

# Refactoring the `nls()` function in R

John C. Nash \*

Arkajyoti Bhattacharjee †

2022-5-29

## Contents

<b>Abstract</b>	<b>1</b>
<b>The <code>nls()</code> function: strengths and shortcomings</b>	<b>2</b>
Feature: Convergence and termination tests (FIXED)	2
Feature: Failure when Jacobian is computationally singular	3
Feature: Jacobian computation	4
Feature: Subsetting	4
Feature: <code>na.action</code>	4
Feature: model frame	5
Feature: sources of data	5
Feature: missing start vector and self-starting models	5
Issue: returned output of <code>nls()</code> and its documentation	6
Feature: partially linear models and their specification	8
Issue: code structure and documentation for maintenance	9
Feature: indexed parameters	10
Feature: bounds on parameters	12
<b>Goals of our effort to improve <code>nls()</code></b>	<b>14</b>
Code rationalization and documentation	14
Provide tests and use-case examples	15
<b>Progress so far</b>	<b>16</b>
Reports and articles	16
Problem sets and test infrastructure	16
Code and documentation	16
Strategic choices in a nonlinear model estimation	17
<b>Acknowledgement</b>	<b>18</b>
<b>References</b>	<b>18</b>

## Abstract

Based on work for a Google Summer of Code project “Improvements to `nls()`”, we consider the features and limitations of the R function `nls()` in the context of improving and rationalizing R tools for nonlinear regression.

---

\*retired professor, Telfer School of Management, University of Ottawa

†Department of Mathematics and Statistics, Indian Institute of Technology, Kanpur

An important consideration is the maintainability of the code base that provides the functionality `nls()` claims to offer. Our work suggests that the existing code makes maintenance and improvement very difficult, with legacy applications and examples blocking some important updates. Discussion of these matters is relevant to improving R generally, as well as its nonlinear estimation tools.

## The `nls()` function: strengths and shortcomings

`nls()` is the tool in base R (the distributed software package from [CRAN](https://cran.r-project.org/) for estimating nonlinear statistical models. The function dates to the 1980s and the work related to Bates and Watts (1988) in S (see [https://en.wikipedia.org/wiki/S\\_\(programming\\_language\)](https://en.wikipedia.org/wiki/S_(programming_language))).

The `nls()` function has a remarkable and comprehensive set of capabilities for estimating nonlinear models that are expressed as **formulas**. In particular, we note that it

- handles formulas that include R functions
- allows data to be subset
- permits parameters to be indexed over a set of related data
- produces measures of variability (i.e., standard error estimates) for the estimated parameters
- has related profiling capabilities for exploring the likelihood surface as parameters are changed
- links to a number of pre-coded (“selfStart”) models

With such a range of features and long history, the code has become untidy and overly patched. It is, to our mind, essentially unmaintainable. Moreover, its underlying methods can and should be improved, as we suggest below.

### Feature: Convergence and termination tests (FIXED)

A previous issue with `nls()` that prevented it from providing parameter estimates for zero-residual (i.e., perfect fit) data was corrected thanks to suggestions by one of us.

In the manual page for `nls()` in R 4.0.0 there is the warning

**Do not use `nls` on artificial “zero-residual” data.**

The `nls` function uses a relative-offset convergence criterion that compares the numerical imprecision at the current parameter estimates to the residual sum-of-squares. This performs well on data of the form

$$y = f(x, \theta) + \textit{eps}$$

(with `var(eps) > 0`). It fails to indicate convergence on data of the form

$$y = f(x, \theta)$$

because the criterion amounts to comparing two components of the round-off error.

If you wish to test `nls` on artificial data please add a noise component, as shown in the example below.

This amounts to admitting R cannot solve well-posed problems unless data is polluted with errors.

This issue can be easily resolved. The “termination test” for the **program** rather than for “convergence” of the underlying **algorithm** is the Relative Offset Convergence Criterion (see Bates, Douglas M. and Watts, Donald G. (1981)). This projects the proposed step in the parameter vector on the gradient and estimates how much the sum of squares loss function should decrease. This estimate is divided by the current size of the loss function to avoid scale issues. When we have “converged”, the estimated decrease is very small. However, with small residuals, the sum of squares loss function is (almost) zero and we get the possibility of a zero-divide failure.

Adding a small quantity to the loss function before dividing avoids trouble. In 2021, one of us (J. Nash) proposed that `nls.control()` have an additional parameter `scaleOffset` with a default value of zero. Setting it to a small number – 1.0 is a reasonable choice – allows small-residual problems (i.e., near-exact fits) to be dealt with easily. We call this the **safeguarded relative offset convergence criterion**. The default value gives the legacy behaviour. We are pleased to note that this improvement is now in the R distributed code as of version 4.1.0.

### More general termination tests

The single convergence criterion of `nls()` leaves out some possibilities that could be useful for some problems. The package `nlsr` (John C Nash and Duncan Murdoch (2019)) already offers both the safeguarded relative offset test (`roffset`) as well as a **small sum of squares** test (`smallstest`) that compares the latest evaluated sum of squared (weighted) residuals to `e4` times the initial sum of squares, where `e4 <- (100*.Machine$double.eps)^4` is approximately 2.43e-55.

We note that `nls()` uses a termination test to stop after `maxiter` “iterations”. Unfortunately, the meaning of “iteration” varies among programs and requires careful examination of the code. We prefer to use the number of times the residuals or the jacobian have been computed and put upper limits on these. Our codes exit (terminate) when these limits are reached. Generally we prefer larger limits than the default `maxiter=50` of `nls()`, but that may reflect the more difficult problems we have encountered because users consult us when standard tools have given unsatisfactory results.

### Feature: Failure when Jacobian is computationally singular

This refers to the infamous “singular gradient” termination message of `nls()`. A Google search of

R nls "singular gradient"

gets over 4000 hits that are spread over some years. This could be because the Jacobian is poorly approximated (see **Jacobian computation** below). However, it is common in nonlinear least squares computations that the Jacobian is very close to singular for some values of the model parameters. In such cases we need to find an alternative algorithm to the Gauss-Newton iteration of `nls()`. The most common work-around is the Levenberg-Marquardt stabilization (see Marquardt (1963), Levenberg (1944), Nash (1977)), and versions of it have been implemented in packages `minpack.lm` and `nlsr`. In our work we prepared experimental prototypes of `nls` that incorporate stabilizations, but integration with all the features of `nls()` has proved difficult.

