

Guidelines for S3 Regression Models

Stephen Milborrow

November 24, 2016

Abstract

This document is intended for authors of R functions that build S3 regression models. It describes how these functions should interface to the rest of the world. The intention is to summarize current good practice. Models that follow the guidelines summarized in this document will be compatible with tools that further process the model, such as functions like `plotmo` that plot the regression surface.

1 Introduction

Once an R regression model has been built we usually want to use it for further processing. For example, we may want make predictions from the model. Or we may want to plot the model's residuals (which is what `plotres` does), or plot the model's response as the predictors are varied (which is what `plotmo` does¹).

For an S3 model to be amenable to such further processing it should follow some commonly accepted interface standards. These are obvious to experienced developers, but there are many packages on CRAN that don't follow them. This is possibly because there seems to be no summary of the standards.

What are the standards? This document attempts to give a convenient summary, by way of a checklist and an example. The intention is merely to make explicit de facto good practice.

More generally, the S programming book by Venables and Ripley [4] still seems to be the best place for advice on writing S3 model code, although a little dated. (For R programming in general there are of course more modern books, but here we are talking specifically about S3 regression models.) The ultimate reference is the R core code itself, and examples in Venables and Ripley should be checked against that code for current practice.

¹Both those functions are in the `plotmo` package [3], but this document isn't really about those functions. They are just examples.

2 Checklist for S3 regression models

S3 regression models should adhere to the following guidelines. Some of these may be disregarded in certain situations. This isn't a comprehensive list, but enough for most applications.

1. Give the model a unique class. In particular, `class(model)` shouldn't return "list". In the model-building function, do something like `class(model) <- "foo"`.

In general, the class name should be the name of the model-building function. This means, for example, that if the model-building function is `foo`, the `summary` and `plot` methods will be `summary.foo` and `plot.foo`, as expected.

2. Save the call with the model. In the model-building function, do something like `model$call <- match.call()`. This expands any argument names that the user abbreviated to their full names.
3. Provide both formula and x,y model-building functions. Name the formula method `modelclass.formula` and the x,y method `modelclass.default`. Typically both of these call an underlying workhouse function `modelclass.fit`.
4. For model functions with a formula interface, save the `terms` with the model. (A `terms` object is a model formula with additional attributes, as described on the help page for `terms.object`. Additional background is given in Chambers and Hastie [1] and Venables and Ripley Section 4.2 [4].)
5. For model functions with an x,y interface:

- i. Use `x` and `y` as the first two arguments to the model-building function, in that order. Don't call these arguments anything but `x` and `y`, unless that isn't meaningful for your model.
- ii. The x,y interface should be as similar as possible to the formula interface. Where possible, `summary`, `predict`, and friends should work in the same way for models built with the x,y interface and the formula interface.

One acceptable difference between the formula and x,y functions is as follows: The formula interface should convert factors in `x` to indicator columns before doing the regression; the x,y interface should reject factors with an error message.

In the formula interface, conversion of factors comes with the standard use of `model.matrix` (Section 3). In the x,y interface, using `as.matrix` as described below will correctly reject factors and other unsuitable data.

- iii. Be kind to the user and allow `x` and `y` to be `data.frames`, `vectors` (if one-dimensional), or `matrices`. That is, automatically convert to a `matrix` internally in the model function; don't force the user to pre-convert the data. Issue a clear error message when this conversion can't be made.

We suggest `as.matrix` is used for the conversion to `matrix`. Note that `as.matrix` converts *all* columns to strings if there are *any* factors or strings in the input. So to check that the input was converted correctly, you need only check that the first element is numeric, because either all or none of the

converted matrix elements will be numeric. Note that `as.matrix` is efficient in that it will simply return `x` if `x` is already a matrix (it doesn't make a copy of `x`).

Alternatively you can use `data.matrix` to make the conversion. This will convert factors in a `data.frame` to their internal numeric representation. This conversion implicitly assumes that any factors are ordered with equally spaced levels, which isn't true in general. Therefore for most models geared towards continuous data, it's better to issue an error message than to silently make such conversions, i.e., use `as.matrix` rather than `data.matrix` unless you have a special reason not to.

- iv. Consider saving `x` and `y` with the model. If you do, save them in fields named `x` and `y`. Don't use those names for anything else saved with the model.

We recommend `y` is saved as a one-column matrix (not as a vector), with the response name as the column name of the matrix. This allows functions that process the model to easily access the response name for use in plot labels etc.

If `subset` is supported, save `x` and `y` before taking the subset, and also save `subset`. Likewise for `weights`.

A word of explanation. If the data or environment isn't saved with the model, functions like `plotmo` can't unambiguously access the data used to build the model. Saving just the `call` isn't sufficient in general, although it works in the common scenario where the user is working from the command line. (Note that we are talking here about models with an `x,y` interface. For formula-based models, the `call` and `terms` fields suffice.)

