inf</code> , which corresponds to the sup-norm of the t-statistic. Other options are <code>lp=q</code> for a positive integer q. |

| | |
|------------|---|
| bins | A vector. If $\text{bins}=\text{c}(p,s)$, it sets the piecewise polynomial of degree p with s smoothness constraints for data-driven (IMSE-optimal) selection of the partitioning/binning scheme. The default is $\text{bins}=\text{c}(0,0)$, which corresponds to the piecewise constant. |
| bynbins | a vector of the number of bins for partitioning/binning of x , which is applied to the binscatter estimation for each group. If a single number is specified, it is applied to the estimation for all groups. If $\text{bynbins}=\text{T}$ or $\text{bynbins}=\text{NULL}$ (default), the number of bins is selected via the companion function binsregselect in a data-driven way whenever possible. <i>Note:</i> If a vector with more than one number is supplied, it is understood as the number of bins applied to binscatter estimation for each subgroup rather than the range for selecting the number of bins. |
| binspos | position of binning knots. The default is $\text{binspos}=\text{"qs"}$, which corresponds to quantile-spaced binning (canonical binscatter). The other options are "es" for evenly-spaced binning, or a vector for manual specification of the positions of inner knots (which must be within the range of x). |
| pselect | vector of numbers within which the degree of polynomial p for point estimation is selected. If the selected optimal degree is p , then piecewise polynomials of degree $p+1$ are used to conduct pairwise group comparison. <i>Note:</i> To implement the degree or smoothness selection, in addition to pselect or sselect , $\text{bynbins}=\#$ must be specified. |
| sselect | vector of numbers within which the number of smoothness constraints s for point estimation is selected. If the selected optimal smoothness is s , then piecewise polynomials with $s+1$ smoothness constraints are used to conduct pairwise group comparison. If not specified, for each value p supplied in the option pselect , only the piecewise polynomial with the maximum smoothness is considered, i.e., $s=p$. |
| binsmethod | method for data-driven selection of the number of bins. The default is $\text{binsmethod}=\text{"dpi"}$, which corresponds to the IMSE-optimal direct plug-in rule. The other option is: "rot" for rule of thumb implementation. |
| nbinsrot | initial number of bins value used to construct the DPI number of bins selector. If not specified, the data-driven ROT selector is used instead. |
| samebinsby | if true, a common partitioning/binning structure across all subgroups specified by the option by is forced. The knots positions are selected according to the option binspos and using the full sample. If nbins is not specified, then the number of bins is selected via the companion command binsregselect and using the full sample. |
| randcut | upper bound on a uniformly distributed variable used to draw a subsample for bins/degree/smoothness selection. Observations for which $\text{runif}() \leq \#$ are used. $\#$ must be between 0 and 1. By default, $\max(5000, 0.01n)$ observations are used if the samples size $n > 5000$. |
| nsims | number of random draws for hypothesis testing. The default is $\text{nsims}=500$, which corresponds to 500 draws from a standard Gaussian random vector of size $[(p+1)*J - (J-1)*s]$. Setting at least $\text{nsims}=2000$ is recommended to obtain the final results. |
| simsgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the supremum (infimum or L_p metric) operation needed to construct hypothesis testing procedures. The default is $\text{simsgrid}=20$, which corresponds to 20 evenly-spaced evaluation points within each bin for approximating |

| | |
|---------------------------|---|
| | the supremum (infimum or Lp metric) operator. Setting at least <code>simsgrid=50</code> is recommended to obtain the final results. |
| <code>simsseed</code> | seed for simulation. |
| <code>vce</code> | procedure to compute the variance-covariance matrix estimator. For least squares regression and generalized linear regression, the allowed options are the same as that for <code>binsreg</code> or <code>binsqreg</code> . For quantile regression, the allowed options are the same as that for <code>binsqreg</code> . |
| <code>cluster</code> | cluster ID. Used for compute cluster-robust standard errors. |
| <code>asyvar</code> | if true, the standard error of the nonparametric component is computed and the uncertainty related to control variables is omitted. Default is <code>asyvar=FALSE</code> , that is, the uncertainty related to control variables is taken into account. |
| <code>dfcheck</code> | adjustments for minimum effective sample size checks, which take into account number of unique values of x (i.e., number of mass points), number of clusters, and degrees of freedom of the different stat models considered. The default is <code>dfcheck=c(20,30)</code> . See Cattaneo, Crump, Farrell and Feng (2023c) for more details. |
| <code>masspoints</code> | how mass points in x are handled. Available options: <ul style="list-style-type: none"> • "on" all mass point and degrees of freedom checks are implemented. Default. • "noadjust" mass point checks and the corresponding effective sample size adjustments are omitted. • "nolocalcheck" within-bin mass point and degrees of freedom checks are omitted. • "off" "noadjust" and "nolocalcheck" are set simultaneously. • "veryfew" forces the function to proceed as if x has only a few number of mass points (i.e., distinct values). In other words, forces the function to proceed as if the mass point and degrees of freedom checks were failed. |
| <code>weights</code> | an optional vector of weights to be used in the fitting process. Should be <code>NULL</code> or a numeric vector. For more details, see lm . |
| <code>subset</code> | optional rule specifying a subset of observations to be used. |
| <code>numdist</code> | number of distinct for selection. Used to speed up computation. |
| <code>numclust</code> | number of clusters for selection. Used to speed up computation. |
| <code>estmethodopt</code> | a list of optional arguments used by <code>rq</code> (for quantile regression) or <code>glm</code> (for fitting generalized linear models). |
| <code>...</code> | optional arguments to control bootstrapping if <code>estmethod="qreg"</code> and <code>vce="boot"</code> . See boot.rq . |

Value

| | |
|---------------------------|---|
| <code>stat</code> | A matrix. Each row corresponds to the comparison between two groups. The first column is the test statistic. The second and third columns give the corresponding group numbers. The null hypothesis is $\mu_i(x) \leq \mu_j(x)$, $\mu_i(x) = \mu_j(x)$, or $\mu_i(x) \geq \mu_j(x)$ for group i (given in the second column) and group j (given in the third column). The group number corresponds to the list of group names given by <code>opt\$byvals</code> . |
| <code>pval</code> | A vector of p-values for all pairwise group comparisons. |
| <code>imse.var.rot</code> | Variance constant in IMSE expansion, ROT selection. |

| | |
|----------------------------|--|
| <code>imse.bsqr.rot</code> | Bias constant in IMSE expansion, ROT selection. |
| <code>imse.var.dpi</code> | Variance constant in IMSE expansion, DPI selection. |
| <code>imse.bsqr.dpi</code> | Bias constant in IMSE expansion, DPI selection. |
| <code>opt</code> | A list containing options passed to the function, as well as <code>N.by</code> (total sample size for each group), <code>Ndist.by</code> (number of distinct values in <code>x</code> for each group), <code>Nclust.by</code> (number of clusters for each group), and <code>nbins.by</code> (number of bins for each group), and <code>byvals</code> (number of distinct values in <code>by</code>). |

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepk@gmail.com>.

References

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: [On Binscatter](#). Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: [Nonlinear Binscatter Methods](#). Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: [Binscatter Regressions](#). Working Paper.

See Also

[binsreg](#), [binsqreg](#), [binsglm](#), [binsregselect](#), [binstest](#).