Let us consider an ‘singular gradient’ example <https://stats.stackexchange.com/questions/13053/singular-gradient-error-in-nls-with-correct-starting-values>. `nlsr::nlxb()` displays the singular values of the Jacobian, so the `nls()` error message is correct, but unhelpful in obtaining a solution.

```
reala = -3; realb = 5; realc = 0.5; realr = 0.7; realm = 1 # underlying parameters
x = 1:11 # x values; 11 data points
y=reala + realb * realr^(x - realm) + realc * x # linear+exponential function
testdat = data.frame(x, y) # save in a data frame
strt <- list(a = -3, b = 5, c = 0.5, r = 0.7, m = 1) # give programs a good start
jform <- y ~ a + b * r^(x - m) + c * x # Formula
library(nlsr)
linexp2 <- try(nlxb(jform, data = testdat, start = strt, trace = F))
linexp2 # Note singular values of Jacobian in rightmost column
```

```
## nlsr object: x
## residual sumsqares = 0 on 11 observations
## after 1 Jacobian and 1 function evaluations
## name      coeff      SE      tstat      pval      gradient      JSingval
## a          -3         NA         NA         NA         0         26.49
## b           5         NA         NA         NA         0         9.615
## c          0.5         NA         NA         NA         0         2.466
```

```
## r          0.7          NA          NA          NA          0          0.1098
## m          1          NA          NA          NA          0          1.398e-16
```

```
# We get an error with nls()
linexp=try(nls(jform, data = testdat, start = strt, trace = F))
```

```
## Error in nlsModel(formula, mf, start, wts, scaleOffset = scOff, nDcentral = nDcntr) :
##   singular gradient matrix at initial parameter estimates
```

The Jacobian is essentially singular as shown by its singular values. Note that these are **displayed** in a single column in the output to provide a compact layout, but the values do **NOT** correspond to the individual parameters in whose row they appear; they are a property of the whole problem.

## Feature: Jacobian computation

`nls()`, with the `numericDeriv()` function, computes the Jacobian as the “gradient” attribute of the residual vector. This is implemented as a mix of R and C code, but we have created a rather more compact version entirely in R.

The Jacobian used by `nlsr::nlxb()` is computed from analytic expressions, while that of `nls()` (and, by extension, `minpack.lm::nlsLM`) mostly use a numerical approximation. This can lead to different behaviour:

- A pasture regrowth problem (Huet et al. (2004), page 1, based on Ratkowsky (1983)) has a poorly conditioned Jacobian and `nls()` fails with “singular gradient”. Worse, numerical approximation to the Jacobian give the error “singular gradient matrix at initial parameter estimates” for `minpack.lm::nlsLM` so that the Marquardt stabilization is unable to take effect, while the analytic derivatives of `nlsr::nlxb` give a solution.
- Karl Schilling (private communication) provided an example where a model specified with the formula  $y \sim a * (x \wedge b)$  causes `nlsr::nlxb` to fail because the partial derivative w.r.t.  $b$  is  $a * (x^b * \log(x))$ . If there is data for which  $x = 0$ , this is undefined. In such cases, we observed that `nls()` and `minpack.lm::nlsLM` found a solution, though it can be debated whether such lucky accidents can be taken as an advantage.

Note that the selfStart models in the base R `./src/library/stats/R/zzModels.R` file provide the Jacobian in the “gradient” attribute of the “one-sided” formula that defines each model, and these Jacobians may be the analytic forms. The `nls()` function, after computing the “right hand side” or `rhs` of the residual, then checks to see if the “gradient” attribute is defined, and, if not, uses `numericDeriv` to compute a Jacobian into that attribute. This code is within the `nlsModel()` or `nlsModel.pliner()` functions. The use of analytic Jacobians almost certainly contributes to the good performance of `nls()` on selfStart models.

## Feature: Subsetting

`nls()` accepts an argument `subset`. This acts through the mediation of `model.frame` and is not clearly obvious in the source code files `/src/library/stats/R/nls.R` and `/src/library/stats/src/nls.C`.

While the implementation of `subset` at the level of the call to `nls()` has a certain attractiveness, it does mean that the programmer of the solver needs to be aware of the source (and value) of objects such as the data, residuals and Jacobian. By preference, we would implement subsetting by means of zero-value weights, with observation counts (and degrees of freedom) computed via the numbers of non-zero weights. Alternatively, we would extract a working dataframe from the relevant elements in the original.

## Feature: na.action

`na.action` is an argument to the `nls()` function, but it does not appear obviously in the source code, often being handled behind the scenes after referencing the option `na.action`. A useful, but possibly dated, description is given in: <https://stats.idre.ucla.edu/r/faq/how-does-r-handle-missing-values/>.

The typical default action, which can be seen by using the command `getOption("na.action")` is `na.omit`. This option essentially omits from computations any observations containing missing values (i.e. any row of a data frame containing an NA). `na.exclude` does much the same for computations, but keeps the rows with NA elements so that predictions are in the correct row position. We recommend that workers actually test output to verify behaviour is as wanted.

A succinct description of the issue is given in: <https://stats.stackexchange.com/questions/492955/should-i-use-na-omit-or-na-exclude-in-a-linear-model-in-r>

The only benefit of `na.exclude` over `na.omit` is that the former will retain the original number of rows in the data. This may be useful where you need to retain the original size of the dataset - for example it is useful when you want to compare predicted values to original values. With `na.omit` you will end up with fewer rows so you won't as easily be able to compare.

`na.pass` simply passes on data "as is", while `na.fail` will essentially stop if any missing values are present.

Our concern with `na.action` is that users may be unaware of the effects of an option setting they may not even be aware has been set. Is `na.fail` a safer default?

## Feature: model frame

`model` is an argument to the `nls()` function, which is documented as

**model** logical. If true, the model frame is returned as part of the object. Default is FALSE.

Indeed, the argument only gets used when `nls()` is about to return its result object, and the element `model` is NULL unless the calling argument `model` is TRUE. (Using the same name could be confusing.) However, the model frame is used within the function code in the form of the object `mf`. We feel that users could benefit from more extensive documentation and examples of its use.

