

NYU R Workshop

Thinking in R

<http://drewdimmery.com/r-workshop-2013/>

Drew Dimmery



NEW YORK UNIVERSITY

February 8, 2013

Outline

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- 1 Outline
- 2 Why R?
- 3 Where to start
- 4 Where to get help
 - Documentation
 - CRAN
 - Journal of Statistical Software
- 5 Misc. Quirks
- 6 Using third party code
- 7 Data Types
 - Character
 - Numeric
 - Logical
 - Objects
- 8 Data Structures
 - Vectors
 - Lists
 - Matrices
 - Dataframes
- 9 Control Flow
 - Loops
 - If/Else
 - Functions
 - Useful Functions
- 10 Working with Data
 - Input
 - Output
- 11 Plotting
- 12 Tying it all Together
- 13 Good Habits

Why R?

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Free!
- › Open!
- › Wide array of statistical procedures
- › Extensible
- › Interface-able
- › Large volumes of online support

How to setup a workflow

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › <http://www.r-project.org/>
- › Traditional Console and text editor
 - › Notepad++
 - › vim
 - › Emacs
 - › Kate, Textmate, Sublime, etc
- › **Rstudio**
- › Other IDEs (Eclipse, Netbeans, etc)

Rstudio

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

The screenshot displays the RStudio interface with three main panels:

- Source Editor:** Contains R code for a least squares regression. The code defines a function `loc.lin` that takes a matrix `X` and a vector `y` as input. It calculates the least squares estimate `beta` using the normal equations $\beta = (X^T X)^{-1} X^T y$. The function also returns the residuals `res` and the variance-covariance matrix of the coefficients `vc`. The code includes comments in Italian and a call to `loc.lin` with specific data.
- Console:** Shows the execution of the `loc.lin` function. It displays the values of `beta`, `res`, and `vc`. The output is:

```
[1] "a" "b" "c"  
[[1]] first = c(1.2,2), second = c(7^4)  
$first  
[1] 1 2 3  
  
$second  
[1] "a"  
[[1]] first = c(7^4,"b","c"), second = c(09.8673)  
$first  
[1] "a" "b" "c"  
  
$second  
[1] 99.867  
  
[[1]] first = c(7^4,"b","c"), second = c(09.8673)$first  
[1] "a" "b" "c"  
[[1]] first = c(7^4,"b","c"), second = c(09.8673)[[1]]  
[[1]] "a" "b" "c"  
##second[1] "b" "c"  
- [[1]] first = 7  
- [[1]] second = 3  
print("hello world")  
[1] "hello world"
```
- Plots Panel:** Displays a scatter plot of data points (black dots) and a fitted curve (red line). The x-axis is labeled 'x' and ranges from 0 to 5. The y-axis is unlabeled but ranges from approximately -1 to 1. The data points show a non-linear trend, and the fitted curve follows the general shape of the data.

So your code won't run. What now?

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Documentation
Example: ?lm
- › Google
- › CRAN: Reference Manual, vignettes
Example: <http://cran.r-project.org/web/packages/AER/index.html>
- › Journal of Statistical Software
Example: <http://www.jstatsoft.org/v27/i02>
- › Stack Overflow
- › Listservs

How to read R documentation

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › <http://stat.ethz.ch/R-manual/R-patched/library/stats/html/lm.html>
- › **Call:** `?lm`
- › **Sections:**
 - › Description
 - › Usage
 - › Arguments
 - › Details
 - › Value
 - › Author(s)
 - › References
 - › See Also
 - › Examples

› AER on CRAN: <http://cran.r-project.org/web/packages/AER/index.html>

AER: Applied Econometrics with R

Functions, data sets, examples, demos, and vignettes for the book Christian Kleibler and Achim Zeileis (2008), Applied Econometrics with R, Springer-Verlag, New York. ISBN 978-0-387-77316-2. (See the vignette for a package overview.)