Examples

```
x <- runif(500); y <- sin(x)+rnorm(500); t <- 1*(runif(500)>0.5)
## Binned scatterplot
binspwc(y,x, by=t)
```

binsqreg

Data-Driven Binscatter Quantile Regression with Robust Inference Procedures and Plots

Description

`binsqreg` implements `binscatter` quantile regression with robust inference procedures and plots, following the results in [Cattaneo, Crump, Farrell and Feng \(2023a\)](#) and [Cattaneo, Crump, Farrell and Feng \(2023b\)](#). `Binscatter` provides a flexible way to describe the quantile relationship between two variables, after possibly adjusting for other covariates, based on partitioning/binning of the independent variable of interest. The main purpose of this function is to generate binned scatter plots with curve estimation with robust pointwise confidence intervals and uniform confidence band. If the binning scheme is not set by the user, the companion function [binsregselect](#) is used to implement `binscatter` in a data-driven way. Hypothesis testing about the function of interest can be conducted via the companion function [binstest](#).

Usage

```
binsqreg(y, x, w = NULL, data = NULL, at = NULL, quantile = 0.5,
  deriv = 0, dots = NULL, dotsgrid = 0, dotsgridmean = T,
  line = NULL, linegrid = 20, ci = NULL, cigrid = 0, cigridmean = T,
  cb = NULL, cbgrid = 20, polyreg = NULL, polyreggrid = 20,
  polyregcigrid = 0, by = NULL, bycolors = NULL, bysymbols = NULL,
  bylpatterns = NULL, legendTitle = NULL, legendoff = F, nbins = NULL,
  binspos = "qs", binsmethod = "dpi", nbinsrot = NULL, pselect = NULL,
  sselect = NULL, samebinsby = F, randcut = NULL, nsims = 500,
  simsgrid = 20, simsseed = NULL, vce = "nid", cluster = NULL,
  asyvar = F, level = 95, noplot = F, dfcheck = c(20, 30),
  masspoints = "on", weights = NULL, subset = NULL, plotxrange = NULL,
  plotyrange = NULL, qregopt = NULL, ...)
```

Arguments

| | |
|--------------|--|
| y | outcome variable. A vector. |
| x | independent variable of interest. A vector. |
| w | control variables. A matrix, a vector or a formula . |
| data | an optional data frame containing variables in the model. |
| at | value of w at which the estimated function is evaluated. The default is at="mean", which corresponds to the mean of w. Other options are: at="median" for the median of w, at="zero" for a vector of zeros. at can also be a vector of the same length as the number of columns of w (if w is a matrix) or a data frame containing the same variables as specified in w (when data is specified). Note that when at="mean" or at="median", all factor variables (if specified) are excluded from the evaluation (set as zero). |
| quantile | the quantile to be estimated. A number strictly between 0 and 1. |
| deriv | derivative order of the regression function for estimation, testing and plotting. The default is deriv=0, which corresponds to the function itself. |
| dots | a vector or a logical value. If dots=c(p, s), a piecewise polynomial of degree p with s smoothness constraints is used for point estimation and plotting as "dots". The default is dots=c(0, 0), which corresponds to piecewise constant (canonical binscatter). If dots=T, the default dots=c(0, 0) is used unless the degree p or smoothness s selection is requested via the option pselect or sselect (see more details in the explanation of pselect and sselect). If dots=F is specified, the dots are not included in the plot. |
| dotsgrid | number of dots within each bin to be plotted. Given the choice, these dots are point estimates evaluated over an evenly-spaced grid within each bin. The default is dotsgrid=0, and only the point estimates at the mean of x within each bin are presented. |
| dotsgridmean | If true, the dots corresponding to the point estimates evaluated at the mean of x within each bin are presented. By default, they are presented, i.e., dotsgridmean=T. |
| line | a vector or a logical value. If line=c(p, s), a piecewise polynomial of degree p with s smoothness constraints is used for plotting as a "line". If line=T is specified, line=c(0, 0) is used unless the degree p or smoothness s selection is requested via the option pselect or sselect (see more details in the explanation of pselect and sselect). If line=F or line=NULL is specified, the line is not included in the plot. The default is line=NULL. |

| | |
|---------------|--|
| linegrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>line=c(p,s)</code> option. The default is <code>linegrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for fitting/plotting the line. |
| ci | a vector or a logical value. If <code>ci=c(p,s)</code> a piecewise polynomial of degree <code>p</code> with <code>s</code> smoothness constraints is used for constructing confidence intervals. If <code>ci=T</code> is specified, <code>ci=c(1,1)</code> is used unless the degree <code>p</code> or smoothness <code>s</code> selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If <code>ci=F</code> or <code>ci=NULL</code> is specified, the confidence intervals are not included in the plot. The default is <code>ci=NULL</code> . |
| cigrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>ci=c(p,s)</code> option. The default is <code>cigrid=1</code> , which corresponds to 1 evenly-spaced evaluation point within each bin for confidence interval construction. |
| cigridmean | If true, the confidence intervals corresponding to the point estimates evaluated at the mean of <code>x</code> within each bin are presented. The default is <code>cigridmean=T</code> . |
| cb | a vector or a logical value. If <code>cb=c(p,s)</code> , a the piecewise polynomial of degree <code>p</code> with <code>s</code> smoothness constraints is used for constructing the confidence band. If the option <code>cb=T</code> is specified, <code>cb=c(1,1)</code> is used unless the degree <code>p</code> or smoothness <code>s</code> selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If <code>cb=F</code> or <code>cb=NULL</code> is specified, the confidence band is not included in the plot. The default is <code>cb=NULL</code> . |
| cbgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>cb=c(p,s)</code> option. The default is <code>cbgrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for confidence interval construction. |
| polyreg | degree of a global polynomial regression model for plotting. By default, this fit is not included in the plot unless explicitly specified. Recommended specification is <code>polyreg=3</code> , which adds a cubic (global) polynomial fit of the regression function of interest to the binned scatter plot. |
| polyreggrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>polyreg=p</code> option. The default is <code>polyreggrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for confidence interval construction. |
| polyregcigrid | number of evaluation points of an evenly-spaced grid within each bin used for constructing confidence intervals based on polynomial regression set by the <code>polyreg=p</code> option. The default is <code>polyregcigrid=0</code> , which corresponds to not plotting confidence intervals for the global polynomial regression approximation. |
| by | a vector containing the group indicator for subgroup analysis; both numeric and string variables are supported. When <code>by</code> is specified, <code>binsreg</code> implements estimation and inference for each subgroup separately, but produces a common binned scatter plot. By default, the binning structure is selected for each subgroup separately, but see the option <code>samebinsby</code> below for imposing a common binning structure across subgroups. |
| bycolors | an ordered list of colors for plotting each subgroup series defined by the option <code>by</code> . |
| bysymbols | an ordered list of symbols for plotting each subgroup series defined by the option <code>by</code> . |