## Feature: sources of data

`nls()` can be called without specifying the `data` argument. In this case, it will search in the available environments (i.e., workspaces) for suitable data objects. We do NOT like this approach, but it is "the R way". R allows users to leave many objects in the default (`.GlobalEnv`) workspace. Moreover, users have to actively suppress saving this workspace (`.RData`) on exit, and any such file in the path when R is launched will be loaded. The overwhelming proportion of R users in our acquaintance avoid saving the workspace because of the danger of lurking data and functions which may cause unwanted results.

Nevertheless, to provide compatible behaviour with `nls()`, competing programs need to ensure equivalent behaviour, but users should test that the operation is as they intend.

## Feature: missing start vector and self-starting models

Nonlinear estimation algorithms are almost all iterative and need a set of starting parameters. `nls()` offers a special class of modeling functions called **selfStart** models. There are a number of these in base R (`./src/library/stats/R/zzModels.R`) and others in R packages such as CRAN package `nlraa` (Miguez (2021)), as well as the now-archived package `NRAIA`. Unfortunately, some **selfStart** codes entangle the calculation of starting values for parameters with the particulars of the `nls()` code. Though there is a `getInitial()` function, this is not easy to use to simply compute the initial parameter estimates, in part because it may call `nls()`. Such circular references are, in our view, dangerous.

Other concerns are revealed by the example below. Here the `SSlogis` selfStart function is used to generate a set of initial parameters for a 3-parameter logistic curve. The form used by `SSlogis` is  $y \sim Asym/(1 + \exp((xmid - tt)/scal))$ , but we show how the starting parameters for this model can be transformed to those of another form of the model, namely,  $y \sim b1/(1 + b2 * \exp(-b3 * t))$ .

The code for `SSlogis()` is in `./src/library/stats/R/zzModels.R`. This R function includes analytic expressions for the Jacobian ("gradient"). These could be useful to R users, especially if documented.

Moreover, we wonder why the programmers have chosen to save so many quantities in “hidden” variables, i.e., with names preceded by “.”. These are then not displayed by the `ls()` command, making them difficult to access.

In the event that a `selfStart` model is not available, `nls()` sets all the starting parameters to 1. This is, in our view, tolerable, but could be improved by using a set of values that are all slightly different, which, in the case of the example model  $y \sim a * \exp(-b * x) + c * \exp(-d * x)$  would avoid a singular Jacobian because  $b$  and  $d$  were equal in value. A sequence like 1.0, 1.1, 1.2, 1.3 for the four parameters could be provided quite simply instead of all 1’s.

### Strategic issues in selfStart models

Because the Gauss-Newton algorithm is unreliable from many sets of starting of parameters, `selfStart` models are a part of the `nls()` infrastructure rather than an accessory. However, creating such functions is a lot of work, and their documentation (file `./src/library/stats/man/selfStart.Rd`) is quite complicated. We believe that the focus should be placed on getting good initial parameters, that is `getInitial()` function, though avoiding the current calls back to `nls()`. Interactive tools, such as “visual fitting” (Nash and Velleman (1996)) might be worth considering.

### Issue: returned output of nls() and its documentation

The output of `nls()` is an object of class “nls” which has the following structure:

A list of:

<code>m</code>	an <code>nlsModel</code> object incorporating the model.
<code>data</code>	the expression that was passed to <code>nls</code> as the data argument. The actual data values are present in the environment of the <code>m</code> components, e.g., <code>environment(m\$conv)</code> .
<code>call</code>	the matched call with several components, notably <code>algorithm</code> .
<code>na.action</code>	the “na.action” attribute (if any) of the model frame.
<code>dataClasses</code>	the “dataClasses” attribute (if any) of the “terms” attribute of the model frame.
<code>model</code>	if <code>model = TRUE</code> , the model frame.
<code>weights</code>	if <code>weights</code> is supplied, the weights.
<code>convInfo</code>	a list with convergence information.
<code>control</code>	the control list used, see the <code>control</code> argument.

There are also two deprecated items if `algorithm = "port"` fit only. These are `convergence` (a code = 0 for convergence) and `message`. These are available from `convInfo`.

The `nls` object contains some elements that are awkward to produce by other algorithms, but some information that would be useful is not presented obviously. Moreover, the complexity of the object is a challenge to users.

In the following, we use `result` as the returned object from `nls()`.

The `data` return element is an R symbol. To actually access the data from this element, we need to use the syntax

```
eval(parse(text=result$data))
```

However, if the call is made with `model=TRUE`, then there is a returned element `model` which contains the data, and we can list its contents using

```
ls(result$model)
```

and if there is an element called `xdata`, it can be accessed as `result$model$xdata`.

Let us compare the `nls()` result with that from `nlsr::nlxb()`, with ostensibly solves the same problem:

<code>coefficients</code>	A named vector giving the parameter values at the supposed solution.
<code>ssquares</code>	The sum of squared residuals at this set of parameters.

<code>resid</code>	The residual vector at the returned parameters.
<code>jacobian</code>	The jacobian matrix (partial derivatives of residuals w.r.t. the parameters) at the returned parameters.
<code>feval</code>	The number of residual evaluations (sum of squares computations) used.
<code>jeval</code>	The number of Jacobian evaluations used.

However, actually looking at the structure of a returned result gives a list of 11 items. The extra 5 are:

```
$ lower      : num [1:3] -Inf -Inf -Inf
$ upper      : num [1:3] Inf Inf Inf
$ maskidx    : int(0)
$ weights    : NULL
$ formula    : Class 'formula' language y ~ Asym/(1 + exp((xmid - tt)/scal))
.. ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
- attr(*, "class")= chr "nlshr"
```

The result object from `nlshr::nlxb()` is still much smaller than the one `nls()` returns. Moreover, `nlxb` explicitly returns the sum of squares as well as the residual vector and Jacobian. The counts of evaluations are also returned. (Note that the singular values of the Jacobian are actually computed in the `print` and `summary` methods for the result.) We noted several potential updates to the `nlshr` documentation as well as that for `nls()` as we made our comparisons.

### Weights in returned functions from `nls()`

The functions `resid()` (an alias for `residuals()`) and `fitted()` and `lhs()` are UNWEIGHTED. But if we return `ans` from `nls()` or `minpack.lm::nlsLM` or our new `nlsh` (interim package), then `ans$m$resid()` is WEIGHTED. This needs clear documentation so users can get the answers intended.

