

# PACKAGE NP FAQ

JEFFREY S. RACINE

This set of frequently asked questions is intended to help users who are encountering unexpected or undesired behavior when trying to use the `np` package.

Many of the underlying C routines have been extensively tested over the past two decades. However, the R ‘hooks’ needed to call these routines along with processing of the data required for a seamless user experience may produce unexpected or undesired results in some settings.

Kindly report any such issues to me, and please include your code and data so that I can help track down any such issues (racinej@mcmaster.ca). And, of course, if you encounter an issue that you think might be of interest to others, kindly email me the relevant information and I will incorporate it into this FAQ.

## 1. VERSION COVERED

This FAQ refers to the most recent version, which as of this writing is 0.30-4. Kindly update your version should you not be using the most current (from within R, `update.packages()` ought to do it, though also see 3 below.). See the appendix in this file for cumulative changes between this and previous versions of the `np` package.

## 2. FREQUENTLY ASKED QUESTIONS

### (1) *How do I cite the np package?*

Once you have installed the `np` package (`install.packages("np")`), if you load the `np` package and type `citation("np")` you will be presented with the following information.

```
> citation("np")
```

To cite `np` in publications use:

Tristen Hayfield and Jeffrey S. Racine (2008). Nonparametric Econometrics: The `np` Package. *Journal of Statistical Software* 27(5). URL <http://www.jstatsoft.org/v27/i05/>.

A BibTeX entry for LaTeX users is

---

*Date:* January 27, 2010.

```
@Article{,
  title = {Nonparametric Econometrics: The np Package},
  author = {Tristen Hayfield and Jeffrey S. Racine},
  journal = {Journal of Statistical Software},
  year = {2008},
  volume = {27},
  number = {5},
  url = {http://www.jstatsoft.org/v27/i05/},
}
```

- (2) *I have never used R before. Can you direct me to some introductory material that will guide me through the basics...*

There are many excellent introductions to the R environment with more on the way. First, I would recommend going directly to the R website (<http://www.r-project.org>) and looking under Documentation/Manuals (<http://cran.r-project.org/manuals.html>) where you will discover a wealth of documentation for R users of all levels. See also the R task views summary page (<http://cran.nedmirror.nl/web/views/index.html>) for information grouped under field of interest. A few documents that I mention to my students which are tailored to econometricians include Cribari-Neto & Zarkos (1999) [1], Racine & Hyndman (2002) [11] and Farnsworth (2006) [3], to name but a few.

Those looking for exemplar datasets outside of those contained in the `np` package are directed to the `Ecdat` [2] and `AER` [7] packages.

Often the best resource is right down the hall. Ask a colleague whether they use or know anyone who uses R, then offer to buy that person a coffee and along the way drop something like “I keep hearing about the R project...I feel like such a Luddite...”

- (3) *How do I keep all R packages on my system current?*

Run the command `update.packages(checkBuilt=TRUE,ask=FALSE)`, which will not only update all packages that are no longer current, but will also update all packages built under outdated installed versions of R, if appropriate.

- (4) *It seems that there are a lot of packages that must be installed in order to conduct econometric analysis (tseries, lmtest, np, etc.). Is there a way to avoid having to individually install each package individually?*

Certainly. The Comprehensive R Archive Network (CRAN) is a network of ftp and web servers around the world that store identical, up-to-date, versions of code and documentation for R. The CRAN task view for computational econometrics might

be of particular interest to econometricians. The econometric task view provides an excellent summary of both parametric and nonparametric econometric packages that exist for the R environment and provides one-stop installation for these packages.

See [cran.r-project.org/web/views/Econometrics.html](http://cran.r-project.org/web/views/Econometrics.html) for further information.

To automatically install a task view, the `ctv` package first needs to be installed and loaded, i.e.,

```
install.packages("ctv")
library("ctv")
```

The econometric task view can then be installed via `install.views()` and updated via `update.views()` (which first assesses which of the packages are already installed and up-to-date), i.e.,

```
install.views("Econometrics")
```

or

```
update.views("Econometrics")
```

- (5) *Is there a ‘gentle guide’ to the `np` package that contains some easy to follow examples?*

Perhaps the most gentle introduction is contained in the `np` package itself in the form of a ‘vignette’. To view the vignette run R, install the `np` package (`install.packages("np")`), then type `vignette("np",package="np")` to view or print the vignette (this vignette is essentially the article that appeared in the Journal of Statistical Software that describes the `np` package (Hayfield & Racine [5])).

In addition, you might be interested in the nonparametric econometrics primer (Racine (2008) [10]) which is available for download from my website.

- (6) *What is the difference between the `np` package and the previous stand-alone programs you wrote, namely, `N ©`, `NPREG ©`, and `NPDEN ©`?*

The `np` package is built from the same C library that underlies its predecessors `N ©`, `NPREG ©`, and `NPDEN ©`. In fact, R calls the compiled C code that underlies its predecessors (one of the beauties of R is that you can obtain the benefits of compiled code (i.e., speed) yet have access to the rich superset of R routines and R packages built by others). Therefore, there is no penalty in run-time when using R versus the stand alone precompiled binary programs `N ©`, `NPREG ©`, and `NPDEN ©` (unless of course the compiler or compiler flags differ from those used to build its predecessors).

- (7) *How can I read data stored in various formats such as Stata, SAS, Minitab, SPSS etc. into the R program?*

Install the foreign library via `install.packages("foreign")` then do something like

```
mydat <- read.dta("datafile.dta"),
```

where `datafile.dta` is the name of your Stata data file. Note that, as of version 0.8-34, the foreign package function `read.dta` supports reading files directly over the Internet making for more portable code. For instance, one could do something like

```
mydat <- read.dta(file="http://www.principlesofeconometrics.com/stata/mroz.dta")
```

as one could always do with, say, `read.table()`.

- (8) *I want to use so-and-so's semiparametric/nonparametric method, however, the np package does not include this particular method...*

This is why we have included the function `npksum()`, which exists so that you can create your own kernel objects and take advantage of underlying kernel methods implemented in the np package without having to write, say, C or Fortran code.