| | |
|-------------|--|
| bylpatterns | an ordered list of line patterns for plotting each subgroup series defined by the option by. |
| legendTitle | String, title of legend. |
| legendoff | If true, no legend is added. |
| nbins | number of bins for partitioning/binning of x . If nbins=T or nbins=NULL (default) is specified, the number of bins is selected via the companion command binsregselect in a data-driven, optimal way whenever possible. If a vector with more than one number is specified, the number of bins is selected within this vector via the companion command binsregselect . |
| binspos | position of binning knots. The default is binspos="qs", which corresponds to quantile-spaced binning (canonical binscatter). The other options are "es" for evenly-spaced binning, or a vector for manual specification of the positions of inner knots (which must be within the range of x). |
| binsmethod | method for data-driven selection of the number of bins. The default is binsmethod="dpi", which corresponds to the IMSE-optimal direct plug-in rule. The other option is: "rot" for rule of thumb implementation. |
| nbinsrot | initial number of bins value used to construct the DPI number of bins selector. If not specified, the data-driven ROT selector is used instead. |
| pselect | vector of numbers within which the degree of polynomial p for point estimation is selected. Piecewise polynomials of the selected optimal degree p are used to construct dots or line if dots=T or line=T is specified, whereas piecewise polynomials of degree $p+1$ are used to construct confidence intervals or confidence band if ci=T or cb=T is specified. <i>Note:</i> To implement the degree or smoothness selection, in addition to pselect or sselect, nbins=# must be specified. |
| sselect | vector of numbers within which the number of smoothness constraints s for point estimation is selected. Piecewise polynomials with the selected optimal s smoothness constraints are used to construct dots or line if dots=T or line=T is specified, whereas piecewise polynomials with $s+1$ constraints are used to construct confidence intervals or confidence band if ci=T or cb=T is specified. If not specified, for each value p supplied in the option pselect, only the piecewise polynomial with the maximum smoothness is considered, i.e., $s=p$. |
| samebinsby | if true, a common partitioning/binning structure across all subgroups specified by the option by is forced. The knots positions are selected according to the option binspos and using the full sample. If nbins is not specified, then the number of bins is selected via the companion command binsregselect and using the full sample. |
| randcut | upper bound on a uniformly distributed variable used to draw a subsample for bins/degree/smoothness selection. Observations for which <code>runif()<=#</code> are used. # must be between 0 and 1. By default, <code>max(5000, 0.01n)</code> observations are used if the samples size $n > 5000$. |
| nsims | number of random draws for constructing confidence bands. The default is nsims=500, which corresponds to 500 draws from a standard Gaussian random vector of size $[(p+1)*J - (J-1)*s]$. Setting at least nsims=2000 is recommended to obtain the final results. |
| simsgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the supremum operation needed to construct confidence bands. The default is simsgrid=20, which corresponds to 20 evenly-spaced evaluation points within each bin for approximating the supremum operator. Setting at least simsgrid=50 is recommended to obtain the final results. |

| | |
|------------|--|
| simsseed | seed for simulation. |
| vce | <p>Procedure to compute the variance-covariance matrix estimator (see summary.rq for more details). Options are</p> <ul style="list-style-type: none"> • "iid" which presumes that the errors are iid and computes an estimate of the asymptotic covariance matrix as in KB(1978). • "nid" which presumes local (in quantile) linearity of the the conditional quantile functions and computes a Huber sandwich estimate using a local estimate of the sparsity. • "ker" which uses a kernel estimate of the sandwich as proposed by Powell (1991). • "boot" which implements one of several possible bootstrapping alternatives for estimating standard errors including a variate of the wild bootstrap for clustered response. See boot.rq for further details. |
| cluster | cluster ID. Used for compute cluster-robust standard errors. |
| asyvar | if true, the standard error of the nonparametric component is computed and the uncertainty related to control variables is omitted. Default is asyvar=FALSE, that is, the uncertainty related to control variables is taken into account. |
| level | nominal confidence level for confidence interval and confidence band estimation. Default is level=95. |
| noplot | if true, no plot produced. |
| dfcheck | adjustments for minimum effective sample size checks, which take into account number of unique values of x (i.e., number of mass points), number of clusters, and degrees of freedom of the different statistical models considered. The default is dfcheck=c(20,30). See Cattaneo, Crump, Farrell and Feng (2023c) for more details. |
| masspoints | <p>how mass points in x are handled. Available options:</p> <ul style="list-style-type: none"> • "on" all mass point and degrees of freedom checks are implemented. Default. • "noadjust" mass point checks and the corresponding effective sample size adjustments are omitted. • "nolocalcheck" within-bin mass point and degrees of freedom checks are omitted. • "off" "noadjust" and "nolocalcheck" are set simultaneously. • "veryfew" forces the function to proceed as if x has only a few number of mass points (i.e., distinct values). In other words, forces the function to proceed as if the mass point and degrees of freedom checks were failed. |
| weights | an optional vector of weights to be used in the fitting process. Should be NULL or a numeric vector. For more details, see lm . |
| subset | optional rule specifying a subset of observations to be used. |
| plotxrange | a vector. plotxrange=c(min,max) specifies a range of the x-axis for plotting. Observations outside the range are dropped in the plot. |
| plotyrange | a vector. plotyrange=c(min,max) specifies a range of the y-axis for plotting. Observations outside the range are dropped in the plot. |
| qregopt | a list of optional arguments used by rq . |
| ... | optional arguments to control bootstrapping. See boot.rq . |

Value

| | |
|---------------------------|--|
| <code>bins_plot</code> | A ggplot object for binscatter plot. |
| <code>data.plot</code> | <p>A list containing data for plotting. Each item is a sublist of data frames for each group. Each sublist may contain the following data frames:</p> <ul style="list-style-type: none"> • <code>data.dots</code> Data for dots. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.line</code> Data for line. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.ci</code> Data for CI. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>ci.l</code> and <code>ci.r</code>, left and right boundaries of each confidence intervals. • <code>data.cb</code> Data for CB. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>cb.l</code> and <code>cb.r</code>, left and right boundaries of the confidence band. • <code>data.poly</code> Data for polynomial regression. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.polyci</code> Data for confidence intervals based on polynomial regression. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>polyci.l</code> and <code>polyci.r</code>, left and right boundaries of each confidence intervals. • <code>data.bin</code> Data for the binning structure. It contains: <code>bin.id</code>, ID for each bin; <code>left.endpoint</code> and <code>right.endpoint</code>, left and right endpoints of each bin. |
| <code>imse.var.rot</code> | Variance constant in IMSE, ROT selection. |
| <code>imse.bsq.rot</code> | Bias constant in IMSE, ROT selection. |
| <code>imse.var.dpi</code> | Variance constant in IMSE, DPI selection. |
| <code>imse.bsq.dpi</code> | Bias constant in IMSE, DPI selection. |
| <code>cval.by</code> | A vector of critical values for constructing confidence band for each group. |
| <code>opt</code> | A list containing options passed to the function, as well as <code>N.by</code> (total sample size for each group), <code>Ndist.by</code> (number of distinct values in <code>x</code> for each group), <code>Nclust.by</code> (number of clusters for each group), and <code>nbins.by</code> (number of bins for each group), and <code>byvals</code> (number of distinct values in <code>by</code>). The degree and smoothness of polynomials for dots, line, confidence intervals and confidence band for each group are saved in <code>dots</code> , <code>line</code> , <code>ci</code> , and <code>cb</code> . |

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepku@gmail.com>.

References

- Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: [On Binscatter](#). Working Paper.
- Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: [Nonlinear Binscatter Methods](#). Working Paper.
- Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: [Binscatter Regressions](#). Working Paper.

See Also

[binsregselect](#), [binstest](#).