### Interim output from the “port” algorithm

As the `nls()` **man** page states, when the “port” algorithm is used with the `trace` argument `TRUE`, the iterations display the objective function value which is 1/2 the sum of squares (or deviance). It is likely that the trace display is embedded in the Fortran of the `nlminb` routine that is called to execute the “port” algorithm, but the discrepancy is nonetheless unfortunate for users.

### Failure to return best result achieved

If `nls()` reaches a point where it cannot continue but has not found a point where the relative offset convergence criterion is met, it may simply exit, especially if a “singular gradient” (singular Jacobian) is found. However, this may occur AFTER the function has made considerable progress in reducing the sum of squared residuals. Here is an abbreviated example:

```
time=c( 1,  2,  3,  4,  6 , 8, 10, 12, 16)
conc = c( 0.7, 1.2, 1.4, 1.4, 1.1, 0.8, 0.6, 0.5, 0.3)
NLSdata <- data.frame(time,conc)
NLSstart <-c(lrc1=-2,lrc2=0.25,A1=150,A2=50) # a starting vector (named!)
NLSformula <-conc ~ A1*exp(-exp(lrc1)*time)+A2*exp(-exp(lrc2)*time)
tryit <- try(nls(NLSformula, data=NLSdata, start=NLSstart, trace=TRUE))
```

```
## 61216.    (3.56e+03): par = (-2 0.25 150 50)
## 2.1757    (2.23e+01): par = (-1.9991 0.31711 2.6182 -1.3668)
## 1.6211    (7.14e+00): par = (-1.9605 -2.6203 2.5753 -0.55599)
## Error in nls(NLSformula, data = NLSdata, start = NLSstart, trace = TRUE) :
##   singular gradient
print(tryit)
```

```
## [1] "Error in nls(NLSformula, data = NLSdata, start = NLSstart, trace = TRUE) : \n singular gradien"
```



```
## attr("class")
## [1] "try-error"
## attr("condition")
## <simpleError in nls(NLSformula, data = NLSdata, start = NLSstart, trace = TRUE): singular gradient>
```

Note that the sum of squares has been reduced from 61216 to 1.6211, but unless `trace` is invoked, the user will not get any information about this. This almost trivial change to the `nls()` function could be useful to R users.

## Feature: partially linear models and their specification

Specifying a model to a solver should, ideally, use the same syntax across solver tools. Unfortunately, R allows multiple approaches within different modelling tools, and within `nls()` itself.

One obvious case is that nonlinear modeling tools are a superset of linear ones. Yet the explicit model  $y \sim a \cdot x + b$  does not work with the linear modeling function `lm()`, which requires this model to be specified as  $y \sim x$ .

Within `nls()`, consider the following FOUR different calling sequences for the same problem, though an intuitive choice, labelled `fm2a`, will not work. In this failed attempt, putting the `Asym` parameter in the model causes the `plinear` algorithm to try to add another term to the model. We believe this is unfortunate, and would like to see a consistent syntax. At the time of writing we do not foresee a resolution for this issue. In the example we have NOT evaluated the commands to save space.

```
DNase1 <- subset(DNase, Run == 1) # select the data
## using a selfStart model - do not specify the starting parameters
fm1 <- nls(density ~ SSlogis(log(conc), Asym, xmid, scal), DNase1)
summary(fm1)

## using conditional linearity - leave out the Asym parameter
fm2 <- nls(density ~ 1/(1 + exp((xmid - log(conc))/scal)),
          data = DNase1, start = list(xmid = 0, scal = 1),
          algorithm = "plinear")
summary(fm2)

## without conditional linearity
fm3 <- nls(density ~ Asym/(1 + exp((xmid - log(conc))/scal)),
          data = DNase1,
          start = list(Asym = 3, xmid = 0, scal = 1))
summary(fm3)

## using Port's nl2sol algorithm
fm4 <- try(nls(density ~ Asym/(1 + exp((xmid - log(conc))/scal)),
             data = DNase1, start = list(Asym = 3, xmid = 0, scal = 1),
             algorithm = "port"))
summary(fm4)

## using conditional linearity AND Asym does NOT work
fm2a <- try(nls(density ~ Asym/(1 + exp((xmid - log(conc))/scal)),
              data = DNase1, start = list(Asym=3, xmid = 0, scal = 1),
              algorithm = "plinear", trace = TRUE))
summary(fm2a)
```



## Issue: code structure and documentation for maintenance

The `nls()` code is structured in a way that inhibits both maintenance and improvement. In particular, the iterative setup is such that introduction of Marquardt stabilization is not easily available. To obtain performance a lot of the code was in C with consequent calls and returns that complicate the code. Over time, R has become much more efficient on modern computers, and the need to use compiled C and Fortran is less critical. Moreover, the burden for maintenance could be reduced by moving code entirely to R.

While R and its packages are generally very well documented for usage, that documentation rarely extends as far as it might in directions helpful for code maintenance. The paucity of such documentation is exacerbated by the mixed R/C/Fortran code base, where seemingly trivial differences in things like representation of numbers or base index for arrays lead to traps for unwary programmers.

For `nls()` the concern is demonstrated by a short email query from John Nash to Doug Bates and his reply. This is NOT a criticism of Prof. Bates (or any other) work, but a reflection on how difficult it is to develop code in this subject area and to keep it maintainable. We have experienced similar loss of understanding for some of our own codes.

```
...
https://gitlab.com/nashjc/improvenls/-/blob/master/Croucher-expandednlsnoc.R
... has the test problem and the expanded code. Around line 367 is where we are
scratching our heads. The function code (from nlsModel()) is in the commented
lines below the call. This is
```

```
# > setPars
# function(newPars) {
#   setPars(newPars)
#   resid <- .swts * (lhs - (rhs <- getRHS())) # envir = thisEnv {2 x}
#   dev   <- sum(resid^2) # envir = thisEnv
#   if(length(gr <- attr(rhs, "gradient")) == 1L) gr <- c(gr)
#   QR <- qr(.swts * gr) # envir = thisEnv
#   (QR$rank < min(dim(QR$qr))) # to catch the singular gradient matrix
# }
```

I'm anticipating that we will be able to set up a (possibly inefficient) code with documentation that will be easier to follow and test, then gradually figure out how to make it more efficient.