With the options available, you could create new nonparametric tests or even new kernel estimators. For instance, the convolution kernel option would allow you to replicate, say, the least squares cross-validation function for kernel density estimation found in `npudensbw()`. The function `npksum()` uses highly-optimized C code that strives to minimize its memory footprint, while there is low overhead involved when using repeated calls to this function.

See, by way of illustration, the example in the `npksum()` help file that conducts leave-one-out cross-validation for a local constant regression estimator via calls to the R function `nlm()`, and compares this to the `npregbw()` function.

If you wish to have a method incorporated into a future version of the np package, the best way to achieve this is to successfully code up the method using `npksum()`, briefly document it and write up an example, then send it to us. We will then, at our discretion, do our best to adapt and incorporate this into a future version and of course give credit where credit is due.

- (9) *Cross-validation takes forever, and I can't wait that long...*

This is the most common complaint from frustrated users coming to terms with numerically demanding statistical methods. I am fond of saying 'if you want the wrong answer, I can give it to you right away', but this wears thin quickly.

- (a) Some background may be in order. Cross-validation methods have run times that are exponential in the number of observations (of computational order  $n^2$  hence a doubling of the sample size will increase run time by a factor of four).

The solution I favor is to run the code in a parallel computing environment. The underlying C code for `np` is MPI-aware (MPI denotes the ‘message passing interface’, a popular parallel programming library that is an international standard), and we intend to develop the R `np` package in a parallel environment via the `Rmpi` wrapper.

- (b) Alternatively, you can use the method outlined in Racine (1993) [9]. The method is based on the fact that the unknown constant  $c_j$  (the ‘scale factor’) in the formula  $c_j \sigma_j n^{-1/(2p+r)}$  is independent of the sample size, so one can conduct bandwidth selection on random subsets and do this for a large number of subsets then take the mean/median over these subsets and feed the scale factor into the final routine for the entire sample. Below you will find simple R code that replicates the method using numerical search and resampling without replacement rather than the grid method outlined in [9] (both have equivalent properties but this is perhaps simpler to implement using the `np` package).

```
## Generate a moderately large dataset

set.seed(12345)
n <- 100000
x1 <- runif(n)
x2 <- runif(n)

y <- 1 + x1 + sin(pi*x2) + rnorm(n,sd=.1)

## Set the number of resamples and the subsample size

num.res <- 50
n.sub <- 250

## Create a storage matrix

bw.mat <- matrix(NA,nrow=num.res,ncol=2)

## Get the scale factors for resamples from the full sample of size n.sub

options(np.messages=FALSE)

for(i in 1:num.res) {

  cat(paste(" Replication", i, "of", num.res, "...\\r"))
```

```

bw.mat[i,] <- npregbw(y~x1+x2,regtype="ll",
                      subset=sample(n,n.sub))$sfactor$x
}

## A function to compute the median of the columns of a matrix

colMedians <- function(data) {
  colmed <- numeric(ncol(data))
  for(i in 1:ncol(data)) {
    colmed[i] <- median(data[,i])
  }
  return(colmed)
}

## Take the median scale factors

bw <- colMedians(bw.mat)

## The final model for the full dataset

model.res <- npreg(y~x1+x2,bws=bw,regtype="ll",bwscaling=TRUE)

```

- (c) Barring this, you can set the search tolerances to be a bit less terse (at the expense of potential accuracy, i.e., becoming trapped in local minima) by setting, say, `tol=0.1` and `ftol=0.1` in the respective bandwidth routine (see the docs for examples). Also, you can set `nmulti=1` which overrides the number of times the search procedure restarts from different random starting values (the default is to restart  $k$  times where  $k$  is the number of variables). Be warned, however, that this is *definitely not recommended* and should *be avoided at all costs* for all but the most casual examination of a relationship. One ought to use multistarting for any final results and never override default search tolerances *unless increasing multistarts beyond the default*. Results based upon exhaustive search *often differ dramatically* from that based on limited search achieved by overriding default search tolerances.
- (d) For those who like to tinker and who work on a \*NIX system with the gcc compiler suite, you can change the default compiler switches used for building R packages which may generate some modest improvements in run time. The default compiler switches are

```
-g -O2
```

and are set in the R-`*/etc/Makeconf` file (where `*/` refers to your R version number, e.g. R-2.10.1). You can edit this file and change these defaults to

```
-O3 -ffast-math -fexpensive-optimizations -fomit-frame-pointer
```

then reinstall the `np` package and you may experience some improvements in run time. Note that the `-g` flag turns on production of debugging information which can involve some overhead, so we are disabling this feature. This is not a feature used by the typical applied researcher but if you envision requiring this it is clearly trivial to re-enable debugging. I typically experience in the neighborhood of a 0-5% reduction in run time for data-driven bandwidth selection on a variety of systems depending on the method being used, though mileage will of course vary.

- (e) For those using Mac OS X Tiger/Leopard/Snow Leopard you can use the 64 bit version of R instead of the 32 bit (default) version of R. This results in roughly a 20-30% reduction in run time for many `np` routines. See <http://r.research.att.com/> for more information.