Examples

```
x <- runif(500); y <- sin(x)+rnorm(500)
## Binned scatterplot
binsqreg(y,x)
```

binsreg

Data-Driven Binscatter Least Squares Regression with Robust Inference Procedures and Plots

Description

binsreg implements binscatter least squares regression with robust inference procedures and plots, following the results in [Cattaneo, Crump, Farrell and Feng \(2023a\)](#) and [Cattaneo, Crump, Farrell and Feng \(2023b\)](#). Binscatter provides a flexible way to describe the mean relationship between two variables, after possibly adjusting for other covariates, based on partitioning/binning of the independent variable of interest. The main purpose of this function is to generate binned scatter plots with curve estimation with robust pointwise confidence intervals and uniform confidence band. If the binning scheme is not set by the user, the companion function [binsregselect](#) is used to implement binscatter in a data-driven (optimal) way. Hypothesis testing about the regression function can be conducted via the companion function [binstest](#).

Usage

```
binsreg(y, x, w = NULL, data = NULL, at = NULL, deriv = 0,
  dots = NULL, dotsgrid = 0, dotsgridmean = T, line = NULL,
  linegrid = 20, ci = NULL, cigrid = 0, cigridmean = T, cb = NULL,
  cbgrid = 20, polyreg = NULL, polyreggrid = 20, polyregcigrid = 0,
  by = NULL, bycolors = NULL, bysymbols = NULL, bylpatterns = NULL,
  legendTitle = NULL, legendoff = F, nbins = NULL, binspos = "qs",
  binsmethod = "dpi", nbinsrot = NULL, pselect = NULL, sselect = NULL,
  samebinsby = F, randcut = NULL, nsims = 500, simsgrid = 20,
  simsseed = NULL, vce = "HC1", cluster = NULL, asyvar = F,
  level = 95, noplot = F, dfcheck = c(20, 30), masspoints = "on",
  weights = NULL, subset = NULL, plotxrange = NULL, plotyrange = NULL)
```

Arguments

| | |
|---------------------------|---|
| <code>y</code> | outcome variable. A vector. |
| <code>x</code> | independent variable of interest. A vector. |
| <code>w</code> | control variables. A matrix, a vector or a formula . |
| <code>data</code> | an optional data frame containing variables used in the model. |
| <code>at</code> | value of <code>w</code> at which the estimated function is evaluated. The default is <code>at="mean"</code> , which corresponds to the mean of <code>w</code> . Other options are: <code>at="median"</code> for the median of <code>w</code> , <code>at="zero"</code> for a vector of zeros. <code>at</code> can also be a vector of the same length as the number of columns of <code>w</code> (if <code>w</code> is a matrix) or a data frame containing the same variables as specified in <code>w</code> (when <code>data</code> is specified). Note that when <code>at="mean"</code> or <code>at="median"</code> , all factor variables (if specified) are excluded from the evaluation (set as zero). |
| <code>deriv</code> | derivative order of the regression function for estimation, testing and plotting. The default is <code>deriv=0</code> , which corresponds to the function itself. |
| <code>dots</code> | a vector or a logical value. If <code>dots=c(p, s)</code> , a piecewise polynomial of degree <code>p</code> with <code>s</code> smoothness constraints is used for point estimation and plotting as "dots". The default is <code>dots=c(0, 0)</code> , which corresponds to piecewise constant (canonical binscatter). If <code>dots=T</code> , the default <code>dots=c(0, 0)</code> is used unless the degree <code>p</code> or smoothness <code>s</code> selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If <code>dots=F</code> is specified, the dots are not included in the plot. |
| <code>dotsgrid</code> | number of dots within each bin to be plotted. Given the choice, these dots are point estimates evaluated over an evenly-spaced grid within each bin. The default is <code>dotsgrid=0</code> , and only the point estimates at the mean of <code>x</code> within each bin are presented. |
| <code>dotsgridmean</code> | If true, the dots corresponding to the point estimates evaluated at the mean of <code>x</code> within each bin are presented. By default, they are presented, i.e., <code>dotsgridmean=T</code> . |
| <code>line</code> | a vector or a logical value. If <code>line=c(p, s)</code> , a piecewise polynomial of degree <code>p</code> with <code>s</code> smoothness constraints is used for plotting as a "line". If <code>line=T</code> is specified, <code>line=c(0, 0)</code> is used unless the degree <code>p</code> or smoothness <code>s</code> selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If <code>line=F</code> or <code>line=NULL</code> is specified, the line is not included in the plot. The default is <code>line=NULL</code> . |
| <code>linegrid</code> | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>line=c(p, s)</code> option. The default is <code>linegrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for fitting/plotting the line. |
| <code>ci</code> | a vector or a logical value. If <code>ci=c(p, s)</code> a piecewise polynomial of degree <code>p</code> with <code>s</code> smoothness constraints is used for constructing confidence intervals. If <code>ci=T</code> is specified, <code>ci=c(1, 1)</code> is used unless the degree <code>p</code> or smoothness <code>s</code> selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If <code>ci=F</code> or <code>ci=NULL</code> is specified, the confidence intervals are not included in the plot. The default is <code>ci=NULL</code> . |
| <code>cigrid</code> | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the <code>ci=c(p, s)</code> option. The default is <code>cigrid=1</code> , which corresponds to 1 evenly-spaced evaluation point within each bin for confidence interval construction. |
| <code>cigridmean</code> | If true, the confidence intervals corresponding to the point estimates evaluated at the mean of <code>x</code> within each bin are presented. The default is <code>cigridmean=T</code> . |

| | |
|--------------|--|
| cb | a vector or a logical value. If $cb=c(p, s)$, a the piecewise polynomial of degree p with s smoothness constraints is used for constructing the confidence band. If the option $cb=T$ is specified, $cb=c(1, 1)$ is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). If $cb=F$ or $cb=NULL$ is specified, the confidence band is not included in the plot. The default is $cb=NULL$. |
| cbgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the $cb=c(p, s)$ option. The default is $cbgrid=20$, which corresponds to 20 evenly-spaced evaluation points within each bin for confidence interval construction. |
| polyreg | degree of a global polynomial regression model for plotting. By default, this fit is not included in the plot unless explicitly specified. Recommended specification is $polyreg=3$, which adds a cubic (global) polynomial fit of the regression function of interest to the binned scatter plot. |
| polyreggrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the point estimate set by the $polyreg=p$ option. The default is $polyreggrid=20$, which corresponds to 20 evenly-spaced evaluation points within each bin for confidence interval construction. |
| polyregcgrid | number of evaluation points of an evenly-spaced grid within each bin used for constructing confidence intervals based on polynomial regression set by the $polyreg=p$ option. The default is $polyregcgrid=0$, which corresponds to not plotting confidence intervals for the global polynomial regression approximation. |
| by | a vector containing the group indicator for subgroup analysis; both numeric and string variables are supported. When <code>by</code> is specified, <code>binsreg</code> implements estimation and inference for each subgroup separately, but produces a common binned scatter plot. By default, the binning structure is selected for each subgroup separately, but see the option <code>samebinsby</code> below for imposing a common binning structure across subgroups. |
| bycolors | an ordered list of colors for plotting each subgroup series defined by the option <code>by</code> . |
| bysymbols | an ordered list of symbols for plotting each subgroup series defined by the option <code>by</code> . |
| bylpatterns | an ordered list of line patterns for plotting each subgroup series defined by the option <code>by</code> . |
| legendTitle | String, title of legend. |
| legendoff | If true, no legend is added. |
| nbins | number of bins for partitioning/binning of x . If $nbins=T$ or $nbins=NULL$ (default) is specified, the number of bins is selected via the companion command <code>binsregselect</code> in a data-driven, optimal way whenever possible. If a vector with more than one number is specified, the number of bins is selected within this vector via the companion command <code>binsregselect</code> . |
| binspos | position of binning knots. The default is $binspos="qs"$, which corresponds to quantile-spaced binning (canonical <code>binscatter</code>). The other options are <code>"es"</code> for evenly-spaced binning, or a vector for manual specification of the positions of inner knots (which must be within the range of x). |
| binsmethod | method for data-driven selection of the number of bins. The default is $binsmethod="dpi"$, which corresponds to the IMSE-optimal direct plug-in rule. The other option is: <code>"rot"</code> for rule of thumb implementation. |