The equivalent from `minpack.lm` is

```
setPars = function(newPars) {
  setPars(newPars)
  assign("resid", .swts * (lhs - assign("rhs", getRHS(),
    envir = thisEnv)), envir = thisEnv)
  assign("dev", sum(resid^2), envir = thisEnv)
  assign("QR", qr(.swts * attr(rhs, "gradient")), envir = thisEnv)
  return(QR$rank < min(dim(QR$qr)))
}
```

In both there is the recursive call, which must have a purpose I don't understand.

The reply was

```
I'm afraid that I don't know the purpose of the recursive call either.
I know that I wrote the code to use a closure for the response, covariates, etc.,
but I don't recall anything like a recursive call being necessary.
...
```

## Feature: indexed parameters

The **man** file for `nls()` includes the following example of a situation in which parameters are indexed. It also uses the “plinear” option as an added complication. Here we use a truncated version of the example to save display space.

```
## The muscle dataset in MASS is from an experiment on muscle
## contraction on 21 animals. The observed variables are Strip
## (identifier of muscle), Conc (Cacl concentration) and Length
## (resulting length of muscle section).
if(! requireNamespace("MASS", quietly = TRUE)) stop("Need MASS pkg")
mm<- MASS::muscle[1:12,] # take only 1st few values of Strip
str(mm)

## 'data.frame': 12 obs. of 3 variables:
## $ Strip : Factor w/ 21 levels "S01","S02","S03",...: 1 1 1 1 2 2 2 2 3 3 ...
## $ Conc : num 1 2 3 4 1 2 3 4 0.25 0.5 ...
## $ Length: num 15.8 20.8 22.6 23.8 20.6 26.8 28.4 27 7.2 15.4 ...

mm<-droplevels(mm)
str(mm)

## 'data.frame': 12 obs. of 3 variables:
## $ Strip : Factor w/ 3 levels "S01","S02","S03": 1 1 1 1 2 2 2 2 3 3 ...
## $ Conc : num 1 2 3 4 1 2 3 4 0.25 0.5 ...
## $ Length: num 15.8 20.8 22.6 23.8 20.6 26.8 28.4 27 7.2 15.4 ...

nlev <- nlevels(mm)
withAutoprint({
  ## The non linear model considered is
  ## Length = alpha + beta*exp(-Conc/theta) + error
  ## where theta is constant but alpha and beta may vary with Strip.
  with(mm, table(Strip)) # 2, 3 or 4 obs per strip
  nl <- nlevels(mm$Strip)
  ## We first use the plinear algorithm to fit an overall model,
  ## ignoring that alpha and beta might vary with Strip.
  musc.1 <- nls(Length ~ cbind(1, exp(-Conc/th)), mm,
    start = list(th = 1), algorithm = "plinear")
  summary(musc.1)

  ## Then we use nls' indexing feature for parameters in non-linear
  ## models to use the conventional algorithm to fit a model in which
  ## alpha and beta vary with Strip. The starting values are provided
  ## by the previously fitted model.
  ## Note that with indexed parameters, the starting values must be
  ## given in a list (with names):
  ## ?? but why use b here AND in the new formula??
  b <- coef(musc.1)

  musc.2 <- nls(Length ~ a[Strip] + b[Strip]*exp(-Conc/th), data=mm,
    start = list(a = rep(b[2], nl), b = rep(b[3], nl), th = b[1]))
  summary(musc.2)
})

## > with(mm, table(Strip))
## Strip
## S01 S02 S03
```

```

## 4 4 4
## > nl <- nlevels(mm$Strip)
## > musc.1 <- nls(Length ~ cbind(1, exp(-Conc/th)), mm, start = list(th = 1),
## +      algorithm = "plinear")
## > summary(musc.1)
##
## Formula: Length ~ cbind(1, exp(-Conc/th))
##
## Parameters:
##      Estimate Std. Error t value Pr(>|t|)
## th      0.624      0.222   2.81  0.0203 *
## .lin1  25.684      1.441  17.82 2.5e-08 ***
## .lin2 -26.631      6.147  -4.33  0.0019 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.96 on 9 degrees of freedom
##
## Number of iterations to convergence: 5
## Achieved convergence tolerance: 4.71e-07
##
## > b <- coef(musc.1)
## > musc.2 <- nls(Length ~ a[Strip] + b[Strip] * exp(-Conc/th), data = mm,
## +      start = list(a = rep(b[2], nl), b = rep(b[3], nl), th = b[1]))
## > summary(musc.2)
##
## Formula: Length ~ a[Strip] + b[Strip] * exp(-Conc/th)
##
## Parameters:
##      Estimate Std. Error t value Pr(>|t|)
## a1  22.9277      0.5154  44.49 1.1e-07 ***
## a2  27.8606      0.5150  54.09 4.1e-08 ***
## a3  28.3426      1.0771  26.31 1.5e-06 ***
## b1 -44.1558     12.3900  -3.56  0.016 *
## b2 -43.8084     12.3148  -3.56  0.016 *
## b3 -32.8421      2.1553 -15.24 2.2e-05 ***
## th   0.5548      0.0808   6.86  0.001 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.798 on 5 degrees of freedom
##
## Number of iterations to convergence: 6
## Achieved convergence tolerance: 1.94e-06

```

Note that the answers for the parameters are NOT indexed. e.g. `coef(musc.2)` is a single level vector of parameters. We do not see `a[1]`, `a[2]`, `a[3]` but `a1`, `a2`, `a3`. This is because the model must integrate all the parameters because `th` is a common parameter across the index `Strip`.

We believe this structure is quite likely to cause confusion and error, and propose an alternative approach below.

## Feature: bounds on parameters

There are many situations where the context of a problem constrains the values of parameters. For example, one of us was asked many years ago to estimate a model. One parameter returned a negative value. The client complained “But that is supposed to be the number of grain elevators in Saskatchewan”. That is, the number should have been a positive integer, and likely less than a few thousand.

`nls()` can impose bounds on parameters, but only if the `port` algorithm is used. Unfortunately, the manual states

The algorithm = “port” code appears unfinished, and does not even check that the starting value is within the bounds. Use with caution, especially where bounds are supplied.