- (10) *Is there a way to figure out roughly how long cross-validation will take on a large sample?*

Certainly. You can run cross-validation on subsets of your data of increasing size and time each run, then estimate a double-log model of sample size on run time (run time can be approximated by a linear log-log model) and then use this to predict approximate run-time. The following example demonstrates this for a simple model, but you can modify it trivially for your data. Note that the larger is `n.max` the more accurate it will likely be. Note that we presume your data is in no particular order (if it is, you perhaps ought to shuffle it first). We plot the log-log model fit and prediction along with that expressed in hours.

```
## Set the upper bound (n.max > 100) for the sub-samples on which you
## will run cross-validation (perhaps n.max = 1000 (or 2000) ought to
## suffice). For your application, n will be your sample size
```

```
n <- 2000
n.max <- 1000
```

```
x <- runif(n)
y <- 1 + x + rnorm(n)
```

```
n.seq <- seq(100,n.max,by=100)
```

```

time.seq <- numeric(length(n.seq))

for(i in 1:length(n.seq)) {
  time.seq[i] <- system.time(npregbw(y~x,subset=seq(1:n.seq[i])))[3]
}

## Now fit a double-log model and generate/plot actual values plus
## prediction for n (i.e., approximate run time in hours)

log.time <- log(time.seq)
log.n <- log(n.seq)

model <- lm(log.time~log.n)

n.seq.aug <- c(n.seq,n)
time.fcst <- exp(predict(model,newdata=data.frame(log.n=log(n.seq.aug))))

par(mfrow=c(2,1))

plot(log(n.seq.aug),log(time.fcst),type="b",
      xlab="log(Sample Size)",
      ylab="log(Run Time)",
      main="Approximate Run Time (log seconds)")

plot(n.seq.aug,time.fcst/3600,type="b",
      xlab="Sample Size (n)",
      ylab="Hours",
      main="Approximate Run Time (hours)",
      sub=paste("Predicted run time for n =", n, "observations:",
                signif(time.fcst[length(time.fcst)]/3600, digits=2),
                "hours"))

```

- (11) *I notice your density estimator supports manual bandwidths, the normal-reference rule-of-thumb, likelihood cross-validation and least squares cross-validation, but not plug-in rules. How can I use various plug-in rules for estimating a univariate density?*

For the univariate case this is straightforward as the default R installation supports a variety of univariate plug-in bandwidths (see `?bw.nrd` for details). For example, `bs.SJ` computes the bandwidths outlined in

Sheather, S. J. and Jones, M. C. (1991) A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society series B*, **53**, 683-690.



Incorporating these univariate plug-in bandwidth selectors into the univariate density estimation routines in `np` is straightforward as the following code snippet demonstrates. They will certainly be faster than the likelihood and least-squares cross-validation approaches. Also, you may not wish the density estimate for all sample realizations, but rather for a shorter grid. The following example demonstrates both and is closer in spirit to the `density()` function in base R.

```
x <- rnorm(10000)
xeval <- seq(min(x),max(x),length=100)
f <- npudens(tdat=x,edat=xeval,bws=bw.SJ(x))
```

You might even be tempted to use these in the multivariate case via `bws=c(bw.SJ(x1),bw.SJ(x2),...)` though these would be optimal for univariate densities and most certainly not for a multivariate density. However, for exploratory purposes these may be of interest to some users.

- (12) *I wrote a program using np and it does not work as I expected...*

There exist a rather extensive set of examples contained in the docs. You can run these examples by typing `example("npfunctionname")` where `npfunctionname` is, say, `w`, as in `example("w")`. These examples all pass quality control and produce the expected results, so first see whether your problem fits into an existing example, and if not, carefully follow the examples listed in a given function for syntax issues etc.

If you are convinced that the problem lies with `np` (there certainly will be undiscovered ‘features’, i.e., bugs), then kindly send me your code and data so that I can replicate and help resolve the issue.

- (13) *Under Mac OS X, when I run a command no progress is displayed...*

This should no longer occur for `np` versions 0.30-0 and up. For previous versions, this reflected a peculiarity of console input/output (I/O) under Mac OS X. Note, however, that if you run R in a terminal rather than `Rgui` you will get the full \*NIX<sup>1</sup> experience.

- (14) *When some routines are running under MS Windows, R appears to be ‘not responding.’ It appears that the program is not ‘hung’, rather is simply computing. The previous stand-alone program (N ©) always displayed useful information...*

This should no longer occur for `np` versions 0.30-0 and up. For previous versions, this reflected a peculiarity of the R Windows GUI, and was not specific to the `np` package.

---

<sup>1</sup>\*NIX is often used to describe UNIX and other UNIX-like platforms (i.e., UNIX, BSD, and GNU/Linux distributions). I harbor strong preferences for \*NIX computing platforms.

From the R Windows FAQ...

“When using Rgui the output to the console seems to be delayed. This is deliberate: the console output is buffered and re-written in chunks to be faster and less distracting. You can turn buffering off or on from the ‘Misc’ menu or the right-click menu: `<Ctrl-W>` toggles the setting.”

- (15) *Some results take a while, and my MS Windows computer is sluggish while R is running...*

You can easily change the priority of your R job on the fly, just as you might under \*NIX. Pull up the task manager (`<Ctrl>-<Alt>-<Del>`), go to the process list, and find the process `Rgui.exe` (or `R.exe` if you are running `Rterm`), select this process by left clicking on it, then right clicking will bring up a menu, select **Set priority**, then change priority to **low** and hit `<ok>`. For lengthy jobs this will make your life much smoother, and you can, say, run multiple jobs in low priority with no sluggishness whatsoever for your other applications (useful for queuing a number of long jobs). Alternatively, you could permanently change the default priority of R under MS Windows by modifying the properties of your R desktop icon.

- (16) *A variable must be cast as, say, a factor in order for `np` to recognize this as an unordered factor. How can I determine whether my data is already cast as a factor?*

Use the `class()` function. For example, define `x <- factor(c("male","female"))`, then type `class(x)`.

- (17) *When I use `plot()` (`npplot()`) existing graphics windows are overwritten. How can I display multiple graphics plots in separate graphics windows so that I can see results from previous runs and compare that to my new run?*

Use the `dev.new()` command in between each call to `plot()`. This will leave the existing graphics window open and start a new one. The command `dev.list()` will list all graphics windows, and the command `dev.set(integer.foo)` will allow you to switch from one to another and overwrite existing graphics windows should you so choose.

- (18) *Sometimes `plot()` fails to use my variable names...*

This should not occur unless you are using the data frame method and not naming your variables (e.g., you are doing something like `data.frame(ordered(year))`). To correct this, name your variables in the respective data frame, as in

```
data <- data.frame(year=ordered(year),gdp)
```

so that the ordered factor appears as ‘year’ and not ‘ordered.year’

- (19) *Sometimes `plot()` appends my variable names with `.ordered` or `.factor`...*

See also 18 above.