| | |
|------------|--|
| nbinsrot | initial number of bins value used to construct the DPI number of bins selector. If not specified, the data-driven ROT selector is used instead. |
| pselect | vector of numbers within which the degree of polynomial p for point estimation is selected. Piecewise polynomials of the selected optimal degree p are used to construct dots or line if <code>dots=T</code> or <code>line=T</code> is specified, whereas piecewise polynomials of degree $p+1$ are used to construct confidence intervals or confidence band if <code>ci=T</code> or <code>cb=T</code> is specified. <i>Note:</i> To implement the degree or smoothness selection, in addition to <code>pselect</code> or <code>sselect</code> , <code>nbins=#</code> must be specified. |
| sselect | vector of numbers within which the number of smoothness constraints s for point estimation is selected. Piecewise polynomials with the selected optimal s smoothness constraints are used to construct dots or line if <code>dots=T</code> or <code>line=T</code> is specified, whereas piecewise polynomials with $s+1$ constraints are used to construct confidence intervals or confidence band if <code>ci=T</code> or <code>cb=T</code> is specified. If not specified, for each value p supplied in the option <code>pselect</code> , only the piecewise polynomial with the maximum smoothness is considered, i.e., $s=p$. |
| samebinsby | if true, a common partitioning/binning structure across all subgroups specified by the option <code>by</code> is forced. The knots positions are selected according to the option <code>binspos</code> and using the full sample. If <code>nbins</code> is not specified, then the number of bins is selected via the companion command <code>binsregselect</code> and using the full sample. |
| randcut | upper bound on a uniformly distributed variable used to draw a subsample for bins/degree/smoothness selection. Observations for which <code>runif()<=#</code> are used. <code>#</code> must be between 0 and 1. By default, <code>max(5000, 0.01n)</code> observations are used if the samples size $n > 5000$. |
| nsims | number of random draws for constructing confidence bands. The default is <code>nsims=500</code> , which corresponds to 500 draws from a standard Gaussian random vector of size $[(p+1)*J - (J-1)*s]$. Setting at least <code>nsims=2000</code> is recommended to obtain the final results. |
| simsgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the supremum operation needed to construct confidence bands. The default is <code>simsgrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for approximating the supremum operator. Setting at least <code>simsgrid=50</code> is recommended to obtain the final results. |
| simsseed | seed for simulation. |
| vce | Procedure to compute the variance-covariance matrix estimator. Options are <ul style="list-style-type: none"> • "const" homoskedastic variance estimator. • "HC0" heteroskedasticity-robust plug-in residuals variance estimator without weights. • "HC1" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc1</code> weights. Default. • "HC2" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc2</code> weights. • "HC3" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc3</code> weights. |
| cluster | cluster ID. Used for compute cluster-robust standard errors. |
| asyvar | If true, the standard error of the nonparametric component is computed and the uncertainty related to control variables is omitted. Default is <code>asyvar=FALSE</code> , that is, the uncertainty related to control variables is taken into account. |

| | |
|------------|---|
| level | nominal confidence level for confidence interval and confidence band estimation. Default is level=95. |
| noplot | if true, no plot produced. |
| dfcheck | adjustments for minimum effective sample size checks, which take into account number of unique values of x (i.e., number of mass points), number of clusters, and degrees of freedom of the different statistical models considered. The default is <code>dfcheck=c(20,30)</code> . See Cattaneo, Crump, Farrell and Feng (2023c) for more details. |
| masspoints | how mass points in x are handled. Available options: <ul style="list-style-type: none"> • "on" all mass point and degrees of freedom checks are implemented. Default. • "noadjust" mass point checks and the corresponding effective sample size adjustments are omitted. • "nolocalcheck" within-bin mass point and degrees of freedom checks are omitted. • "off" "noadjust" and "nolocalcheck" are set simultaneously. • "veryfew" forces the function to proceed as if x has only a few number of mass points (i.e., distinct values). In other words, forces the function to proceed as if the mass point and degrees of freedom checks were failed. |
| weights | an optional vector of weights to be used in the fitting process. Should be NULL or a numeric vector. For more details, see lm . |
| subset | Optional rule specifying a subset of observations to be used. |
| plotxrange | a vector. <code>plotxrange=c(min,max)</code> specifies a range of the x-axis for plotting. Observations outside the range are dropped in the plot. |
| plotyrange | a vector. <code>plotyrange=c(min,max)</code> specifies a range of the y-axis for plotting. Observations outside the range are dropped in the plot. |

Value

| | |
|-----------|--|
| bins_plot | A ggplot object for binscatter plot. |
| data.plot | A list containing data for plotting. Each item is a sublist of data frames for each group. Each sublist may contain the following data frames: <ul style="list-style-type: none"> • <code>data.dots</code> Data for dots. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.line</code> Data for line. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. • <code>data.ci</code> Data for CI. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>ci.l</code> and <code>ci.r</code>, left and right boundaries of each confidence intervals. • <code>data.cb</code> Data for CB. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>cb.l</code> and <code>cb.r</code>, left and right boundaries of the confidence band. • <code>data.poly</code> Data for polynomial regression. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; and <code>fit</code>, fitted values. |

| | |
|----------------------------|---|
| | <ul style="list-style-type: none"> • <code>data.polyci</code> Data for confidence intervals based on polynomial regression. It contains: <code>x</code>, evaluation points; <code>bin</code>, the indicator of bins; <code>isknot</code>, indicator of inner knots; <code>mid</code>, midpoint of each bin; <code>polyci.l</code> and <code>polyci.r</code>, left and right boundaries of each confidence intervals. • <code>data.bin</code> Data for the binning structure. It contains: <code>bin.id</code>, ID for each bin; <code>left.endpoint</code> and <code>right.endpoint</code>, left and right endpoints of each bin. |
| <code>imse.var.rot</code> | Variance constant in IMSE, ROT selection. |
| <code>imse.bsqr.rot</code> | Bias constant in IMSE, ROT selection. |
| <code>imse.var.dpi</code> | Variance constant in IMSE, DPI selection. |
| <code>imse.bsqr.dpi</code> | Bias constant in IMSE, DPI selection. |
| <code>cval.by</code> | A vector of critical values for constructing confidence band for each group. |
| <code>opt</code> | A list containing options passed to the function, as well as <code>N.by</code> (total sample size for each group), <code>Ndist.by</code> (number of distinct values in <code>x</code> for each group), <code>Nclust.by</code> (number of clusters for each group), and <code>nbins.by</code> (number of bins for each group), and <code>byvals</code> (number of distinct values in <code>by</code>). The degree and smoothness of polynomials for dots, line, confidence intervals and confidence band for each group are saved in <code>dots</code> , <code>line</code> , <code>ci</code> , and <code>cb</code> . |

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepk@gmail.com>.