Version: 1.1-9
 Depends: R (≥ 2.10.0), stats, car (≥ 2.0-1), Formula (≥ 0.2-0), ltest, sandwich, strucchange, survival, zoo
 Imports: stats
 Suggests: book, dplyr, effects, foreign, lme4, KernSmooth, lattice, MASS, nlmixt, nlme, nnet, np, plus, pls, quantreg, ROCR, sampleSelection, scatterplot3d, systemfit, rgl, traucreg, tseries, xgboost
 Published: 2011-12-10
 Author: Christian Kleibler [aut], Achim Zeileis [aut, cre]
 Maintainer: Achim Zeileis <Achim.Zeileis@R-project.org>
 License: GPL-2
 Citation: [AER citation info](#)
 In view of: [Econometrics](#), [Survival](#), [TimeSeries](#)
 CRAN checks: [AER results](#)

Download:

Package source: [AER_1.1-9.tar.gz](#)

MacOS X binary: [AER_1.1-9.tgz](#)

Windows binary: [AER_1.1-9.zip](#)

Reference manual: [AER.pdf](#)

Vignettes: [Applied Econometrics with R Package Vignette and Errata](#)
[Sutative Example: Linear Regression for Economics Journals Data](#)

News ChangeLog: [NEWS](#)

Old sources: [AER archive](#)

Reverse dependencies:

Reverse depends: [rdd](#), [toyomic](#)

Reverse suggests: [cmfReg](#), [ltest](#), [micEconCES](#), [nlmixt](#), [pls](#), [REEMtree](#), [sandwich](#)

› `plm` in JSS: <http://www.jstatsoft.org/v27/i02/>



Journal of Statistical Software

July 2008, Volume 27, Issue 2.

<http://www.jstatsoft.org/>

Panel Data Econometrics in R: The `plm` Package

Yves Croissant

Université Lumière Lyon 2

Giovanni Millo

University of Trieste and Generali SpA

Abstract

Panel data econometrics is obviously one of the main fields in the profession, but most of the models used are difficult to estimate with R. `plm` is a package for R which intends to make the estimation of linear panel models straightforward. `plm` provides functions to cope with some of the typical problems associated with economic data, first of all that of unobserved heterogeneity.

Keywords: panel data, covariance matrix estimators, generalized method of moments, R.

1. Introduction

Panel data econometrics is a continuously developing field. The increasing availability of data observed on cross-sections of units (like households, firms, countries etc.) and over time has given rise to a number of estimation approaches exploiting this double dimensionality to cope with some of the typical problems associated with economic data, first of all that of unobserved heterogeneity.

Timewise observation of data from different observational units has long been common in other fields of statistics (where they are often termed *longitudinal* data). In the panel data field as well as in others, the econometric approach is nevertheless peculiar with respect to experimental contexts, as it is emphasizing model specification and testing and tackling a

Things that can be confusing

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › At the prompt, > means you're good to go, + means a parenthesis or bracket is open.
- › Case sensitive
- › Use / in path names.
- › R uses variables – there is no "sheet" here, like in Stata
- › **R is a programming language**

Using third party code

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › **Relevant commands:** `install.packages`, `library`
- › How does one find and install an appropriate package?
- › **Example:**

```
1 ?covariance
2 ??covariance
3 install.packages("sandwich")
4 library("sandwich")
5 ?vcovHC
```

- › In reality, Google to see if an appropriate package exists.

Data Types

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Character – strings
- › Double / Numeric – numbers
- › Logical – true/false
- › Factor – unordered categorical variables
- › Objects – it's complicated

Character

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

› **Example:** `str<-"This is a string"`

› **String concatenation:**

```
1 > paste("This", "is", "a", "string", sep=" ")
2 [1] "This is a string"
```

› **Numbers can be strings:** `as.character(99)`

› **Possible to determine type:** `class(str)`

Numeric

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

› Most of the time you'll work with these.