- (20) *I specify a variable as `factor()` or `ordered()` in a data frame, then call this when I conduct bandwidth selection. However, when I try to plot the resulting object, it complains*

```
Error in eval(expr, envir, enclos) : object "variable" not found...
```

This arises because `plot()` (`npplot()`) tries to retrieve the variable from the environment but you have changed the definition when you called the bandwidth selection routine (e.g., `npregbw(y~x,data=dataframe)`).

To correct this, simply call `plot()` with the argument `data=dataframe` where `dataframe` is the name of your data frame.

- (21) *My np code gives errors when I attempt to run it...*

First, it is good practice to name all arguments (see the docs for examples) as in `npregbw(formula=y~x)` (i.e., explicitly call formula for functions that use named formulas). This will help the code return a potentially helpful error message.

Next, follow the examples listed at the end of each function help page closely (i.e., `?npreg` then scroll down to **Examples:**). See also 12 above.

- (22) *I have (numeric) 0/1 dummy variables in my parametric model. Can I just pass them to the np functions as I would in a parametric setting?*

In general, definitely not – you need to correctly classify each variable as type `factor` and treat it as one variable only. By way of example, suppose in your data you have created dummy variables for year, for example, `dummy06` which equals 1 for 2006, 0 otherwise, `dummy07` which equals 1 for 2007, 0 otherwise etc. We create these by habit for parametric models. But, the underlying variable is simply year, which equals 2006, 2007, and so forth.

In `np` (and `R` in general), you get to economize by just telling the function that the variable ‘year’ is ordered, as in `ordered(year)`, where `year` is a vector containing elements 2006, 2007 etc. Of course, seasoned `R` users would appreciate that this is in fact the simple way to do it with a parametric model as well.

You would *never*, therefore, just pass dummy variables to an `np` function as you would for linear parametric models. The *only* exception is where you have only one 0/1 dummy for one variable, say ‘sex’, and in this case you still would have to enter this as `factor(sex)` so that the `np` function recognizes this as a factor (otherwise it would treat it as continuous and use a kernel function that is inappropriate for a factor).

- (23) *I have a categorical variable, ignored the advice in items 16 and 22, and R terminates with the following error.*

```
** Fatal Error in routine kernel_bandwidth() ** variable 0 appears to be constant  
** Program terminated abnormally!
```

This can only occur in the case where a variable is ‘pathological’ in that it has  $IQR=0$  but  $std>0$  (see the section titled ‘Changes from Version 0.30-1 to 0.30-2 [19-Apr-2009]’ for further details). This should no longer occur for np versions 0.30-2 and up.

- (24) Rgui appears to crash. There must be a bug in the program.

Try running your code using the terminal (i.e., `Rterm` in Windows or `R` in a Mac OS X terminal) and see whether you get the message in item 23.

- (25) *Can I skip creating a bandwidth object and enter a bandwidth directly?*

Certainly, though I would advise doing so for exploratory data analysis only. For example, attach a dataset via

```
data(cps71)  
attach(cps71)  
then enter, say,
```

```
plot(age,logwage,main="Manual Bandwidth Example")  
lines(age,fitted(npreg(logwage~age,bws=1)),col="blue",lty=1)  
lines(age,fitted(npreg(logwage~age,bws=2)),col="red",lty=2)  
lines(age,fitted(npreg(logwage~age,bws=3)),col="green",lty=3)  
legend(20,15,  
      c("h=1","h=2","h=3"),  
      col=c("blue","red","green"),  
      lty=c(1,2,3))
```

to plot the local constant estimator with bandwidths of 1, 2, and 3 years. Note that the age variable is already sorted in this dataset. If your data is not sorted you will need to do so prior to plotting so that your `lines` command works properly. Or see 26 below for a multivariate example.

- (26) *When I estimate my gradients and there are two or more covariates and then extract them with the `gradients()` function, they are not ‘smooth’, though if I plot a model with the `gradients=TRUE` option, they are. The `gradients()` function must be broken...*

The function `plot()` (`npplot()`) plots ‘partial’ means and gradients. In other words, it plots  $x_1$  versus  $\hat{g}(x_1, \bar{x}_2)$  for the partial mean, where  $\bar{x}_2$  is, say, the median/modal value of  $x_2$ . It also plots  $x_1$  versus  $\partial \hat{g}(x_1, \bar{x}_2) / \partial x_1$  for the gradient.

Note that we are controlling for the values of the other covariate(s). This is in effect what people expect when they play with linear parametric models of the form  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$  since, given the additive nature of the model,  $\partial y / \partial x_1 = \beta_1$  (i.e., does not vary with  $x_2$ ).

The example below shows how you could manually generate the partial gradients (and means) for your data where the sample realizations form the evaluation data for  $x_1$  (unlike `npplot()` which uses an evenly spaced grid). Note we use the function `uocquantile()` to generate a vector that holds  $x_2$  constant at its median/modal value (i.e., the 0.5 quantile) in the evaluation data. The function `uocquantile()` can compute quantiles for ordered, unordered, and continuous data (see `?uocquantile` for details).

```
n <- 100

x1 <- runif(n)
x2 <- runif(n)
y <- x1^2+x2^2 + rnorm(n,sd=.1)

data.train <- data.frame(x1,x2,y)

bw <- npregbw(y~x1+x2,
              data=data.train,
              regtype="ll",
              bwmethod="cv.aic")

data.eval <- data.frame(x1 = sort(x1),
                       x2 = rep(uocquantile(x2,.5),n))

model <- npreg(bws=bw,
              data=data.train,
              newdata=data.eval,
              gradients=TRUE)

plot(data.eval[,1],model$grad[,1],xlab="X1",ylab="Gradient",type="l")
```

- (27) *I use `plot()` (`npplot()`) to plot, say, a density and the resulting plot looks like an inverted density rather than a density...*

This can occur when the data-driven bandwidth is dramatically undersmoothed. Data-driven (i.e., automatic) bandwidth selection procedures are not guaranteed always to produce good results due to perhaps the presence of outliers or the rounding/discretization of continuous data, among others. By default, `npplot()` takes the two extremes of the data (minimum, maximum i.e., actual data points) then creates an equally spaced grid of evaluation data (i.e., not actual data points in general) and computes the density for these points. Since the bandwidth is extremely small, the density estimate at these evaluation points is correctly zero, while those for the sample realizations (in this case only two, the min and max) are non-zero, hence we get two peaks at the edges of the plot and a flat bowl equal to zero everywhere else.