References

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: [On Binscatter](#). Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: [Nonlinear Binscatter Methods](#). Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: [Binscatter Regressions](#). Working Paper.

See Also

[binsregselect](#), [binstest](#).

Examples

```
x <- runif(500); y <- sin(x)+rnorm(500)
## Binned scatterplot
binsreg(y,x)
```

| | |
|---------------|---|
| binsregselect | <i>Data-Driven IMSE-Optimal Partitioning/Binning Selection for Binscatter</i> |
|---------------|---|

Description

binsregselect implements data-driven procedures for selecting the number of bins for binscatter estimation. The selected number is optimal in minimizing integrated mean squared error (IMSE).

Usage

```
binsregselect(y, x, w = NULL, data = NULL, deriv = 0, bins = NULL,
  pselect = NULL, sselect = NULL, binspos = "qs", nbins = NULL,
  binsmethod = "dpi", nbinsrot = NULL, simsgrid = 20, savegrid = F,
  vce = "HC1", useeffn = NULL, randcut = NULL, cluster = NULL,
  dfcheck = c(20, 30), masspoints = "on", weights = NULL,
  subset = NULL, norotnorm = F, numdist = NULL, numclust = NULL)
```

Arguments

| | |
|------------|---|
| y | outcome variable. A vector. |
| x | independent variable of interest. A vector. |
| w | control variables. A matrix, a vector or a formula . |
| data | an optional data frame containing variables used in the model. |
| deriv | derivative order of the regression function for estimation, testing and plotting. The default is deriv=0, which corresponds to the function itself. |
| bins | a vector. bins=c(p,s) set a piecewise polynomial of degree p with s smoothness constraints for data-driven (IMSE-optimal) selection of the partitioning/binning scheme. By default, the function sets bins=c(0,0), which corresponds to piecewise constant (canonical binscatter). |
| pselect | vector of numbers within which the degree of polynomial p for point estimation is selected. <i>Note:</i> To implement the degree or smoothness selection, in addition to pselect or sselect, nbins=# must be specified. |
| sselect | vector of numbers within which the number of smoothness constraints s for point estimation is selected. If not specified, for each value p supplied in the option pselect, only the piecewise polynomial with the maximum smoothness is considered, i.e., s=p. |
| binspos | position of binning knots. The default is binspos="qs", which corresponds to quantile-spaced binning (canonical binscatter). The other option is binspos="es" for evenly-spaced binning. |
| nbins | number of bins for degree/smoothness selection. If nbins=T or nbins=NULL (default) is specified, the function selects the number of bins instead, given the specified degree and smoothness. If a vector with more than one number is specified, the command selects the number of bins within this vector. |
| binsmethod | method for data-driven selection of the number of bins. The default is binsmethod="dpi", which corresponds to the IMSE-optimal direct plug-in rule. The other option is: "rot" for rule of thumb implementation. |

| | |
|------------|--|
| nbinsrot | initial number of bins value used to construct the DPI number of bins selector. If not specified, the data-driven ROT selector is used instead. |
| simsgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the supremum (infimum or Lp metric) operation needed to construct confidence bands and hypothesis testing procedures. The default is <code>simsgrid=20</code> , which corresponds to 20 evenly-spaced evaluation points within each bin for approximating the supremum (infimum or Lp metric) operator. |
| savegrid | if true, a data frame produced containing grid. |
| vce | procedure to compute the variance-covariance matrix estimator. Options are <ul style="list-style-type: none"> • "const" homoskedastic variance estimator. • "HC0" heteroskedasticity-robust plug-in residuals variance estimator without weights. • "HC1" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc1</code> weights. Default. • "HC2" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc2</code> weights. • "HC3" heteroskedasticity-robust plug-in residuals variance estimator with <code>hc3</code> weights. |
| useeffn | effective sample size to be used when computing the (IMSE-optimal) number of bins. This option is useful for extrapolating the optimal number of bins to larger (or smaller) datasets than the one used to compute it. |
| randcut | upper bound on a uniformly distributed variable used to draw a subsample for bins/degree/smoothness selection. Observations for which <code>runif()<=#</code> are used. <code>#</code> must be between 0 and 1. |
| cluster | cluster ID. Used for compute cluster-robust standard errors. |
| dfcheck | adjustments for minimum effective sample size checks, which take into account number of unique values of <code>x</code> (i.e., number of mass points), number of clusters, and degrees of freedom of the different statistical models considered. The default is <code>dfcheck=c(20,30)</code> . See Cattaneo, Crump, Farrell and Feng (2023c) for more details. |
| masspoints | how mass points in <code>x</code> are handled. Available options: <ul style="list-style-type: none"> • "on" all mass point and degrees of freedom checks are implemented. Default. • "noadjust" mass point checks and the corresponding effective sample size adjustments are omitted. • "nolocalcheck" within-bin mass point and degrees of freedom checks are omitted. • "off" "noadjust" and "nolocalcheck" are set simultaneously. • "veryfew" forces the function to proceed as if <code>x</code> has only a few number of mass points (i.e., distinct values). In other words, forces the function to proceed as if the mass point and degrees of freedom checks were failed. |
| weights | an optional vector of weights to be used in the fitting process. Should be <code>NULL</code> or a numeric vector. For more details, see 1m . |
| subset | optional rule specifying a subset of observations to be used. |
| norotnorm | if true, a uniform density rather than normal density used for ROT selection. |
| numdist | number of distinct for selection. Used to speed up computation. |
| numclust | number of clusters for selection. Used to speed up computation. |

Value

| | |
|-----------------------------|---|
| <code>nbinsrot.poly</code> | ROT number of bins, unregularized. |
| <code>nbinsrot.regul</code> | ROT number of bins, regularized. |
| <code>nbinsrot.uknot</code> | ROT number of bins, unique knots. |
| <code>nbinsdpi</code> | DPI number of bins. |
| <code>nbinsdpi.uknot</code> | DPI number of bins, unique knots. |
| <code>prot.poly</code> | ROT degree of polynomials, unregularized. |
| <code>prot.regul</code> | ROT degree of polynomials, regularized. |
| <code>prot.uknot</code> | ROT degree of polynomials, unique knots. |
| <code>pdpi</code> | DPI degree of polynomials. |
| <code>pdpi.uknot</code> | DPI degree of polynomials, unique knots. |
| <code>srot.poly</code> | ROT number of smoothness constraints, unregularized. |
| <code>srot.regul</code> | ROT number of smoothness constraints, regularized. |
| <code>srot.uknot</code> | ROT number of smoothness constraints, unique knots. |
| <code>sdpi</code> | DPI number of smoothness constraints. |
| <code>sdpi.uknot</code> | DPI number of smoothness constraints, unique knots. |
| <code>imse.var.rot</code> | Variance constant in IMSE expansion, ROT selection. |
| <code>imse.bsqr.rot</code> | Bias constant in IMSE expansion, ROT selection. |
| <code>imse.var.dpi</code> | Variance constant in IMSE expansion, DPI selection. |
| <code>imse.bsqr.dpi</code> | Bias constant in IMSE expansion, DPI selection. |
| <code>int.result</code> | Intermediate results, including a matrix of degree and smoothness (<code>deg_mat</code>), the selected numbers of bins (<code>vec.nbinsrot.poly</code> , <code>vec.nbinsrot.regul</code> , <code>vec.nbinsrot.uknot</code> , <code>vec.nbinsdpi</code> , <code>vec.nbinsdpi.uknot</code>), and the bias and variance constants in IMSE (<code>vec.imse.b.rot</code> , <code>vec.imse.v.rot</code> , <code>vec.imse.b.dpi</code> , <code>vec.imse.v.dpi</code>) under each rule (ROT or DPI), corresponding to each pair of degree and smoothness (each row in <code>deg_mat</code>). |
| <code>opt</code> | A list containing options passed to the function, as well as total sample size <code>n</code> , number of distinct values <code>Ndist</code> in <code>x</code> , and number of clusters <code>Nclust</code> . |
| <code>data.grid</code> | A data frame containing grid. |

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepku@gmail.com>.