If memory use is a concern, give the user an option such as `keep=TRUE` to save `x` and `y`. (There isn't a standard name for this argument — different functions use different names. In our opinion, please *don't* follow the precedent set by `lm` and name the argument `x` or `y`; that can cause confusion.²)

Note that saving `x` and `y` doesn't use as much memory as one might expect, because R will merely create references to `x` and `y`, not make copies of them. On the other hand, R's automatic garbage collection won't be able to release the memory used by `x` and `y` until the model is deleted.

6. Provide a predict method for the model. The first two arguments for the predict method should be `object` and `newdata`.

The default `newdata` should be `NULL` and this should be treated as if the user specified the data used to build the model. If that isn't possible unless `keep` (or similar) was set when building the model, issue an error message to that effect.

The third argument for the predict method should be `type`, unless that isn't meaningful for your model. Make `"response"` one of the options for `type`, possibly the default, unless that isn't meaningful for the model. Apply the `type` argument even with the default `newdata=NULL`; if that isn't possible, issue an error message rather than silently returning bad results.

²This is an exception to the rule that models should conform to the `lm` way of doing things. Note also that `lm.fit` shouldn't be used as an example of an `x,y` interface — because, for example, `predict` can't be used to make predictions on `lm.fit` models. Instead use a `".default"` function as illustrated in Section 3.

Provide defaults for the other arguments where possible so the user can call `predict` with minimum bother. Be kind to the user and allow `newdata` to be a matrix or a `data.frame`. (From `plotmo`'s perspective this is more than just being kind, it's necessary for `plotmo`'s default internal call to `predict`.) You can use `as.matrix` for the conversion to `matrix`, as described in section iii above.

7. If the model supports prediction or confidence levels, allow the user to access those in the same way as `predict.lm`, i.e., when the appropriate arguments are specified, `predict` should return a matrix with column names `fit`, `lwr`, and `upr`.
8. It is good practice to provide the standard model functions. A basic list is `case.names`, `coef`, `fitted`, `model.matrix`, `na.action`, `plot`, `print`, `print.summary`, `residuals`, `summary`, `update`, `variable.names`, and `weights`. Not all of those may apply to your model. Some of them come for free if the model is built in the standard way (the default methods in the `stats` package will automatically work for the model).

Note that `coefficients`, `fitted.values`, and `resid` methods are unnecessary, since the standard functions for these dispatch to `coef`, `fitted`, and `residuals`. For inference the following should be added where applicable: `deviance`, `df.residual`, `logLik`, `nobs`, and `vcov`.

9. Don't use `missing()`. Accomplish the same thing by making the default value of the argument `NA` or some other special value, and checking for that value internally. (The use of `missing` in a function complicates code that calls the function — it has to include two different calls to the function, one with the argument and one without. This can get out of hand if `missing` is used on more than argument.)
10. Allow the user to abbreviate argument names and values. Use `match.arg` or similar to match arguments that take strings.

3 Example S3 Model

[Friedrich Leisch's tutorial](#) [2] is a good introduction to building R packages, and is recommended for a broader context on some of the ideas discussed in this document.

The minimal `linmod` code in the tutorial, although ideal for the purposes of the tutorial, has limitations that can create issues with functions that further process the model. For example

```
data(trees)
fit1 <- linmod(Volume~., data = trees)
predict(fit1, newdata = data.frame(Girth = 10, Height = 80))
```

gives

```
Error in eval(expr, envir, enclos) : object 'Volume' not found
```

and

```
fit2 <- linmod(cbind(Intercept = 1, trees[,1:2]), trees[,3])
predict(fit2, newdata = trees[,1:2])
```

gives

```
Error in x %*% coef(object) : requires numeric/complex matrix/vector arguments
```

Tools like `plotmo` can sometimes be modified to work around these issues, but a more general solution is to modify `linmod`. Figure 1 shows one way of doing that. We can try the code with a few examples:

```
library(plotmo) # for plotres
data(trees)

fit1 <- linmod(Volume~., data=trees)      # formula interface
plotres(fit1) # plot the model's residuals

fit2 <- linmod(trees[,1:2], trees[,3])    # x,y interface
plotres(fit2)
```

Functions like `print.linmod` in Friedrich Leisch's tutorial don't need to be modified for tools like `plotmo`, and don't appear in the code in Figure 1.

The new `linmod.formula` saves the model terms, not just the formula. The new `predict.linmod` accepts a `data.frame` or a `matrix`. This is what users would expect, and is necessary for `plotmo`'s internal call to `predict`. Note also that `linmod.default` doesn't require the user to manually add an intercept column. There are also a few minor changes to the model fields for closer compatibility with `lm`.