This can also happen when your data is heavily discretized and you treat it as continuous. In such cases, treating the data as ordered may result in more sensible estimates.

- (28) *Can `npksum()` compute analytical derivatives with respect to a continuous variable?*

As of version 0.20-0 and up, yes it can, using the `operator = "derivative"` argument, which is put to its paces in the following code snippet (this supports multiple arguments including `"integral"` and `"convolution"` in addition to `"normal"`, the default).

```
Z <- seq(-2.23,2.23,length=100)
Zc <- seq(-4.47,4.47,length=100)

par(mfrow=c(2,2))

plot(Z,main="Kernel",ylab="K()",npksum(txdat=0,exdat=Z,bws=1,
    ckertype="epanechnikov",ckerorder=2,operator="normal")$ksum,
    col="blue",type="l")

plot(Z,main="Kernel Derivative",ylab="K()",npksum(txdat=0,exdat=Z,bws=1,
    ckertype="epanechnikov",ckerorder=2,operator="derivative")$ksum,
    col="blue",type="l")

plot(Z,main="Kernel Integral",ylab="K()",npksum(txdat=0,exdat=Z,bws=1,
    ckertype="epanechnikov",ckerorder=2,operator="integral")$ksum,
    col="blue",type="l")

plot(Zc,main="Kernel Convolution",ylab="K()",npksum(txdat=0,exdat=Zc,bws=1,
    ckertype="epanechnikov",ckerorder=2,operator="convolution")$ksum,
    col="blue",type="l")
```

An alternative to computing analytical derivatives is to compute them numerically using finite-differences. One simply computes the kernel sum evaluating the sum with variable  $j$  set at  $x_j - h_j/2$  and calls this, say,  $ksum_{j1}$ , then again set at  $x_j + h_j/2$  and call this  $ksum_{j2}$ , then compute  $\nabla = (ksum_{j2} - ksum_{j1})/h_j$ . This method has been used for both theoretical and applied work and produces consistent estimates of the derivatives, as of course do the analytical derivatives, providing that  $h \rightarrow 0$  as  $n \rightarrow \infty$  (which will *not* be the case in some settings, i.e., in the presence of irrelevant covariates and so forth). The following example provides a simple demonstration. See 26 above for multivariate partial regression when using this method.

```
## In this example we consider the local constant estimator computed
## using npksum, and then use npksum to compute numerical derivatives
## using finite-difference methods, then finally compare them with the
## analytical ones.

data(cps71)
attach(cps71)

## Grab the cross-validated bandwidth

bw <- npregbw(logwage~age)
h <- bw$bw[1]

## Evaluate the local constant regression at x-h/2, x+h/2...

ksum.1 <- npksum(txdat=age, exdat=age-h/2,tydat=logwage,bws=bw)$ksum/
  npksum(txdat=age,exdat=age-h/2,bws=bw)$ksum

ksum.2 <- npksum(txdat=age, exdat=age+h/2,tydat=logwage,bws=bw)$ksum/
  npksum(txdat=age,exdat=age+h/2,bws=bw)$ksum

## Compute the numerical gradient...

grad.numerical <- (ksum.2-ksum.1)/h

## Compare with the analytical gradient...

grad.analytical <- gradients(npreg(bws=bw,gradient=TRUE))

## Plot the resulting estimates...
```

```
plot(age,grad.numerical,type="l",col="blue",lty=1,ylab="gradient")
lines(age,grad.analytical,type="l",col="red",lty=2)
legend(20,-0.05,c("Numerical","Analytic"),col=c("blue","red"),lty=c(1,2))
```

- (29) *Can I use the `npcmstest()` function that implements the consistent test for correct specification of parametric regression models as described in Hsiao, Li, & Racine (2007) [6] to test for correct specification of the semiparametric partially linear model?*

As Brennan Thompson points out, yes, you can.

To test a parametric linear specification against a semiparametric partially linear alternative, i.e.,

$$H_0 : y = X'\beta + Z'\gamma + u$$

$$H_1 : y = X'\beta + g(Z) + u,$$

you could use `npcmstest()` as follows:

```
lmodel <- lm(y~X+Z,y=TRUE,x=TRUE)
uhat <- resid(lmodel)
npcmstest(xdat=Z,ydat=uhat,model=lmodel)
```

A slightly better way (as discussed in Li & Wang (1998) [8]) would be to use a ‘mixed’ residual, i.e.,  $\hat{u}_i = y_i - X_i'\tilde{\beta} - Z_i'\hat{\gamma}$  in the test, where  $\tilde{\beta}$  is the semiparametric estimator of  $\beta$  (based on the semiparametric partially linear model), and  $\hat{\gamma}$  is the OLS estimator of  $\gamma$  based on the linear model. This could lead to potential power gains due to the improved efficiency of  $\hat{\beta}$  under the alternative.

- (30) *I am using `npcmstest()` on a `glm()` object (the document says `glm()` objects are supported and I am estimating a Logit model) but it returns an error saying*

**Error in eval(expr, envir, enclos) : y values must be 0 <= y <= 1.**

`npcmstest()` supports conditional mean models with continuous outcomes (`glm()` objects are supported so that models that are nonlinear in parameters can be estimated). The test is based on residual bootstrapping to generate resamples for  $Y$ . In particular, a resample for the residual vector ( $\hat{\epsilon}^*$ ) is added to the models’ fit (i.e.,  $Y^* = \hat{Y} + \hat{\epsilon}^*$ ) to generate a resample under the null. This excludes binary outcome models and the like because you would have generated a resample for  $Y$  that no longer contains zeros and ones, hence the error message.