This is clearly unsatisfactory for general use. Inclusion of bounds on parameters in nonlinear least squares computations is relatively straightforward, and it is part of the default method for package `nlsr`. Package `minpack.lm` also includes provision for bounds, but the following script, which we do not run here for reasons of space, shows that the solutions found are not optimal for the Hobbs test problem, even in its scaled form.

```
#### bounds with formula specification of problem
## Use the Hobbs Weed problem
weed <- c(5.308, 7.24, 9.638, 12.866, 17.069, 23.192, 31.443,
          38.558, 50.156, 62.948, 75.995, 91.972)
tt <- 1:12
weedframe <- data.frame(y=weed, tt=tt)
st <- c(b1=1, b2=1, b3=1) # a default starting vector (named!)
## Unscaled model
wmodu <- y ~ b1/(1+b2*exp(-b3*tt))
## Scaled model
wmods <- y ~ 100*b1/(1+10*b2*exp(-0.1*b3*tt))
## We can provide the residual and Jacobian as functions
# Scaled Hobbs problem
shobbs.res <- function(x){ # scaled Hobbs weeds problem -- residual
  # This variant uses looping
  if(length(x) != 3) stop("shobbs.res -- parameter vector n!=3")
  y <- c(5.308, 7.24, 9.638, 12.866, 17.069, 23.192, 31.443,
          38.558, 50.156, 62.948, 75.995, 91.972)
  tt <- 1:12
  res <- 100.0*x[1]/(1+x[2]*10.*exp(-0.1*x[3]*tt)) - y
}

shobbs.jac <- function(x) { # scaled Hobbs weeds problem -- Jacobian
  jj <- matrix(0.0, 12, 3)
  tt <- 1:12
  yy <- exp(-0.1*x[3]*tt)
  zz <- 100.0/(1+10.*x[2]*yy)
  jj[tt,1] <- zz
  jj[tt,2] <- -0.1*x[1]*zz*zz*yy
  jj[tt,3] <- 0.01*x[1]*zz*zz*yy*x[2]*tt
  attr(jj, "gradient") <- jj
  jj
}

##----- Hobbsbounded unique code starts here -----
require(minpack.lm)
traceval<-FALSE
cat("Infeasible start test\n")
```

```

start1inf <- c(b1=4, b2=4, b3=4) # b3 OUT OF BOUNDS for next few tries
## Infeasible start! No warning message!
anlM2i <- try(nlsLM(wmodu, start=start1inf, data=weedframe, lower=c(0,0,0), upper=c(2, 6, 3),
                  trace=traceval))

# feasible start i.e. on or within bounds
start1<-st
anlM1 <- try(nlsLM(wmodu, start=start1, data=weedframe, lower=c(0,0,0), upper=c(2, 6, 3),
                  trace=traceval))
print(anlM1)

# Hobbs scaled problem with bounds, formula specification

require(nlsr)
anlxb1 <- nlxb(wmods, start=start1, data=weedframe, lower=c(0,0,0),
               upper=c(2,6,3), trace=traceval)
print(anlxb1)

## nlsLM seems NOT to work properly with bounds
anlsLM1b <- nlsLM(wmods, start=start1, data=weedframe, lower=c(0,0,0),
                  upper=c(2,6,3), trace=traceval)
print(anlsLM1b)
newst<-coef(anlsLM1b)
print(newst)
anlsLM1bx <- nlsLM(wmods, start=newst, data=weedframe, lower=c(0,0,0),
                  upper=c(2,6,3))
print(anlsLM1bx)

cat("Single number bounds:\n")
anlsLM1b1 <- try(nlsLM(wmods, start=start1, data=weedframe, lower=0,
                      upper=3))
print(anlsLM1b1)
anlsLM1b1a <- try(nlsLM(wmods, start=start1, data=weedframe, lower=c(0,0,0),
                       upper=c(3,3,3)))
print(anlsLM1b1a)

anlxb1b1<- try(nlxb(wmods, start=start1, data=weedframe, lower=0,
                   upper=3))
print(anlxb1b1)

anlm1b <- nls.lm(par=start1, fn=shobbs.res, jac=shobbs.jac, lower=c(0,0,0),
                 upper=c(2,6,3))
print(anlm1b)

# functional presentation of problem -- seems to work
anlm1bx <- nls.lm(par=start1, fn=shobbs.res, jac=shobbs.jac, lower=c(0,0,0),
                  upper=c(2,6,3))
print(anlm1bx)

anlfb1bx <- nlfb(start=start1, resfn=shobbs.res, jacfn=shobbs.jac, lower=c(0,0,0),
                 upper=c(2,6,3))
print(anlfb1bx)

```

## Philosophical considerations

Bounds on parameters raise some interesting and difficult questions about how uncertainty in parameter estimates should be computed or reported. That is, the traditional “standard errors” are generally taken to imply symmetric intervals about the point estimate in which the parameter may be expected to be found with some probability under certain assumptions. Bounds change those assumptions drastically. Hence `nlsr::nlxb()` does not compute standard errors nor their derived statistics when bounds are active.

## Goals of our effort to improve `nls()`

### Code rationalization and documentation

We want

- to provide a packaged version of `nls()` (call it `nlsalt`) coded entirely in R that matches the version in base R or what is packaged in `nls pkg` as described in Nash and Bhattacharjee (2022).
- to streamline the overall `nls()` infrastructure. By this we mean a re-factoring of the routines so they are better suited to maintenance of both the existing `nls()` methods and features as well as new features we or others would like to add.
- to explain what we do, either in comments or separate maintainer documentation. Since we are complaining about the lack of explanatory material for the current code, we feel it is incumbent on us to provide such material for our own work, and if possible for the existing code.

### Rationalization of formula specifications

We have seen that `nls()` uses a different formula specification from the default if the `plinear` algorithm is used. This is unfortunate, since the user cannot then simply add `algorithm="plinear"` to the call. This makes errors more likely.

While not consistent with the `lm()` function and specification of linear models, we believe a complete set of parameters should be specified for nonlinear models. That is, the model formula should be written down as it would be computed, rather than for estimation. Thus  $y \sim a \cdot x + b$  rather than  $y \sim x$ .

Moreover, we want to avoid the partially linear parameters appearing in the result object as `.lin1`, `.lin2`, etc., where their actual position in the model is not explicit.

A possible workaround to allow consistency would be to specify the partially linear parameters when the `plinear` option is specified. For example, we could use `algorithm="plinear(Asym)"` while keeping the general model formula from the default specification. This would allow for the output of different `algorithm` options to be consistent. However, we have not yet tried to code such a change.

