## Data Science: Lets All Be Ninjas

Matthew J. Denny

mzd5530@psu.edu - www.mjdenny.com - @MatthewJDenny

 $www.mjdenny.com/ICPSR\_Data\_Science\_2015.html$ 

July 27, 2015



・ロト ・ 同ト ・ ヨト ・ ヨト

Sac

L SCIENCE

NSTAT

- 1. Means lots of things...
- 2. Data: Collection, curation, exploration, prediction.

うして ふゆう ふほう ふほう ふしつ

- 3. Rewards programming skills.
- 4. Yay for outside options!

### Skills

- 1. Data management.
  - ▶ Funky data, many sources, replicability.
- 2. Working at large scale and high speed.
  - ▶ Approximate methods, HPC, resource management.

(日) (日) (日) (日) (日) (日) (日)

- 3. Description and interpretation.
  - ▶ Visualization, summarization, causal inference.
- 4. Collaboration and open source.
  - ► Fluency, Extensibility, distributability.

## What We Will Cover

- 1. Scientific Programming in R.
- 2. Remote Access cluster resources.
- 3. Bash, RStudio, Github.
- 4. High Performance Computing parallelization and performant programming.

- 5. Big Data Memory management.
- 6. Hardware.
- 7. R and C++.
- 8. Web Scraping.
- 9. Package Development.

## 1. General Programming in R.

- 2. Super Friends!
- 3. Super Tools!

# 1. General Programming in R.

#### Motivation – The Gardener

- ▶ How many plants to water?
- Which plants to water?



## > for( ) and while( ) loops.

- ▶ if( ) statements.
- ▶ Nested loops



#### Preliminaries

```
# create a vector
my_vector <- c(1:10)
print(my_vector)
```

# get the length of the vector length(my\_vector)

▲ロト ▲周ト ▲ヨト ▲ヨト 三日 - のく⊙

# comparison operators
5 < 6
5 > 6
5 == 5
5 != 6
5 <= 5</pre>

## The for( ) Loop



▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

```
▶ Do something N times.
```

```
my_vector <- c(20:30)
for(i in 1:length(my_vector)){
    my_vector[i] <- sqrt(my_vector[i])
}</pre>
```

my\_vector

[1] 4.472136 4.582576 4.690416 4.795832 4.898979
[6] 5.000000 5.099020 5.196152 5.291503 5.385165
[11] 5.477226

#### The while( ) Loop



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ = のへぐ

- Do something until a condition is met.
- Useful if you do not know the number of iterations ahead of time.

```
my_vector <- c(20:30)
counter <- 1
while(counter <= length(my_vector)){
    my_vector[counter] <- sqrt(my_vector[counter])
    counter <- counter + 1
}</pre>
```

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

my\_vector



▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 臣 のへで

- Do something if some condition is met.
  Can be built into a loop
- ► Can be built into a loop.

```
my_vector <- c(20:30)</pre>
```

```
for(i in 1:length(my_vector)){
    if(my_vector[i] == 25){
        print("The square root of 25 is 5!")
    }
}
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三回 ● ○へ⊙

 Do something if some condition is not met.

うして ふゆ く は く は く む く し く

```
my_vector <- c(20:30)
for(i in 1:length(my_vector)){
    if(my_vector[i] == 25){
        print("I am 25!")
    }else{
        print("I am not 25!")
    }
}</pre>
```

Traversing A Matrix

#### > matrix(1:25,5,5)[,1] [,2] [,3] [,4] [,5] [1,] [2,] - 3 [3,] [4,] [5,]

## Nested Loops

## Can loop over entries in higher dimensional data structures.

```
my_matrix <- matrix(1:100,ncol=10,nrow=10)</pre>
```

```
for(i in 1:length(my_matrix[,1])){
    for(j in 1:length(my_matrix[1,])){
        if(my_matrix[i,j] %% 2 == 0){
            my_matrix[i,j] <- 0
        }
    }
}</pre>
```

うして ふゆ く は く は く む く し く

my\_matrix

- Flexible, can store any kind of data including another list.
- ▶ Good for keeping results together.