- (31) *I want to plot the kernel function itself. How can I do this?*

Use the `npksum()` function and switch the evaluation and training roles as in the following example that plots the 2nd, 4th, 6th and 8th order Epanechnikov kernels.

```
Z <- seq(-sqrt(5),sqrt(5),length=100)
```



```

par(mfrow=c(2,2))
plot(Z,ylab="kernel",npksum(txdat=0,exdat=Z,bws=1,ckertype="epanechnikov",
ckerorder=2)$ksum,type="l",main="Epanechnikov [order = 2]")
plot(Z,ylab="kernel",npksum(txdat=0,exdat=Z,bws=1,ckertype="epanechnikov",
ckerorder=4)$ksum,type="l",main="Epanechnikov [order = 4]")
plot(Z,ylab="kernel",npksum(txdat=0,exdat=Z,bws=1,ckertype="epanechnikov",
ckerorder=6)$ksum,type="l",main="Epanechnikov [order = 6]")
plot(Z,ylab="kernel",npksum(txdat=0,exdat=Z,bws=1,ckertype="epanechnikov",
ckerorder=8)$ksum,type="l",main="Epanechnikov [order = 8]")

```

- (32) *In version 0.20-0 and up I can ‘combine’ steps such as bandwidth selection and estimation. But when I do `summary(model)` I don’t get the same summary that I would get from, say, `summary(bw)` and then `summary(model)`. How do I get bandwidth object summaries when combining steps?*

Don’t worry, the bandwidth object exists when you do the combined steps and is easily accessed via `summary(model$bws)` or extracted via `bw <- model$bws` where `model` is the name of your model.

- (33) *I estimated a semiparametric index model via `model <- npindex(y~x1+x2)` but `se(model)` returns `NULL`.*

Single-index models do not have asymptotic standard errors implemented at the moment (and I would never recommend using asymptotic standard errors to begin with). You can, however, get robust bootstrapped standard errors by adding the argument `errors = TRUE` to your call to `npindex` so that `se(model)` returns the vector of standard errors for the estimated conditional mean where `model` is the name of your model. Note that `model$merr` will contain the standard errors returned by `se(model)` while `model$gerr` will return the matrix of standard errors for the gradients *provided* you have set `gradients=TRUE` (furthermore, `model$mean.grad` and `model$mean.gerr` will give the average derivative and its bootstrapped standard errors). See the documentation of `npindex` for further details. The following code snippet demonstrates how one could do this for a simulated dataset.

```

n <- 100
x1 <- runif(n, min=-1, max=1)
x2 <- runif(n, min=-1, max=1)
y <- x1 - x2 + rnorm(n)

bw <- npindexbw(formula=y~x1+x2)
model <- npindex(bws=bw,errors=TRUE,gradients=TRUE)

```

```
se.mean <- model$merr
se.grad <- model$gerr
```

(34) *How do I interpret gradients from the conditional density estimator?*

If you plot a conditional density  $f(y|x)$  when  $x$  is a scalar, with gradients, by default you will get the following:

- (a) A plot of  $\partial f(y = \text{median}|x)/\partial x$  (admittedly not the most useful plot).  
(If  $y$  is discrete the only difference is that you get a plot of  $\partial f(y = (\text{unconditional}) \text{ mode}|x)/\partial x$ ).
- (b) A plot of  $\partial f(y|x = \text{median})/\partial x$ .

If  $x$  is multivariate (for example, 2D) you get:

- (a) A plot of  $\partial f(y = \text{median}|x_1, x_2 = \text{median})/\partial x_1$
- (b) A plot of  $\partial f(y = \text{median}|x_1, x_2 = \text{median})/\partial x_2$
- (c) A plot of  $\partial f(y = \text{median}|x_1 = \text{median}, x_2)/\partial x_1$
- (d) A plot of  $\partial f(y = \text{median}|x_1 = \text{median}, x_2)/\partial x_2$
- (e) A plot of  $\partial f(y|x_1 = \text{median}, x_2 = \text{median})/\partial x_1$
- (f) A plot of  $\partial f(y|x_1 = \text{median}, x_2 = \text{median})/\partial x_2$

(35) *When I run R in batch mode via `R CMD BATCH filename.R` unwanted status messages (e.g., “Multistart 1 of 10”) crowd out desired output. How can I turn off these unwanted status messages?*

After loading the `np` library add the line `options(np.messages=FALSE)` and all such messages will be disabled.

(36) *I am getting an ‘unable to allocate...’ message after repeatedly interrupting large jobs.*

Repeated interruption of large jobs can reduce available memory under R. This occurs because memory is allocated dynamically, and memory that has been allocated is not freed when you interrupt the job (the routine will clean up after itself only if it is allowed to complete all computations - when you interrupt you never reach the point in the code where the memory is freed). If this becomes an issue simply restart R (i.e., exit then run a fresh R session).

(37) *I have a large number of variables, and when using the formula interface I get an error message stating ‘invoked with improper formula’. I have double checked and everything looks fine. Furthermore, it works fine with the data frame interface.*

The issue is that converting formulas into character strings in R appears to be limited to 500 characters. We are not aware of a simple workaround so we simply advise that you use the data frame interface when this occurs.