A further complication of model specification in R is that some formulas require the user to surround a term with `I( )` to inhibit interpretation of arithmetic operators. We would prefer that formulas be usable without this complication, which has limited documentation and examples.

It is also of interest to consider code that will detect partial linearity, rather than expecting the user to provide the special structure and settings. Work in this area appears to be quite limited, with “Automatic Linearity Detection, Report NA-13-04” (2013) seemingly a sole contribution we can find. Unfortunately, the mention of a “Linear and Nonlinear Discoverer” in Zhang, Cheng, and Liu (2011) is about a class of models that have almost nothing in common with the present issue. The algorithm in “Automatic Linearity Detection, Report NA-13-04” (2013) will need investigating further, but we suspect that it is complicated enough that it will not easily be adapted to a program like `nls()`. Possibly a simpler algorithm may be found.

## Rationalization of indexed models

Indexed models clearly have a place in some areas of research. We need to be able to process models, as the example above, where the model formula is `Length ~ a[Strip] + b[Strip] * exp(-Conc/th)` and `Strip` is the index. Data `Length` and `Conc` are available for all observations. Parameter `th` applies to every fitted point, but `a[Strip]` and `b[Strip]` depend on the index, which could be an experimental run, or some other categorization.

Given that `[ ]` are the R method to index arrays, their presence in a model formula signals indexing. The question is how to perform the correct calculations.

While tedious (??maybe dummies or fastdummies packages make it less so?), we could set up the model as a sum of similar models with dummy variables to turn “on” or “off” the particular case, transforming the model parameters to simple, non-indexed ones (as does `nls()`).

Wrapper functions to prepare the estimation calculations and present the results in a form that matches the modelling formula would avoid adding complexity to `nls()` and other tools.

## Provide tests and use-case examples

We need suitable tests and use-case examples in order:

- to ensure our new `nlsalt` or related packages work properly, in particular, giving results comparable to or better than the `nls()` in base R or `nls pkg`;
- to test individual solver functions to ensure they work across the range of calling mechanisms, that is, different ways of supplying inputs to the solver(s);
- to pose “silly” inputs to nonlinear least squares solvers (in R) to see if these bad input exceptions are caught by the programs.

## A test runner program ?? NLSCompare

?? Arkajyoti – do you want to expand?

When we have a “new” or trial solver function, we would like to know if it gives acceptable results on a range of sample problems of different types, starting parameters, input conditions, constraints, subsets, weights or other settings. Ideally we want to be able to get a summary that is easy to read and assess. For example, one approach would be to list the names of a set of tests with a red, green or yellow dot beside the name for FAILURE, SUCCESS, or “NOT APPLICABLE”. In the last category would be a problem with constraints that the solver is not designed to handle.

To accomplish this, we need a suitable “runner” program that can be supplied with the name of a solver or solvers and a list of test problem cases. Problems generally have a base setup – a specification of the function to fit as a formula, some data and a default starting set of parameters. Other cases can be created by imposing bounds or mask constraints, subsets of the data, and different starts.

How to set up this “runner” and its supporting infrastructure is non-trivial. While the pieces are not as complicated as the inter-related parts of the solvers, especially `nls()`, the categorization of tests, their documentation, and the structuring to make running them straightforward and easy requires much attention to detail.

Some considerations for our test scripts:

- Is it useful to have a “base” script for each family of test problem, with numbered particular cases? That is, if we run the scripts in order, we can avoid some duplication of code and data.
- While we have developed some tags to document the test problem families and cases, we believe that such tags (essentially summary documentation) will continue to need revision as different tools and problems are included in scope of `nlsCompare`.
- Similarly, we expect that there will be ongoing review of the structure of the result files.



## Progress so far

Much of our work is reported in <https://github.com/nashjc/RNonlinearLS> .

## Reports and articles

- The present article is the central report of the Google Summer of Code 2021 project “Improvements to `nls()`”.
- `VarietyInNonlinearLeastSquaresCodes.Rmd`: a review of the different algorithms and the many choices in their implementation for nonlinear least squares. This is still a DRAFT at 2021-8-20. ?? fix
- `TestsDoc.Rmd` (??citation) is a survey of testing tools in R. It has more general possibilities and fits into the subject of regression testing, in which case a more extensive literature review will be needed. Note that this document reflects the work in the the “Problem sets and test infrastructure” below.
- `DerivsNLS.Rmd` ??citation: a document to explain different ways in which Jacobian information is supplied to nonlinear least squares computation in R. File `ExDerivs.R` is a DRAFT of a script to provide examples.
- The BibTex bibliography `ImproveNLS.bib` (??citation or reference) was set up for use with all documents in this project, but has wider application to nonlinear least squares projects in general.
- `MachineSummary.Rmd` (??citation – blog post??) is an investigation of ways to report the characteristics and identity of machines running tests. `MachID.R` offers a suggested concise summary function to identify a particular computational system used for tests.
- Nash and Bhattacharjee (2022) is an explanation of the construction of the `nls pkg` from the code in R-base.
- As the 2021 Summer of Code period was ending, one of us (JN) was invited to prepare a review of optimization in R. Ideas from the present work have been instrumental in the creation of Nash (2022).

## Problem sets and test infrastructure

?? need to clean up

We have several test problems and variants thereof in the `inst/scripts/` directory of the `nlsCompare` package available on Github (<https://github.com/ArkaB-DS/nlsCompare>). We direct the reader to that package for documentation of the test infrastructure, in particular the problems and methods files (`problems.csv` and `methods.csv`) and the various functions invoked by `run.R` to produce an output file in CSV form also.

Towards the end of the project, we have focused our attention on the `nlsCompare` package, which looks at evaluating and comparing functions for nonlinear least squares problems. We use many of the same tests as checks that are or will be built into our packages for such problems e.g., `nlsalt`. These use the `testthat` structure and may include verification of outputs that are specific to a package, such as the upper triangular matrix `R` of the QR decomposition that has been computed for a Jacobian. Since such an object will not be computed by all methods, testing it in `nlsCompare` makes no sense, and in that package we concentrate on the minimum sum of squares and the associated model parameters.