› `num<-99.867`

› `round(str, digits=2)`

› Many useful (and not often useful) operations:

› `sin`

› `exp`

› `log`

› `factorial`

› `choose`

› `BesselJ`

› `pi` (not a function)

› **And many more**

Logical

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › The logical type allows us to make statements about truth

```
1 > 2==4
2 [1] FALSE
3 > str != num
4 [1] TRUE
5 > "34" == 34
6 [1] TRUE
```

- › Equal to: ==
- › Not equal to: !=
- › Greater than >
- › Greater than or equal to: >=
- › Less than: <
- › Less than or equal to: <=
- › Not: !
- › Logical and: &
- › Logical or: |
- › Any are true: any
- › All are true: all

Objects

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Many functions will return objects rather than a single datatype
- › Example (no, this won't work, since the data doesn't exist):

```
1 > out.lm <- lm(Y~X)
2 > class(out.lm)
3 [1] "lm"
```

- › Objects can have other data embedded inside them.

```
1 > out.lm$rank
2 [1] 2
```

- › Sometimes operations will work differently for different sorts of objects.

Other ways to “hold” data

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › **Vectors**
- › **Lists**
- › **Matrices**
- › **Dataframes**

Vectors

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › **Everything** is a vector (almost).
- › Just like a vector in math

```
1 > as.vector(4)
2 [1] 4
3 > 4
4 [1] 4
```

- › We can combine vectors with `c`

```
1 > vec <- c("a", "b", "c")
2 > vec
3 [1] "a" "b" "c"
4 > c(2, 3, c(4, 5))
5 [1] 2 3 4 5
6
7 > c(1, 2, 3, 4) + c(1, 2)
8 [1] 2 4 4 6      #Recycles shorter vector
```

Vectors

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Vectors are indexed as you would probably expect – just like in math. Or you can name elements yourself.

```
1 > vec[1]
2 [1] "a"
3 > names(vec) <- c("first", "second", "third")
4 > names(vec)
5 [1] "first" "second" "third"
6 > vec$first
7 [1] "a"
8 > vec["first"]
9 [1] "a"
```

- › Vectors can include missing values via NA

```
1 > vec[1] <- NA
2 > vec
3 [1] NA b c
4 > is.na(vec)
5 [1] TRUE FALSE FALSE
6 > vec[!is.na(vec)] #vec[complete.cases(vec)]
7 [1] "b" "c"
```

Lists

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Lists are sort of like a vector, but they allow for mixing of types and lengths in arbitrary ways

```
1 > listie <- list(first = vec , second = num)
2 > listie
3 $first
4 [1] "a" "b" "c"
5
6 $second
7 [1] 99.867
8
9 > listie[[1]]
10 [1] "a" "b" "c"
11
12 > listie$first
13 [1] "a" "b" "c"
```

- › One could even save an object to one element of a list.

Matrices

NYU R
Workshop
Drew Dimmery

Outline
Why R?
Where to start
Where to get help
Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › We can create matrices too. They're just like normal matrices.

$$A = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$$

A_{ij}

$A_{1,2} = 3$

$A_{1,\cdot} = (1,3)$

```
1 > A <- matrix(c(1,2,3,4),nrow=2,nrow=2)
2 > A
3      [,1] [,2]
4 [1,]    1    3
5 [2,]    2    4
6 > A[1,2]
7 [1] 3
8 > A[,2]
9 [1] 3 4
10 > A[1,]
11 [1] 1 3
```

Matrices

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › We can also perform matrix operations

```
1 > solve(A) #A^{-1}
2   [,1] [,2]
3 [1,]  -2  1.5
4 [2,]   1 -0.5
5 > 10*A
6   [,1] [,2]
7 [1,]  10  30
8 [2,]  20  40
9 > B<-diag(c(1,2))
10 > B
11   [,1] [,2]
12 [1,]   1   0
13 [2,]   0   2
14 > A%*%B
15   [,1] [,2]
16 [1,]   1   6
17 [2,]   2   8
```