(38) *How can I estimate additive semiparametric models?*

Generalized additive semiparametric models (see Hastie and Tibshirani (1990) [4]) are implemented in the `gam` package (though they do not support categorical variables). The `gam` package function `gam()` (the `mgcv` package also contains a similar function by the same name) uses iteratively reweighted least squares and either smoothing splines (`s()`) or local polynomial regression fitting (`'loess'`, `lo()`, with default manual `'span'` of 0.75, the parameter which controls the degree of smoothing). The following code snippet demonstrates the capabilities of the `gam()` function via the `wage1` dataset included in the `np` packages using three numeric regressors.

```
library(gam)
data(wage1)
attach(wage1)
model.gam <- gam(lwage~s(educ)+s(exper)+s(tenure))
par(mfrow=c(2,2))
plot(model.gam,se=T)
detach(wage1)
```

The `mgcv` package also has an implementation of generalized additive models via the identical function name, `gam()`. The biggest difference is that this uses generalized cross validation to select the `'span'` (rather than manually) hence the degree of smoothing becomes part of the method (as is the case for all functions in the `np` package). The following code snippet demonstrates the `mgcv` implementation of the `gam()` function.

```
library(mgcv)
data(wage1)
attach(wage1)
model.gam <- gam(lwage~s(educ)+s(exper)+s(tenure))
par(mfrow=c(2,2))
plot(model.gam,se=T)
detach(wage1)
```

## REFERENCES

- [1] Francisco Cribari-Neto and Spyros G Zarkos. R: Yet another econometric programming environment. *Journal of Applied Econometrics*, 14(3):319–29, May–June 1999. Available at <http://ideas.repec.org/a/jae/japmet/v14y1999i3p319-29.html>.
- [2] Yves Croissant. *Ecdat: Data sets for econometrics*, 2006. R package version 0.1-5.
- [3] Grant V. Farnsworth. Econometrics in R. Technical report, June 2006. Available at <http://cran.r-project.org/doc/contrib/Farnsworth-EconometricsInR.pdf>.

- [4] T. Hastie and R. Tibshirani. *Generalized additive models*. Chapman and Hall, London, 1990.
- [5] Tristen Hayfield and Jeffrey S. Racine. Nonparametric econometrics: The np package. *Journal of Statistical Software*, 27(5), 2008.
- [6] C. Hsiao, Q. Li, and J. Racine. A consistent model specification test with mixed categorical and continuous data. *Journal of Econometrics*, 140:802–826, 2007.
- [7] Christian Kleiber and Achim Zeileis. *Applied Econometrics with R*. Springer-Verlag, New York, 2008. ISBN 978-0-387-77316-2.
- [8] Q. Li and S. Wang. A simple consistent bootstrap test for a parametric regression functional form. *Journal of Econometrics*, 87:145–165, 1998.
- [9] J. S. Racine. An efficient cross-validation algorithm for window width selection for nonparametric kernel regression. *Communications in Statistics*, 22(4):1107–1114, October 1993.
- [10] J. S. Racine. Nonparametric econometrics: A primer. *Foundations and Trends in Econometrics*, 3(1):1–88, 2008.
- [11] J. S. Racine and R. Hyndman. Using R to teach econometrics. *Journal of Applied Econometrics*, 17(2):175–189, 2002.

## CHANGES FROM VERSION 0.30-3 TO 0.30-4 [27-JAN-2010]

- Added function `npsymtest` for entropy-based testing of symmetry described in Maasoumi and Racine (2009), Econometric Reviews
- Added function `b.star` that automates block length selection for the stationary and circular bootstrap
- Cleaned up docs

## CHANGES FROM VERSION 0.30-2 TO 0.30-3 [28-MAY-2009]

- Corrected error in Epanechnikov convolution kernels for fixed and generalized bandwidth objects
- Changed default example in `npscoef`

## CHANGES FROM VERSION 0.30-1 TO 0.30-2 [19-APR-2009]

- `min(std, IQR/1.348)` is the adaptive measure of spread. We now test for the pathological case where `IQR=0` but `std>0` and return `std` in this instance

## CHANGES FROM VERSION 0.30-0 TO 0.30-1 [29-JAN-2009]

- `predict()` now supports bandwidth, density, distribution, conbandwidth, condensity, and condistribution objects
- Consistently allow predictions for categorical values outside of support of training data

Note that predictions based upon unconditional density objects defined over categorical variables that lie outside the support of the training data may no longer be true probabilities (i.e., as defined over the training data and the extended/augmented support – their sum may exceed one) and may therefore require renormalization by the user

- Fixed a numerical issue which could hinder `npregbw()`'s cross validation with higher-order kernels
- Default `nmulti` in `npplregbw()` is now set correctly
- Fixed a bug with the ridging routine in `npscoefbw()`, added ridging to `npscoef()`
- Fixed minor i/o issue with `Multistart 1 of...` using `npscoefbw()`

## CHANGES FROM VERSION 0.20-4 TO 0.30-0 [15-JAN-2009]

- Added basic user-interrupt checking for all underlying C code so that either `<Ctrl-C>` (Rterm) or the 'STOP' icon (Rgui) will interrupt all running processes. This has

a number of desirable side effects in addition to being able to interrupt C-based processes including i) R no longer showing up as ‘not responding’ under the task manager (Windows) or the activity monitor (Mac OS X) and ii) buffered output now being correctly displayed when using Rgui under Windows and Mac OS X

Note that repeated interruption of large jobs can reduce available memory under R - if this becomes an issue (i.e., you get a ‘cannot allocate...’ error under R) simply restart R (i.e., exit then run a fresh R session)

- Added a function `npseed()` that allows the user to set/reset the random seed for all underlying C routines
- Fixed a bug that caused `npplregbw()` to ignore any kernel options for the regression of  $y$  on  $z$
- Refined certain constants used in the normal-reference density bandwidth rule for increased accuracy
- Moved from using the maximum likelihood estimate of variance throughout to the degrees of freedom corrected estimate (all variance estimates now change by the factor  $(n-1)/n$ )

#### CHANGES FROM VERSION 0.20-3 TO 0.20-4 [19-Nov-2008]

- Using an adaptive measure of spread throughout. The scale factor reported for a bandwidth can appear to be small when the sample standard deviation of the associated variable is inflated due to the presence of outliers. Furthermore, supplying a scale factor of, say, 1.06 for density estimation when there are outliers that inflate the standard deviation may oversmooth rather dramatically in the presence of outliers. We now use the measure found in Silverman (1986, equation (3.30)) which is  $\min(\text{standard deviation}, \text{interquartile range}/1.349)$ . This robust choice produces expected results for scale factors in the presence of outliers

