

The purpose of this section is to introduce you to writing your own packages in R, and to get you started on calculating R^2 values.

Last Section

I'm looking forward to seeing your problem sets nicely typed up using **knitr**. While it isn't explicitly required, I strongly urge you to do so.

I usually write R code using RStudio but edit `.Rnw` files using TexStudio, as it is mostly L^AT_EX. Instructions for using TexStudio to compile from `.Rnw` to `.pdf` can be found at:

<http://yihui.name/knitr/demo/editors/>

My build command (which is discussed in the link) is, for example:

```
"C:/Program Files/R/R-3.1.2/bin/x64/Rscript.exe" -e "library(knitr); knit('
%.Rnw')|"C:/Program Files (x86)/MiKTeX 2.9/miktex/bin/pdflatex.exe" -inter
action=nonstopmode %.tex|txs:///view-pdf-internal
```

The user defined build command can be accessed using **Alt + Shift + F1**. What the above commands do is 1) open an R instance on the fly and run the `knit()` function on the open `.Rnw` file, creating a `.tex` file in the same folder. Then 2) run `pdflatex` on the created `.tex` file, creating/updating the `.pdf` file, then 3) opening the internal pdf viewer to see the new/updated `.pdf`.

As with **Git**, a section hour is not enough time to really know how to use **knitr**, so you should read up/Google a bit more on what it can do.

Packages

It's likely that most of you will not develop packages in R for the purposes of sharing code with others. However, you will almost certainly end up writing utility functions along the way that you want to keep for yourself, and packages are a great way to nicely collate and store this code.

For a good exposition on the basics of writing packages, see:

<http://r-pkgs.had.co.nz/>

I was extremely surprised at how simple package development is using RStudio. Let's see how it works!

Getting Started

First things first, install the `devtools` and `roxygen2` packages:

```
install.packages("devtools")
install.packages("roxygen2")
```

Next, in RStudio, go to **File** → **New Project** and choose “New Directory”, then “R Package”. Name your package, uncheck “Create a git repository” (you can do this yourself later) and click “Create Project”.

My version of RStudio creates a folder named whatever-I-chose-for-the-package-name that contains 1) DESCRIPTION, NAMESPACE, and .Rbuildignore files, and 2) man, and R folders. The absolute essentials are a folder named R that contains some R code, a DESCRIPTION file, and a NAMESPACE file.¹

The DESCRIPTION file

Open your DESCRIPTION file in RStudio by clicking on it in the Files view in the bottom right corner. Fill in Title and Description with something appropriate.

The DESCRIPTION file contains important metadata about your package. Who wrote it, what it's for, who can use it, etc.

Authors@R

I like to remove the Author and Maintainer fields and replace them with the Authors@R field which handles both. You can do something like:

```
Authors@R:person("Kendon", "Bell", email = "kmb56@berkeley.edu",  
  role = c("aut", "cre"))
```

This says I am both the author (aut) and the maintainer (cre) of the package - be sure to include the spaces before the second line!

License

If you ever want to share your code, you should read up on the different licenses at:

<http://r-pkgs.had.co.nz/description.html#license>

Otherwise it doesn't matter. Important fields that are not included by default in the RStudio default are Depends, and Imports.

Depends

You should only use Depends to specify a version of R that is required for your package to work.² For example:

```
Depends: R (>= 3.1.2)
```

Imports

Imports is used to load packages that functions in your package will use. For example:

¹See <http://r-pkgs.had.co.nz/description.html#description> for a detailed discussion of the DESCRIPTION file, and <http://r-pkgs.had.co.nz/namespace.html> for a detailed description of the NAMESPACE file.

²See Wickham's website for the few exceptions to this rule.

```
Imports:
  dplyr (>= 0.3.0.1),
  ggplot2 (>= 1.0.0)
```

Again, make sure you have the spaces at the start of the second and third lines. We won't worry about the `NAMESPACE` file for now.

R Code

Now, let's add an R function to the package. Create a `.R` file containing just the OLS function that you created last week and save it in the `/R` folder in your package directory.

```
OLS <- function(y, X) {
  b <- solve(t(X) %*% X) %*% t(X) %*% y
  b
}
```

Now, in RStudio, hit the **Build & Reload** button in the **Build** tab in the top right corner. Ask me questions if you get a error - especially sensitive is the `DESCRIPTION` file. The syntax needs to be perfect. Make sure that, if you have fields that go for more than one line, you indent them with 2 spaces.

Go to the console, type `OLS` (without the `()`), and hit **ENTER**. You should now see the function definition as we wrote it last week. What R is doing when it 'builds' your function is simply sourcing (or running) all the `.R` files in the `/R` folder.

Now, open the `NAMESPACE` file. It should read something like `exportPattern("~[:alpha:]]+")` by default. Remove that line, save the `NAMESPACE` file, and **Build & Reload** again.

Now, you shouldn't see the `OLS` function. What happened??