Matrices

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

› Continued

```
1 > t(A)
2      [,1] [,2]
3 [1,]    1    2
4 [2,]    3    4
5 > rbind(A,B)
6      [,1] [,2]
7 [1,]    1    3
8 [2,]    2    4
9 [3,]    1    0
10 [4,]    0    2
11 > cbind(A,B)
12     [,1] [,2] [,3] [,4]
13 [1,]    1    3    1    0
14 [2,]    2    4    0    2
15 > c(1,2,3)%x%c(1,1) #Kronecker Product
16 [1] 1 1 2 2 3 3
```

Matrices

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

› Naming things

```
1 > rownames(A)
2 NULL
3 > rownames(A) <- c("a", "b")
4 > colnames(A) <- c("c", "d")
5 > A
6   c d
7 a 1 3
8 b 2 4
9 > A[, "d"]
10 a b
11 3 4
12
13 > A[3] #Matrices are vectors
14 [1] 3
```


Dataframes

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › The workhorse.
- › Basically just a matrix that allows mixing of type by column

```
1 > rstudio::viewData(ged)
```

21,860 observations of 40 variables

priogrid_gid	Date_start	Date_end	Date_prec	Event_type	Deaths_A	Deaths_B	Civilian_Deaths	Unknown	Best_est	High_est	Low_est	Deathsplit
182530	2006-10-14	2006-10-14	1	1	1	0	0	0	1	1	1	0
182530	2006-09-02	2006-09-02	1	1	4	0	1	0	5	5	5	0
182530	2005-01-21	2005-01-21	1	1	0	0	0	0	0	1	0	0
166478	2009-10-28	2009-11-04	4	1	0	1	0	0	1	1	1	0
182528	2009-10-06	2009-10-06	1	1	0	1	0	0	1	1	1	0
182528	1999-10-09	1999-10-10	2	1	1	0	1	0	2	2	2	0
182528	2002-02-03	2002-02-03	1	1	1	0	0	0	1	1	1	0
178925	1996-12-24	1996-12-24	1	1	0	7	0	0	7	7	7	0
178925	1998-09-07	1998-09-07	1	1	1	0	0	0	1	1	1	0
181809	2004-05-14	2004-05-14	1	1	1	0	0	0	1	1	1	0
182529	2010-02-15	2010-02-15	1	1	0	1	0	0	1	1	1	0
182529	2010-00-26	2010-00-26	1	1	1	0	0	0	1	1	1	0
181806	1999-01-06	1999-01-07	2	1	5	0	0	0	5	5	5	0
182528	2006-09-02	2006-09-02	1	1	1	0	0	0	1	1	1	0
181895	1993-05-16	1993-05-16	1	1	1	0	0	0	1	1	1	0
182536	2004-10-08	2004-10-08	1	1	0	4	0	0	4	4	4	0
182536	1999-07-17	1999-07-17	1	1	1	0	0	0	1	1	1	0
181895	1998-01-22	1998-01-22	1	1	0	2	0	0	2	2	2	0
181816	1997-07-16	1997-07-16	1	1	0	1	0	1	1	1	1	0

Displayed 1000 rows of 21,860 (20,860 omitted)

Control Flow

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › loops
- › if/else
- › functions
- › Useful Functions

Loops

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › for loops – just a way of saying “do this **for** each element of the index”
- › "this" is defined in what follows the "for" expression.

```
1 > for(i in 1:5) {
2 +   cat(i*10, " ")
3 + }
4 10  20  30  40  50
5
6 > for(i in 1:length(vec)) {
7 +   cat(vec[i], " ")
8 + }
9 a  b  c
10
11 > for(i in vec) {
12 +   cat(i, " ")
13 + }
14 a  b  c
```

If/Else

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › To conditionally execute code, we use if/else statements.

```
1 > if(vec[1]=="a") print("Hello World!")
2 [1] "Hello World!"
3
4 > if(vec[1]=="b") {
5 +   print("Hello World!")
6 + } else {
7 +   print("!dlroW olleH")
8 + }
9 [1] "!dlroW olleH"
```