#### CHANGES FROM VERSION 0.20-2 TO 0.20-3[14-Nov-2008]

- Fixed a typo which caused `predict()` and `plot()` to abort when called on `plregression` objects, and which also prevented `print()` and `summary()` from printing information about the kernels used when called on `plregression` objects
- Fixed a typo which caused partially linear regressions to crash when out-of-sample responses were provided with evaluation data

## CHANGES FROM VERSION 0.20-1 TO 0.20-2 [02-NOV-2008]

- Allow for evaluation outside of discrete support of factors in `npksum()` and fixed a warning in `jksum`
- Fixed a bug which lead to unpredictable behavior when there were more categorical values for the training data than realisations

## CHANGES FROM VERSION 0.20-0 TO 0.20-1 [13-AUG-2008]

- Work-around for scale-factor issues during `npregbw()` cv when changing the training data

## CHANGES FROM VERSION 0.14-3 TO 0.20-0 [28-JUL-2008]

- `npksum()` now supports an expanded set of kernels (including convolution, derivative and integral), which can be selected via the `operator =` argument
- Automatic bandwidth searches are now performed when attempting to evaluate on data without bandwidths. This allows users to combine bandwidth selection and estimation in one step
- The `npsigtest()` interface is brought in line with other functions (S3)
- Significance tests can now be performed on `npreg()` outputs, so `npsigtest(modelname)` is now supported
- Added a vignette and faq. To see the vignette try `vignette("np", package="np")`
- `summary()` on `npconmode()` now properly retrieves names from bandwidth objects
- Fixed the 6th and 8th order epanechnikov kernels
- Fixed some quietness issues
- `npplot()` now returns data upon request for conditional densities
- `npreg()` and `npcdens()` now take the appropriate limits in some pathological cases
- User supplied bandwidths now operate seamlessly with the formula interface

## CHANGES FROM VERSION 0.14-2 TO 0.14-3 [02-MAY-2008]

- Fixed a glitch that only arose when using the `liracine` unordered kernel in the presence of irrelevant variables. The upper bound for numerical search was constrained to be  $(c-1)/c$  [that for the `aitchisonaitken` unordered kernel] but ought to have been 1. The summary output would therefore show a value of lambda hitting the (smaller) upper bound  $(c-1)/1$  when it may have hit the (larger) upper bound 1

## CHANGES FROM VERSION 0.14-1 TO 0.14-2 [11-JAN-2008]

- Relaxed checking tolerances slightly to prevent spurious 'invalid bandwidth' errors
- Empty sections were removed from help files
- `example(foobar)` now works again. This was disabled in 0.14-1 at the request of the R maintainers in order to shorten the duration of R CMD check. All examples remained in the help files but due to the presence of 'dontrun' they were not run when `example(foobar)` is requested. Now a limited subset is run while the full set of examples remain in the documents

## CHANGES FROM VERSION 0.13-1 TO 0.14-1 [18-DEC-2007]

- Now use `optim()` for minimisation in single index and smooth coefficient models
- Fixed bug in klein-spady objective function
- Standard errors are now available in the case of no continuous variables
- Summary should look prettier, print additional information
- Tidied up lingering issues with out-of-sample data and conditional modes
- Fixed error when plotting asymptotic errors with conditional densities
- Fixed a bug in `npplot()` with partially linear regressions and `plot.behavior="data"` or `"plot-data"`
- Maximum default number of multistarts is 5
- Least-squares cross-validation of conditional densities uses a new, much faster algorithm
- New, faster algorithm for least-squares cross-validation for both local-constant and local linear regressions

Note that the estimator has changed somewhat: both cross-validation and the estimator itself use a method of shrinking towards the local constant estimator when singularity would otherwise lead to the breakdown of the estimator. This arises in sparse data settings in conjunction with small bandwidths for one or more regressor

- Optimised smooth coefficient code, added ridging
- Fixed bug in uniform CDF kernel
- Fixed bug where `npindexbw()` would ignore `bandwidth.compute = FALSE` and compute bandwidths when supplied with a preexisting bw object
- Now can handle estimation out of discrete support
- Summary would misreport the values of discrete scale factors which were computed with `bwscaling = TRUE`



## CHANGES FROM VERSION 0.12-1 TO 0.13-1 [03-MAY-2007]

- Bandwidths are now checked for validity based on their variable and kernel types
- np now does a better job of preserving names of some 'y' data
- Names of coefficients returned from `coef()` now match variable names
- Fixed some corner cases in `npksum()` involving the dimensionality of outputs
- Fixed deprecation warnings in R 2.5.0 caused by use of \$ on atomic objects
- Various and sundry bug fixes in `npscoef()`
- `npscoef()` now handles discrete 'z' data
- Predict now accepts the argument 'se.fit', like `predict.lm`
- Fixed bug where incorrect asymptotic standard errors of gradients for regression objects were being displayed in `npplot()`
- Fixed bug where errors of gradients of regression objects were not being returned in matrix form
- `vcov()` now works with partially linear regression objects
- Fixed detection of evaluation responses when using the formula interface
- Pre-computed bandwidth objects are now provided for some of the more computationally burdensome examples
- Added Jeffrey Wooldridge's WAGE1 dataset with qualitative variables (married, female, nonwhite)
- Predictions outside of discrete support for regressions and conditional densities are now allowed
- Fixed sign issue with scaling of standard errors in the single index model
- Fixed error when calculating some bandwidths/scale factors for display purposes
- Bug in passing certain arguments to `npcdensbw()` fixed
- Added predict method for qregression objects
- Proper normalisation for lracine kernel shown in summary
- Fixed output bug ( $\hat{H}$ ) in summary method for sigtest objects
- Fixed regression with plotting of bootstrapped errors in perspective plots
- `npcdist()` no longer incorrectly calls `npcdens()`
- Fixed spacing between var name and p-value in significance test summaries

## VERSION 0.12-1 [19-Nov-2006]

- Initial release of the np package on CRAN