The export line in the `NAMESPACE` file is telling R which functions to load, and you're no longer telling it to load everything using `exportPattern("~[:alpha:]]+")`.³ So, it's not really running *all* the R code, just what is specified in the `NAMESPACE` file. You can revert the file to how it was and save it.

Help File Documentation

Package documentation is very simple using the package `roxygen2`. `roxygen2` allows you to ignore most of the dirty details of how documentation works, which is great.

Note that "documentation" in coding can refer to both the comments a coder makes throughout her functions, and the help files that are produced for users. We're talking about the latter here.

³We will cover regular expressions near the end probably.

Now, lets document our OLS function:

```
#' @title OLS Coefficients
#' @description Calculates OLS coefficients using matrix methods.
#' @param y \code{numeric} or \code{matrix}. The left hand
#'         side variable of the model to be estimated.
#' @param X \code{matrix}. Contains the right hand side variables of the
#'         model to be estimated. Must include a column of 1's if you
#'         wish to include an intercept.
#' @return A \code{matrix} of linear regression coefficients.
OLS <- function(y, X) {
  b <- solve(t(X) %*% X) %*% t(X) %*% y
  b
}
```

Above, we can see the basic structure of **roxygen2** documentation where each line begins with '#', and fields are all prefixed by an '@' symbol. The only essential field in the above is **@title**. There are also several other fields you may use, which you can explore in the **roxygen2** documentation online.

@title, **@description**, and **@details** (omitted here) are differing levels of description of what the function does. These are simply used by typing the field name followed by a space and the content (all 3 can go onto multiple lines). **@details** is usually only used for more complicated functions and will typically contain more than one paragraph.

@param documents the function arguments. This is used by typing the field name, followed by a space and the argument name, followed by a space and a description of the argument.

@return documents the return value.

It is good practice for both **@param** and **@return** to include the expected classes of the variables involved.

There is a bunch of syntax you can use to style the documentation, such as `\code` above. A full list is available at:

<http://cran.r-project.org/doc/manuals/r-release/R-exts.html#Marking-text>

Let's now get **roxygen2** doing its thing! Make sure the package is installed and call:

```
roxygen2::roxygenize()
```

The syntax: `packageName::functionName()` allows you to call a function from a package without explicitly loading the function using `library()`. Using `::` allows you to access functions that are only available internally to the package.

Now, build/reload the package again using **Ctrl/Cmd + Shift + B** or clicking the button.

You now should be able to see your documentation by typing:

```
?OLS
```

in the console. You've just made your first R package! You can now read in your package in any R session using `library(yourPackage)`.

As a last note, `roxygen2` can also automatically generate a `NAMESPACE` file for you if you have some functions you want to keep invisible. You can do this by deleting the `NAMESPACE` file and adding the tag `@export` to the functions you want to be visible to the user.

Calculating R^2

Now, we'll go over how to calculate the R^2 values. We're going to write the code so that we can retrieve the centered R^2 , with *or* without an intercept in our model. Computing R^2 values can be trickier than it seems, particularly since many statistical packages (including R) are not fully transparent about which R^2 they are displaying.

Remember, the centered R^2 is a measure of how much of the variance in our model we can explain relative to a hyperplane where $y = \bar{y}$.⁴

Now the math:

$$R^2 \equiv 1 - \frac{\mathbf{e}'\mathbf{e}}{\mathbf{y}^*\mathbf{y}^*} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} = 1 - \frac{\text{SSR}}{\text{SST}} \quad (1)$$

Remember that $\mathbf{y}^* = \mathbf{A}\mathbf{y}$, the demeaned version of \mathbf{y} .

Now, lets add a demeaning function to our .R file, save it and rebuild the function.

```
demeanMat <- function(n) {
  ones <- rep(1, n)
  diag(n) - (1/n) * ones %*% t(ones)
}
```

Now, the main event! We'll define a function that runs OLS, computes our R^2 value, and returns the results of both in a list.

An object of the list class can store multiple different types of data, not necessarily in a rectangular format. Lists are great for results output like this.

There is a bit of code here to process, but if you refer to the math above you'll see that it all matches up. Be sure to add this function to your package.

```
olsRSq <- function(y, X) {
  n <- nrow(X)
  k <- ncol(X)
  b <- OLS(y,X)
  yh <- X %*% b
  e <- y - yh # yh is y hat, the predicted value for y

  SSR <- t(e) %*% e

  A <- demeanMat(n)
  ys <- A %*% y # this is ystar
  SST.yb <- t(ys) %*% ys # == sum((y - mean(y)^2))
  R2.cen <- 1 - SSR / SST.yb
```

⁴When you see "hyperplane," think "line."

```
  return(list(rsquared.cen=R2.cen, coefficients=b))
}
```