Vectorized If/Else

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › What if we want to conditionally execute something for each element in a vector?

```
1 > new <- vector(length=length(vec))
2 > for(i in 1:length(vec)) {
3 +   if(vec[i]=="a") {
4 +     new[i]<-13
5 +   } else {
6 +     new[i]<-0
7 +   }
8 + }
9 > new
10 [1] 13  0  0
11
12 > new <- ifelse(vec=="a",13,0)
13 > new
14 [1] 13  0  0
```

Functions

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

› $f : X \rightarrow Y$

› Functions in R are essentially the same.

```
1 > add3<-function(X) {
2 + return(X+3)
3 + }
4 > add3(2)
5 [1] 5
6
7 > makeGroups<-function(groups,members=1) {
8 +   return((1:groups)%x%rep(1,members))
9 + }
10 > makeGroups(5)
11 [1] 1 2 3 4 5
12 > makeGroups(5,2)
13 [1] 1 1 2 2 3 3 4 4 5 5
```

› Drew will now get on a function soapbox

Useful functions

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › **Note: Most functions don't do complete case analysis by default (usually option `na.rm=TRUE`)**
- › `print`
- › `complete.cases`
- › `vcov`
- › `cat`
- › `subset`
- › `residuals`
- › `paste`
- › `merge`
- › `vcovHC (from sandwich)`
- › `with`
- › `mean`
- › `ivreg (from AER)`
- › `length`
- › `sum`
- › `countrycode (from countrycode)`
- › `sort`
- › `sd`
- › `order`
- › `var`
- › `unique`
- › `lag`
- › `summary`
- › `rep`
- › `lm`
- › `pdf`
- › `nrow`
- › `model.matrix`
- › `plot`
- › `ncol`
- › `coef`
- › **Tools from `plm`**
- › **Many more.**

Distributional functions

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › ?Distributions
- › **They have a consistent naming scheme.**
- › `rnorm`, `dnorm`, `qnorm`, `pnorm`
- › `rdist` – generate random variable from dist
- › `ddist` – density function of dist
- › `qdist` – quantile function of dist
- › `pdist` – distribution function of dist
- › look at documentation for parameterizations.
- › **Example:** `rnorm(100)` generates 100 samples from a standard normal

Cluster Robust SEs

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

```
1 robust.se <- function(model,
2                       cluster=1:length(model$residuals)) {
3   require(sandwich)
4   require(lmtest)
5   M <- length(unique(cluster))
6   N <- length(cluster)
7   K <- model$rank
8   dfc <- (M/(M - 1)) * ((N - 1)/(N - K))
9   uj <- apply(
10      estfun(model),
11      2,
12      function(x) tapply(x, cluster, sum)
13    )
14   rcse.cov <- dfc * sandwich(model,meat = crossprod(uj)/N)
15   rcse.se <- coefest(model, rcse.cov)
16   return(list(rcse.cov, rcse.se))
17 }
```

Functions and Objects

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Sometimes, passing an object to a function will do something you might not expect.

```
1 > summary(vec)
2   Length      Class      Mode
3     3 character character
4
5 > summary(c(1,2,3,4))
6   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
7   1.00  1.75   2.50   2.50  3.25   4.00
8
9 > summary(data.frame(A))
10          c          d
11 Min.   :1.00   Min.   :3.00
12 1st Qu.:1.25   1st Qu.:3.25
13 Median :1.50   Median :3.50
14 Mean   :1.50   Mean   :3.50
15 3rd Qu.:1.75   3rd Qu.:3.75
16 Max.   :2.00   Max.   :4.00
```