References

- Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: **On Binscatter**. Working Paper.
- Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: **Nonlinear Binscatter Methods**. Working Paper.
- Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: **Binscatter Regressions**. Working Paper.

See Also

[binsreg](#), [binstest](#).

Examples

```
x <- runif(500); y <- sin(x)+rnorm(500)
est <- binsregselect(y,x)
summary(est)
```

binstest

Data-Driven Nonparametric Shape Restriction and Parametric Model Specification Testing using Binscatter

Description

binstest implements binscatter-based hypothesis testing procedures for parametric functional forms of and nonparametric shape restrictions on the regression function of interest, following the results in [Cattaneo, Crump, Farrell and Feng \(2023a\)](#) and [Cattaneo, Crump, Farrell and Feng \(2023b\)](#). If the binning scheme is not set by the user, the companion function [binsregselect](#) is used to implement binscatter in a data-driven way and inference procedures are based on robust bias correction. Binned scatter plots based on different methods can be constructed using the companion functions [binsreg](#), [binsqreg](#) or [binsglm](#).

Usage

```
binstest(y, x, w = NULL, data = NULL, estmethod = "reg",
  family = gaussian(), quantile = NULL, deriv = 0, at = NULL,
  nolink = F, testmodel = NULL, testmodelparfit = NULL,
  testmodelpoly = NULL, testshape = NULL, testshapel = NULL,
  testshaper = NULL, testshape2 = NULL, lp = Inf, bins = NULL,
  nbins = NULL, pselect = NULL, sselect = NULL, binspos = "qs",
  binsmethod = "dpi", nbinsrot = NULL, randcut = NULL, nsims = 500,
  simsgrid = 20, simsseed = NULL, vce = NULL, cluster = NULL,
  asyvar = F, dfcheck = c(20, 30), masspoints = "on", weights = NULL,
  subset = NULL, numdist = NULL, numclust = NULL, estmethodopt = NULL,
  ...)
```

Arguments

| | |
|-----------|--|
| y | outcome variable. A vector. |
| x | independent variable of interest. A vector. |
| w | control variables. A matrix, a vector or a formula . |
| data | an optional data frame containing variables used in the model. |
| estmethod | estimation method. The default is estmethod="reg" for tests based on binscatter least squares regression. Other options are "qreg" for quantile regression and "glm" for generalized linear regression. If estmethod="glm", the option family must be specified. |
| family | a description of the error distribution and link function to be used in the generalized linear model when estmethod="glm". (See family for details of family functions.) |

| | |
|-----------------|--|
| quantile | the quantile to be estimated. A number strictly between 0 and 1. |
| deriv | derivative order of the regression function for estimation, testing and plotting. The default is <code>deriv=0</code> , which corresponds to the function itself. |
| at | value of w at which the estimated function is evaluated. The default is <code>at="mean"</code> , which corresponds to the mean of w . Other options are: <code>at="median"</code> for the median of w , <code>at="zero"</code> for a vector of zeros. <code>at</code> can also be a vector of the same length as the number of columns of w (if w is a matrix) or a data frame containing the same variables as specified in w (when data is specified). Note that when <code>at="mean"</code> or <code>at="median"</code> , all factor variables (if specified) are excluded from the evaluation (set as zero). |
| nolink | if true, the function within the inverse link function is reported instead of the conditional mean function for the outcome. |
| testmodel | a vector or a logical value. It sets the degree of polynomial and the number of smoothness constraints for parametric model specification testing. If <code>testmodel=c(p,s)</code> is specified, a piecewise polynomial of degree p with s smoothness constraints is used. If <code>testmodel=T</code> or <code>testmodel=NULL</code> (default) is specified, <code>testmodel=c(1,1)</code> is used unless the degree p or the smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). |
| testmodelparfit | a data frame or matrix which contains the evaluation grid and fitted values of the model(s) to be tested against. The column contains a series of evaluation points at which the <code>binscatter</code> model and the parametric model of interest are compared with each other. Each parametric model is represented by other columns, which must contain the fitted values at the corresponding evaluation points. |
| testmodelpoly | degree of a global polynomial model to be tested against. |
| testshape | a vector or a logical value. It sets the degree of polynomial and the number of smoothness constraints for nonparametric shape restriction testing. If <code>testshape=c(p,s)</code> is specified, a piecewise polynomial of degree p with s smoothness constraints is used. If <code>testshape=T</code> or <code>testshape=NULL</code> (default) is specified, <code>testshape=c(1,1)</code> is used unless the degree p or smoothness s selection is requested via the option <code>pselect</code> or <code>sselect</code> (see more details in the explanation of <code>pselect</code> and <code>sselect</code>). |
| testshape1 | a vector of null boundary values for hypothesis testing. Each number a in the vector corresponds to one boundary of a one-sided hypothesis test to the left of the form $H_0: \sup_x \mu(x) \leq a$. |
| testshaper | a vector of null boundary values for hypothesis testing. Each number a in the vector corresponds to one boundary of a one-sided hypothesis test to the right of the form $H_0: \inf_x \mu(x) \geq a$. |
| testshape2 | a vector of null boundary values for hypothesis testing. Each number a in the vector corresponds to one boundary of a two-sided hypothesis test of the form $H_0: \sup_x \mu(x) - a = 0$. |
| lp | an L_p metric used for (two-sided) parametric model specification testing and/or shape restriction testing. The default is <code>lp=Inf</code> , which corresponds to the sup-norm of the t -statistic. Other options are <code>lp=q</code> for a positive integer q . |
| bins | a vector. If <code>bins=c(p,s)</code> , it sets the piecewise polynomial of degree p with s smoothness constraints for data-driven (IMSE-optimal) selection of the partitioning/binning scheme. The default is <code>bins=c(0,0)</code> , which corresponds to the piecewise constant. |