## Code and documentation

A summary of our main results:

- `nls pkg`: a packaged version of the ‘`nls()`’ code from R-base. Thanks to Duncan Murdoch
- `nlsalt`: attempt to mirror `nls()` behaviour entirely in R. This is UNFINISHED. The effort showed that the structure of the programs was difficult to follow, undocumented, and unsuited to adding improvements like the Marquardt stabilization. We were able to get a pure-R version of `numericDeriv()` and rework most of the functions of `nlsModel` (but

`notnlsModel.plinear`). This work may continue after the project formally ends, but collaboration is likely needed with workers who have a deep knowledge of R internals. `nls-flowchart.txt` was a start at documentation of the structure of `nls()`.

- the safeguarded relative offset convergence criterion that allows small-residual problems to be handled with `nls()` is now in the R base code.
- `nlsj` is a refactoring of some `nls()` functionality. ?? What are key features?? This is an interim package for exploration of ideas and will NOT, as it currently is presented, become a distributed package.
- `nlsralt`: this is intended to be a modified version of Nash and Murdoch package `nlsr` to examine possible improvements discovered as a result of this project. This is EXPERIMENTAL.

## Strategic choices in a nonlinear model estimation

The key difference in approach between `nls()` and `nlsr::nlxb()` is that `nls()` builds a large infrastructure from which the Gauss-Newton iteration can be executed and other statistical information such as profiles can be computed, while `nlxb()` returns quite limited information, and in its execution computes what is needed on an “as and when” basis. This follows a path that one of us (JN) established almost 50 years ago with the software that became Nash (1979), using setup, solver, and post-solution analysis phases to computation. Here, that last phase is not part of `nlxb()`, but is the domain of other functions in the `nlsr` package. `nls()` also has some similar functions, but they are much more tied into the infrastructure that is created, mostly in the `nlsModel()` function.

The `nls()` approach, as implemented in the base R code and `nlspkg` leads to considerable entangling of the different functions and capabilities. Even after both the preparation for and carrying out of the present project, we do not feel confident to explain the code completely, nor to maintain it. However, we have made some forays into render parts of the code in R, and to understand where difficulties lie.

## Opinion for discussion

To advance the stability and maintainability of R, we believe the program objects (R functions) that are created by tools such as `nls()` should have minimal interactions and side-effects. The aspects of `nls()` that have given us the most trouble are as follows:

- The functions that compute the residuals and jacobians are frequently computed by presuming the current data and parameters for those functions are available in an environment. As long as the correct environment is used, this provides a surprisingly short syntax to invoke the calculations. The danger is that the wrong data is accessed if the internal search finds a valid name that is not the object we want.
- Weights, subsets, and various contextual controls such as that for the `na.action` that tell our code what to do with undefined or missing numbers. Again, this makes for a very simple invocation of the calculation, but the context is hidden from the user. It can be difficult for those of us trying to maintain or improve the code to be certain we have the context correct.
- The mixing of R and C code made for high performance in the past, but without developer documentation, programmers face a lot of work to fix code. We believe in keeping code – at least a runnable reference version – in a single programming language. If necessary, by measuring (“profiling”) the code, we can find bottlenecks and replace just those slower parts of the reference code.
- We believe that a structure that isolates the setup, solve, and post-solution structures for complicated functions reduces the number of objects that must be kept in appropriate alignment at any one stage in a set of calculations.
- We feel that all codes should return the best fit they have found so far, even if there are untoward conditions reached, such as a singular Jacobian. This modification could be made to `nls()` even in the short term.

??fix next bits

## Acknowledgement

Hans Werner Borchers has been helpful in developing the proposal for this project and in comments on this and related work. Heather Turner co-mentored the project and helped guide the progress of the work.??

## References

- “Automatic Linearity Detection, Report NA-13-04.” 2013. Mathematical Institute, University of Oxford. <https://core.ac.uk/download/pdf/9695588.pdf>.
- Bates, D. M., and D. G. Watts. 1988. *Nonlinear Regression Analysis and Its Applications*. Wiley.
- Bates, Douglas M., and Watts, Donald G. 1981. “A Relative Offset Orthogonality Convergence Criterion for Nonlinear Least Squares.” *Technometrics* 23 (2): 179–83.
- Huet, S., A. Bouvier, M.-A. Poursat, and E. Jolivet. 2004. *Statistical Tools for Nonlinear Regression: A Practical Guide with S-PLUS Examples, 2nd Edition*. Berlin & New York: Springer-Verlag.
- John C Nash, and Duncan Murdoch. 2019. *nlsr: Functions for Nonlinear Least Squares Solutions*.
- Levenberg, Kenneth. 1944. “A Method for the Solution of Certain Non-Linear Problems in Least Squares.” *Quarterly of Applied Mathematics* 2: 164–168.
- Marquardt, Donald W. 1963. “An Algorithm for Least-Squares Estimation of Nonlinear Parameters.” *SIAM Journal on Applied Mathematics* 11 (2): 431–41.
- Miguez, Fernando. 2021. *nlraa: Nonlinear Regression for Agricultural Applications*. <https://CRAN.R-project.org/package=nlraa>.
- Nash, John C. 1977. “Minimizing a Nonlinear Sum of Squares Function on a Small Computer.” *Journal of the Institute for Mathematics and Its Applications* 19: 231–37.
- . 1979. *Compact Numerical Methods for Computers: Linear Algebra and Function Minimisation*. Bristol: Adam Hilger.
- . 2022. “Function Minimization and Nonlinear Least Squares in r.” *WIREs Computational Statistics*, no. e1580. <https://doi.org/https://doi.org/10.1002/wics.1580>.
- Nash, John C., and Arkajyoti Bhattacharjee. 2022. “Making a Package from Base R Files.” *R-Bloggers*, January. <https://www.r-bloggers.com/2022/01/making-a-package-from-base-r-files/>.
- Nash, John C., and Paul Velleman. 1996. “Nonlinear Estimation Combining Visual Fitting with Optimization Methods.” In *Proceedings of the Section on Physical and Engineering Sciences of the American Statistical Association*, 256–61. American Statistical Association.
- Ratkowsky, David A. 1983. *Nonlinear Regression Modeling: A Unified Practical Approach*. New York; Basel: Marcel Dekker Inc.
- Zhang, Hao Helen, Guang Cheng, and Yufeng Liu. 2011. “Linear or Nonlinear? Automatic Structure Discovery for Partially Linear Models.” *Journal of the American Statistical Association* 106 (495): 1099–1112. <http://www.jstor.org/stable/23427577>.