- › Why? `?summary`, `?summary.lm`

The *apply family

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › These functions allow one to efficiently perform a large number of actions on data.
- › `apply` – performs actions on the rows or columns of a matrix (1 for rows, 2 for columns)
- › `sapply` – performs actions on every element of a vector
- › `tapply` – performs actions on a vector by group
- › `replicate` – performs the same action a given number of times

apply

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › `apply` – performs actions on the rows or columns of a matrix (1 for rows, 2 for columns)

```
1 > A
2   c d
3 a 1 3
4 b 2 4
5 > apply(A, 1, sum)
6 [1] 4 6
7 > apply(A, 2, mean)
8 [1] 1.5 3.5
```

supply

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › `sapply` – performs actions on every element of a vector

```
1 > vec
2 [1] "a" "b" "c"
3 > sapply(vec, function(x) paste0(x, ".vec"))
4           a           b           c
5 "a.vec" "b.vec" "c.vec"
6
7 #Bad example, though
8 > paste0(vec, ".vec")
9 [1] "a.vec" "b.vec" "c.vec"
```

- › `replicate` is basically just `sapply(1:N, funct)` where `funct` never uses the index.

› `tapply` – performs actions on a vector by group

```
1 > tapply(1:10,  
2 +   makeGroups(5,2),  
3 +   function(x) mean(x)  
4 + )  
5   1   2   3   4   5  
6 1.5 3.5 5.5 7.5 9.5
```

Working with data

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

**Working with
Data**

Plotting

Tying it all
Together

Good Habits

- › Input
- › Output

Input

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

```
1 setwd("data")
2 list.files()
3 dat<-read.csv("dataset.csv")
4
5 library("foreign")
6 dat<-read.dta("dataset.dta")
```

› Useful options:

- › `sep`
- › `na.strings`

› If you have a different format, Google it.

Simulate some fake data

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › We will now just simulate some data for demonstration

```
1 > set.seed(1023)
2 > X<-rnorm(1000,0,5)
3 > Y<-sin(5*X)*exp(abs(X))+rnorm(1000)
4 > dat<-data.frame(X,Y)
5
6 #Simple Linear Regression
7 > dat.lm<-lm(Y~X,data=dat)
8
9 > summary(dat.lm)
```

- › Look at `?summary.lm`

Output

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Multiple options: `xtable`, `apsrtable`, `stargazer`
- › Can be used with, for instance, `booktabs` for prettier \LaTeX

```
1 > install.packages("xtable")
2 > library(xtable)
3 > xtable(dat.lm) #shown below
4
5 #This is the same as
6 > xtable(summary(dat.lm)$coefficients)
7
8 #Look at the following:
9 > summary(dat.lm)$coefficients
10 > class(summary(dat.lm)$coefficients)
11 [1] "matrix"
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.8851	4073.7725	0.00	0.9985
X	2989.0614	828.6452	3.61	0.0003

Plotting

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Plotting is just all about coordinate pairs.
- › `plot(x, y)` plots the pairs of `x` and `y` in order
- › Notable options:
 - › `type` – determines whether you plot points, lines, or whatnot
 - › `pch` – determines plotting character
 - › `xlim` – x limits of the plot (similar for Y)
 - › `xlab` – label on X-axis
 - › `main` – main title
 - › `col` – color
- › A huge number of other options. Read the manual.
- › Some objects will use a special plotting function of their own.
- › Example `plot(dat.lm)`

Tying it all together

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

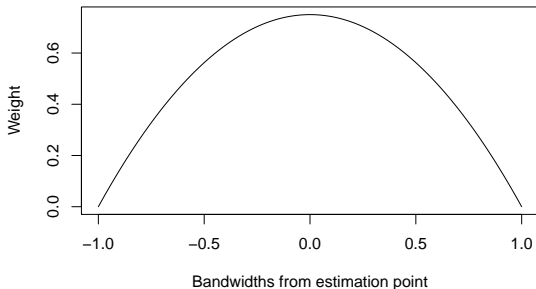
Plotting

**Tying it all
Together**

Good Habits

- › We will now create a function for local linear regression.