4 Limitations of the example S3 model

Production code should include sanity tests that aren't included in our `linmod` example. For example, to prevent confusing downstream error messages, `linmod.fit` should be extended to check that `x` and `y` are numeric, and contain no NAs. From the user's perspective an error message like

```
Error in linmod.fit(x, y) : NA in x
```

is better than the error message issued by the current code

```
Error in qr.default(x) : NA/NaN/Inf in foreign function call (arg 1)
```

And a message like

```
Error in linmod.default(x, y) : non-numeric column in x
```

is better than

```
Error in qr.default(x) : NAs introduced by coercion
```

Similar tests should be made in `predict.linmod`, which should also check that the new data has the correct number of columns. Production code would also handle collinearity properly, ensure `x` and `y` have conformable dimensions, and take care of details like propagating rownames in the input data to the `residuals` and other returned fields.

```

## A simple linear model (new version of linmod from Friedrich Leisch's tutorial).
## Functions like print.linmod in the tutorial don't need to be modified for tools
## like plotmo, and don't appear in the code below.

linmod <- function(...) UseMethod("linmod")

linmod.fit <- function(x, y) # internal function, not for the casual user
{
  # first column of x is the intercept (all 1s)

  y <- as.vector(as.matrix(y))      # necessary when y is a data.frame
  qx <- qr(x)                        # QR-decomposition of x
  coef <- solve.qr(qx, y)            # compute (x'x)^(-1) x'y
  df.residual <- nrow(x) - ncol(x)   # degrees of freedom
  sigma2 <- sum((y - x %*% coef)^2) / df.residual # variance of residuals
  vcov <- sigma2 * chol2inv(qx$qr)   # covar mat is sigma^2 * (x'x)^(-1)
  colnames(vcov) <- rownames(vcov) <- colnames(x)
  fitted.values <- qr.fitted(qx, y)

  fit <- list(coefficients = coef,
              residuals    = y - fitted.values,
              fitted.values = fitted.values,
              vcov         = vcov,
              sigma       = sqrt(sigma2),
              df.residual = df.residual)

  class(fit) <- "linmod"
  fit
}

linmod.default <- function(x, y, ...)
{
  fit <- linmod.fit(cbind("(Intercept)"=1, as.matrix(x)), y)
  fit$call <- match.call()
  fit
}

linmod.formula <- function(formula, data=parent.frame(), ...)
{
  mf <- model.frame(formula=formula, data=data)
  terms <- attr(mf, "terms")
  fit <- linmod.fit(model.matrix(terms, mf), model.response(mf))
  fit$terms <- terms
  fit$call <- match.call()
  fit
}

predict.linmod <- function(object, newdata=NULL, ...)
{
  if(is.null(newdata))
    y <- fitted(object)
  else {
    if(is.null(object$terms)) # x,y interface
      x <- cbind(1, as.matrix(newdata)) # columns must be in same order as orig x
    else { # formula interface
      terms <- delete.response(object$terms)
      x <- model.matrix(terms, model.frame(terms, as.data.frame(newdata)))
    }
    y <- as.vector(x %*% coef(object))
  }
  y
}

```

Figure 1: *A simple linear model. This code can be used for new models by replacing all occurrences of “linmod” with the new model name and by rewriting the function linmod.fit.*

Code that includes the checks mentioned here and the guidelines in Section 2 can be found at www.milbo.org/doc/linmod.R. We suggest that this more complete code is used as a template for new S3 models (rather than the code in Figure 1).

Our code handles the model formula in a very basic way. See [4] for more advanced handling, including `weights` and `subset`, and [5] for a more modern treatment of formulas with extensions.

A thanks goes out to Achim Zeileis for his feedback.

References

- [1] J.M. Chambers and T.J. Hastie. *Statistical Models in S*. Chapman and Hall/CRC, 1991. Cited on page 2.
- [2] Friedrich Leisch. *Creating R Packages: A Tutorial*. Compstat Proceedings in Computational Statistics, 2008. <https://CRAN.R-project.org/doc/contrib/Leisch-CreatingPackages.pdf>. Cited on page 4.
- [3] S. Milborrow. *plotmo: Plot a model's response and residuals*, 2015. R package, <https://CRAN.R-project.org/package=plotmo>. Cited on page 1.
- [4] W.N. Venables and B.D. Ripley. *S Programming*. Springer, 2000. <http://www.stats.ox.ac.uk/pub/MASS3/Sprog>. Cited on pages 1, 2, and 7.
- [5] Achim Zeileis and Yves Croissant. *Extended Model Formulas in R: Multiple Parts and Multiple Responses*. Journal of Statistical Software, 2010. <https://cran.r-project.org/web/packages/Formula/vignettes/Formula.pdf>. Cited on page 7.