Now let's test our package with the `auto.dta` dataset. We'll also confirm our results with `lm()`. You can put something like the following in an `.R` file.

```
library(foreign)
library(misc212)
mydata <- read.dta("auto.dta")
names(mydata) <- c("price", "mpg", "weight")
y <- matrix(mydata$price)
X <- cbind(1, mydata$mpg)
olsRSq(y, X)

## $rsquared.cen
##           [,1]
## [1,] 0.2195829
##
## $coefficients
##           [,1]
## [1,] 11253.0607
## [2,] -238.8943

summary(lm(y ~ X))$r.squared
## [1] 0.2195829
```

Now let's try with `X2`, the data matrix with `mpg` but *no intercept*.

```
X2 <- cbind(mydata$mpg)
olsRSq(y, X2)

## $rsquared.cen
##           [,1]
## [1,] -0.7817091
##
## $coefficients
##           [,1]
## [1,] 253.6302

summary(lm(y ~ 0 + X2))$r.squared
## [1] 0.6718225
```

Our result for R^2 is totally insane. Does this make any sense? Actually, yes! Remember that the centered R^2 compares our model, $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ (with no intercept) to $\mathbf{y} = \bar{\mathbf{y}}$. Since $\mathbf{y} = \bar{\mathbf{y}}$ is actually a much better predictor of \mathbf{y} than $\mathbf{X}\mathbf{b}$, we get that $\text{SSR} > \text{SST}$. Using that $R^2 = 1 - \frac{\text{SSR}}{\text{SST}}$, we can see why we get a negative R^2 .

But, our results don't match the results from `lm()`. This is because `lm()` is being tricky: since we aren't passing it an intercept, it computes the uncentered R_{uc}^2 instead of the centered R^2 , which

makes more sense as a measure of fit for a model with no intercept.⁵ So rather than comparing our model $\mathbf{y} = \mathbf{X}\mathbf{b} + \boldsymbol{\varepsilon}$ to $\mathbf{y} = \bar{\mathbf{y}} + \boldsymbol{\varepsilon}$, it compares it to $\mathbf{y} = \mathbf{0} + \boldsymbol{\varepsilon}$.

Linear algebra puzzles

1. Define vectors $\mathbf{x} = [1 \ 2 \ 3]'$, $\mathbf{y} = [2 \ 3 \ 4]'$, and $\mathbf{z} = [3 \ 5 \ 7]$. Define $\mathbf{W} = [\mathbf{x} \ \mathbf{y} \ \mathbf{z}]$. Calculate \mathbf{W}^{-1} . If you cannot take the inverse, explain why not and adjust \mathbf{W} so that you /can/ take the inverse. *Hint*: the `solve()` function will return the inverse of the supplied matrices.
2. Show, somehow, that $(\mathbf{X}')^{-1} = (\mathbf{X}^{-1})'$.
3. Generate a 3×3 matrix \mathbf{X} , where each element is drawn from a standard normal distribution. Let $\mathbf{A} = \mathbf{I}_3 - \frac{1}{3}\mathbf{B}$ be a demeaning matrix, with $\mathbf{1}$ a 3×3 matrix of ones. First show that \mathbf{A} is idempotent and symmetric. Next show that each row of the matrix \mathbf{XA} is the deviation of each row in \mathbf{X} from its mean. Finally, show that $(\mathbf{XA})(\mathbf{XA})' = \mathbf{XAX}'$, first through algebra and then R code.
4. Demonstrate from random matrices that $(\mathbf{XYZ})^{-1} = \mathbf{Z}^{-1}\mathbf{Y}^{-1}\mathbf{X}^{-1}$.
5. Let \mathbf{X} and \mathbf{Y} be square 20×20 matrices. Show that $\text{tr}(\mathbf{X} + \mathbf{Y}) = \text{tr}(\mathbf{X}) + \text{tr}(\mathbf{Y})$.
6. Generate a diagonal matrix \mathbf{X} , where each element on the diagonal is drawn from $U[10, 20]$. Now generate a matrix \mathbf{B} s.t. $\mathbf{X} = \mathbf{BB}'$. *Hint*: There is a method in R that makes this easy. Does the fact that you can generate \mathbf{B} tell you anything about \mathbf{X} ?
7. Demonstrate that for any scalar c and any square matrix \mathbf{X} of dimension n that $\det(c\mathbf{X}) = c^n \det(\mathbf{X})$.
8. Demonstrate that for an $m \times m$ matrix \mathbf{A} and a $p \times p$ matrix \mathbf{B} that $\det(\mathbf{A} \otimes \mathbf{B}) = \det(\mathbf{A})^p \det(\mathbf{B})^m$. *Hint*: Note that \otimes indicates the Kronecker product⁶. Google the appropriate R function.

Here's a nice problem set hint for those of you who actually read these notes! Check out the `WDI` package in R; it should make downloading the World Bank data much simpler. This is a canned package which is definitely OK to use for the problem set.

⁵For more information, see here: <http://bit.ly/1c0nA6N>.

⁶The Kronecker product is a useful mathematical tool for econometricians, allowing us to more easily describe block-diagonal matrices for use in panel data settings. I wouldn't lose sleep over it, though.