```
1 pdf("locallin.pdf", height=4, width=6)
2 plot(
3   seq(-1, 1, .01),
4   3/4*(1-seq(-1, 1, .01)^2),
5   type="l",
6   xlab="Bandwidths from estimation point",
7   ylab="Weight"
8 )
9 dev.off()
```



Tying it all together

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › W is an $n \times p$ diagonal weighting matrix.
- › Diagonal entries are $\frac{3}{4} \cdot (1 - d^2) \cdot 1_{\{|d| \leq 1\}}$ where $d = \frac{(X-c)}{\text{bandwidth}}$
- › $\hat{\beta}_c = (X'WX)^{-1}X'WY$
- › Covariance matrix: $s^2(X'WX)^{-1}$

```
1 loc.lin<-function(Y,X,c=0,bw=sd(X)/2){
2   d<-(X-c)/bw
3   W<-3/4*(1-d^2)*(abs(d)<1)
4   W<-diag(W)
5   X<-cbind(1,d)
6   b<-solve(t(X)%*%W%*%X)%*%t(X)%*%W%*%Y
7   sigma<-t(Y-X%*%b)%*%W%*%(Y-X%*%b)/(sum(diag(W)>0)-2)
8   sigma<-solve(t(X)%*%W%*%X)*c(sigma)
9   return(c(est=b[1],se=sqrt(diag(sigma))[1]))
10 }
```

Testing it out

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › We generate some random data, estimate our regressions at various points, and then plot it all.
- › $X \sim N(0,5)$
- › $Y_i = \sin(5 \cdot X_i) \cdot e^{|X_i|} + \epsilon_i, i \in 1, \dots, n$
- › $\epsilon_i \sim N(0,1)$

```
1 X<-rnorm(1000,0,5)
2 Y<-sin(5*X)*exp(abs(X))+rnorm(1000)
3 loc.lin(Y,X,c=2,bw=.25)
4
5 X.est<-seq(0,5,.1)
6 dat.llm<-sapply(
7   X.est,
8   function(x)loc.lin(Y,X,c=x,bw=.25)
9 )
10
11 plot(X,Y,xlim=c(0,5),ylim=c(-50,50),pch=20)
12 lines(X.est,dat.llm[1,],col="red")
13 lines(X.est,dat.llm[1,]+1.96*dat.llm[2,],col="pink")
14 lines(X.est,dat.llm[1,]-1.96*dat.llm[2,],col="pink")
```

Tying it all Together

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

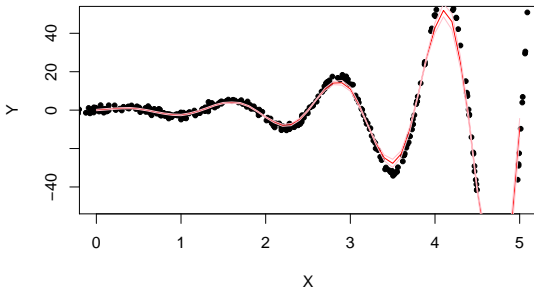
Control Flow

Working with
Data

Plotting

**Tying it all
Together**

Good Habits



```
1 pdf("locallin2.pdf",height=4,width=6)
2 plot(X,Y,xlim=c(0,5),ylim=c(-50,50),pch=20)
3 lines(X.est,df[1,],col="red")
4 lines(X.est,df[1,]+1.96*df[2,],col="pink")
5 lines(X.est,df[1,]-1.96*df[2,],col="pink")
6 dev.off()
```

Conventions and Good Habits

NYU R
Workshop

Drew Dimmery

Outline

Why R?

Where to start

Where to get help

Misc. Quirks

Using third party
code

Data Types

Data Structures

Control Flow

Working with
Data

Plotting

Tying it all
Together

Good Habits

- › Comment your code. You won't remember what you did in 2 months.
- › Make your variables mean something. `tmp` is a horrible name. Don't use it.
- › Always be making functions.
- › Be consistent in style – brackets on the same line as `for`
- › Indent your code – but not too much.
- › If you're just clever enough to write it, you aren't clever enough to debug it.
- › Fiddle with things. That's how you learn.
- › Save your work (and consider offsite backups via Dropbox/git/etc)
- › Break rules when it makes for cleaner, easier to understand code.