| | |
|------------|---|
| nbins | number of bins for partitioning/binning of x . If nbins=T or nbins=NULL (default) is specified, the number of bins is selected via the companion command binsregselect in a data-driven, optimal way whenever possible. If a vector with more than one number is specified, the number of bins is selected within this vector via the companion command binsregselect . |
| pselect | vector of numbers within which the degree of polynomial p for point estimation is selected. If the selected optimal degree is p , then piecewise polynomials of degree $p+1$ are used to conduct testing for nonparametric shape restrictions or parametric model specifications. <i>Note:</i> To implement the degree or smoothness selection, in addition to pselect or sselect, nbins=# must be specified. |
| sselect | vector of numbers within which the number of smoothness constraints s for point estimation is selected. If the selected optimal smoothness is s , then piecewise polynomials of $s+1$ smoothness constraints are used to conduct testing for nonparametric shape restrictions or parametric model specifications. If not specified, for each value p supplied in the option pselect, only the piecewise polynomial with the maximum smoothness is considered, i.e., $s=p$. |
| binspos | position of binning knots. The default is binspos="qs", which corresponds to quantile-spaced binning (canonical binscatter). The other options are "es" for evenly-spaced binning, or a vector for manual specification of the positions of inner knots (which must be within the range of x). |
| binsmethod | method for data-driven selection of the number of bins. The default is binsmethod="dpi", which corresponds to the IMSE-optimal direct plug-in rule. The other option is: "rot" for rule of thumb implementation. |
| nbinsrot | initial number of bins value used to construct the DPI number of bins selector. If not specified, the data-driven ROT selector is used instead. |
| randcut | upper bound on a uniformly distributed variable used to draw a subsample for bins/degree/smoothness selection. Observations for which $\text{runif}() \leq \#$ are used. # must be between 0 and 1. By default, $\max(5000, 0.01n)$ observations are used if the samples size $n > 5000$. |
| nsims | number of random draws for hypothesis testing. The default is nsims=500, which corresponds to 500 draws from a standard Gaussian random vector of size $[(p+1)*J - (J-1)*s]$. Setting at least nsims=2000 is recommended to obtain the final results. |
| simsgrid | number of evaluation points of an evenly-spaced grid within each bin used for evaluation of the supremum (infimum or L_p metric) operation needed to construct hypothesis testing procedures. The default is simsgrid=20, which corresponds to 20 evenly-spaced evaluation points within each bin for approximating the supremum (infimum or L_p metric) operator. Setting at least simsgrid=50 is recommended to obtain the final results. |
| simsseed | seed for simulation. |
| vce | procedure to compute the variance-covariance matrix estimator. For least squares regression and generalized linear regression, the allowed options are the same as that for binsreg or binsqreg . For quantile regression, the allowed options are the same as that for binsqreg . |
| cluster | cluster ID. Used for compute cluster-robust standard errors. |
| asyvar | if true, the standard error of the nonparametric component is computed and the uncertainty related to control variables is omitted. Default is asyvar=FALSE, that is, the uncertainty related to control variables is taken into account. |

| | |
|--------------|---|
| dfcheck | adjustments for minimum effective sample size checks, which take into account number of unique values of x (i.e., number of mass points), number of clusters, and degrees of freedom of the different stat models considered. The default is <code>dfcheck=c(20,30)</code> . See Cattaneo, Crump, Farrell and Feng (2023c) for more details. |
| masspoints | how mass points in x are handled. Available options: <ul style="list-style-type: none"> • "on" all mass point and degrees of freedom checks are implemented. Default. • "noadjust" mass point checks and the corresponding effective sample size adjustments are omitted. • "nolocalcheck" within-bin mass point and degrees of freedom checks are omitted. • "off" "noadjust" and "nolocalcheck" are set simultaneously. • "veryfew" forces the function to proceed as if x has only a few number of mass points (i.e., distinct values). In other words, forces the function to proceed as if the mass point and degrees of freedom checks were failed. |
| weights | an optional vector of weights to be used in the fitting process. Should be NULL or a numeric vector. For more details, see lm . |
| subset | optional rule specifying a subset of observations to be used. |
| numdist | number of distinct for selection. Used to speed up computation. |
| numclust | number of clusters for selection. Used to speed up computation. |
| estmethodopt | a list of optional arguments used by rq (for quantile regression) or glm (for fitting generalized linear models). |
| ... | optional arguments to control bootstrapping if <code>estmethod="qreg"</code> and <code>vce="boot"</code> . See boot.rq . |

Value

| | |
|--------------|--|
| testshapeL | Results for <code>testshapeL</code> , including: <code>testvalL</code> , null boundary values; <code>stat.shapeL</code> , test statistics; and <code>pval.shapeL</code> , p-value. |
| testshapeR | Results for <code>testshaper</code> , including: <code>testvalR</code> , null boundary values; <code>stat.shapeR</code> , test statistics; and <code>pval.shapeR</code> , p-value. |
| testshape2 | Results for <code>testshape2</code> , including: <code>testval2</code> , null boundary values; <code>stat.shape2</code> , test statistics; and <code>pval.shape2</code> , p-value. |
| testpoly | Results for <code>testmodelpoly</code> , including: <code>testpoly</code> , the degree of global polynomial; <code>stat.poly</code> , test statistic; <code>pval.poly</code> , p-value. |
| testmodel | Results for <code>testmodelparfit</code> , including: <code>stat.model</code> , test statistics; <code>pval.model</code> , p-values. |
| imse.var.rot | Variance constant in IMSE, ROT selection. |
| imse.bsq.rot | Bias constant in IMSE, ROT selection. |
| imse.var.dpi | Variance constant in IMSE, DPI selection. |
| imse.bsq.dpi | Bias constant in IMSE, DPI selection. |
| opt | A list containing options passed to the function, as well as total sample size n , number of distinct values N_{dist} in x , number of clusters N_{clust} , and number of bins n_{bins} . |

Author(s)

Matias D. Cattaneo, Princeton University, Princeton, NJ. <cattaneo@princeton.edu>.

Richard K. Crump, Federal Reserve Bank of New York, New York, NY. <richard.crump@ny.frb.org>.

Max H. Farrell, UC Santa Barbara, Santa Barbara, CA. <mhfarrell@gmail.com>.

Yingjie Feng (maintainer), Tsinghua University, Beijing, China. <fengyingjiepk@gmail.com>.

References

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023a: **On Binscatter**. Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023b: **Nonlinear Binscatter Methods**. Working Paper.

Cattaneo, M. D., R. K. Crump, M. H. Farrell, and Y. Feng. 2023c: **Binscatter Regressions**. Working Paper.

See Also

[binsreg](#), [binsqreg](#), [binsglm](#), [binsregselect](#).

Examples

```
x <- runif(500); y <- sin(x)+rnorm(500)
est <- binstest(y,x, testmodelpoly=1)
summary(est)
```

Index

`_PACKAGE` (binsreg-package), [2](#)

`binsglm`, [2](#), [3](#), [8](#), [12](#), [27](#), [31](#)

`binspwc`, [2](#), [8](#)

`binsqreg`, [2](#), [8](#), [11](#), [12](#), [12](#), [27](#), [29](#), [31](#)

`binsreg`, [2](#), [8](#), [11](#), [12](#), [18](#), [27](#), [29](#), [31](#)

`binsreg-package`, [2](#)

`binsregselect`, [2](#), [3](#), [5](#), [6](#), [8](#), [10](#), [12](#), [15](#), [18](#),
[20](#), [21](#), [23](#), [24](#), [27](#), [29](#), [31](#)

`binstest`, [2](#), [3](#), [8](#), [12](#), [18](#), [23](#), [27](#), [27](#)

`boot.rq`, [11](#), [16](#), [30](#)

`family`, [3](#), [9](#), [27](#)

`formula`, [3](#), [9](#), [13](#), [19](#), [24](#), [27](#)

`glm`, [7](#), [11](#), [30](#)

`lm`, [7](#), [11](#), [16](#), [22](#), [25](#), [30](#)

`rq`, [11](#), [16](#), [30](#)

`summary.rq`, [16](#)