```
# Create an empty list
my_list <- vector("list", length = 10)</pre>
```

```
# Create a list from objects
my_list <- list(10, "dog",c(1:10))</pre>
```

# Add a sublist to a list
my\_list <- append(my\_list, list(list(27,14,"cat")))</pre>

#### List Contents

> my\_list [[1]] [1] 10 [[2]] [1] "dog" [[3]] [1] 1 2 3 4 5 6 7 8 9 10 [[4]] [[4]][[1]] [1] 27 [[4]][[2]] [1] 14 [[4]][[3]] [1] "cat"

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ = のへぐ

# 2. Super Friends!

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ = のへぐ

#### **Overview**

- **bash** is a command line terminal.
  - Linux and OS X (Cygwin for Windows).
- ► VPN: access campus network.
- ► **ssh/PuTTY** (Windows) for remote access.

・ロト ・日 ・ モ ・ モ ・ モ ・ シュマ

▶ ftp: for file transfer.

- ► cd: change the current directory.
- ▶ 1s prints contents of current directory.
- ▶ edit/vi/emacs opens a text editor.

0587377979:Desktop matthewjdenny\$ cd RA\_Projects/ 0587377979:RA\_Projects matthewjdenny\$ ls Example Jerry Epstein Jerry Friedman 0587377979:RA\_Projects matthewjdenny\$ cd Example/ 0587377979:Example matthewjdenny\$ ls Example.R 0587377979:Example matthewjdenny\$ edit Example.R

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

- python opens python console.
- **R** opens standard R console.
- cd .. moves us back up a level in the directory structure.

0587377979:Example matthewjdenny\$ cd .. 0587377979:RA\_Projects matthewjdenny\$ cd .. 0587377979:Desktop matthewjdenny\$

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

## Using VPN

- ▶ Provided by university/organization.
- ► For login to local campus resources.
- ▶ Routes all trafic through campus servers



## Using SSH

- ▶ Must have account on remote machine.
- ▶ ssh username@ipaddress
  - ▶ static ip: 128.114.64.8
  - ▶ dynamic: somedomain.dyndns.com
- ▶ prompt to enter password



## FTP Using FileZilla

000	File7illa				
Host: Username: Passwo	rd: Port: Quickconnect				
	^				
Local site: /Users/matthewjdenny/Desktop/	Remote site:				
- Applications					
Desktop					
Documents					
Downloads					
Dropbox					
Google Drive					
Library					
Movies					
Filename A Filesize Filetype	Filename A Filesize Filetype Last modified				
Conferences Direct					
CospRep Direct					
📁 Econ 397 Fall 11 Direct	Not connected to any server				
49 files and 36 directories. Total size: 74,522,858 bytes	Not connected.				
Server/Local file Direction Remote file	Size Priority Status				
Oueued files Failed transfers Successful transfers					
	rza Oueue: empty				
	and Queue, empty				

1. Connect to campus network using VPN

2. **ssh** into remote computing resource

3. Transfer files to/from machine using ftp

4. Navigate directories using **bash** and run analysis in **R** or **Python**.

# 3. Super Tools!

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ = のへぐ

- ▶ Version Control
  - Git or Subversion.
  - Make an account on Github.
- RStudio
  - Lots of setup options.
  - ▶ Version control integration.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ● ◎ ● ●

 System for tracking changes to documents.

うして ふゆ く は く は く む く し く

- Resolving conflicts.
- Reverting changes.
- ► Sharing your work.
- Highly motivating!

► Start by going to: github.com – make an account.

・ロト ・日 ・ モ ・ モ ・ モ ・ シュマ

- ► Then windows.github.com or mac.github.com.
- ► Check out a tutorial.
- ▶ RStudio also has Git integration.

## Git Tracks Changes

00	🚞 matthewjdenny/GERGM		R <sub>M</sub>	
+•	t master ▼     Changes History	Bran	ches	וֹז 🗘 Sync
Filter Repositories	126 commits	mai	n/GERC	IM-package.Rd
GITHUB	3 days ago by matthewjdenny			Bhamidi, Skyler Cranmer, and Bruce Desmarais
Congressional	Update transformation mechanics and tak 2 3 days ago by matthewjdenny			(2015). "Stochastic Weighted Graphs: Flexible Model Specification and Simulation". http://
ContentStructure	pass in parameter to specify transformatio	specify transformatio 29 29 }	arxiv.org/abs/1505.04015 }	
Cross_County	3 days ago by matthewjdenny	30	30	
Email_Preproc	added a caveat 27 days ago by matthewjdenny	31 32		- network <- matrix(runif(100),10,10) - diag(network) <- 0
GERGM_Devel	updated with preparing covaraiates example		31	<pre>+ ####################################</pre>
📮 ISSR_Data_Sci	27 days ago by matthewjdenny	_	32	+ # Preparing an unbounded network without
📮 REmail	updated documentation 27 days ago by matthewijdenny		33	<pre>covariates for gergm estimation # + net &lt;- matrix(rnorm(100,0,20),10,10)</pre>
	removed unnecessary calls 27 days ago by matthewjdenny		34 35	<pre>+ colnames(net) &lt;- rownames(net) &lt;- letters[1:10] + network &lt;- Prenare Network and Covariates(raw network =</pre>
	fixed dimension error in lambda estimation 27 days ago by matthewjdenny	3	36	net,
	updated documentation			<pre>normalization_type = "division")</pre>
27 days ago by matthewjdenny	27 days ago by matthewidenny	33	37	formula <- "network ~ recip + edgeweight"
	fixed making diagonal of network zero 28 days ago by matthewjdenny	35	39	@@ -37,7 +41,7 @@ test <- gergm(formula,

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへの

## Motivation



( = ) (

Sac

- Make sure you have the latest version of R.
- ► Then go to: rstudio.com free academic download.
- ► Check out a tutorial.
- ▶ Spend time getting to know the interface.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ● ◎ ● ●

# www.mjdenny.com/ Data\_Science\_Tools.html

うして ふゆ く は く は く む く し く