

Introduction to Data Science
Jeffrey S. Saltz/Jeffrey M. Stanton

Week 1 Introduction

Copyright 2021: Jeffrey Saltz and Jeffrey Stanton; please do not upload.

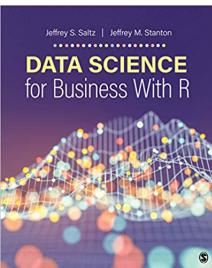
SCHOOL OF INFORMATION STUDIES
school.syr.edu SYRACUSE UNIVERSITY

Outcomes For Today

- Discuss the goals for the course
- Be able to describe what “data science” is
 - Skills needed
 - Data insights
- Get started with R

A Quick Look At the Class

1. What is Data Science
2. Quant overview: Descriptive & Inferential Stats
3. R Skills: Using R functions & libraries
4. Connecting to external data sources
5. Visualization (ggplot)
6. Creating data-driven maps
7. Text Processing/Analysis
8. Linear Modeling
9. Unsupervised Data Mining
10. Supervised Data Mining



Our Backgrounds

Jeffrey S. Saltz

- Computer Science/Information Science/Business Degrees
- Started work life as a programmer
- Led visualization and high-performance computing teams
- Invested in tech companies—venture capital
- Ran credit card risk platform—credit and fraud
- Now do research on “how to do data science”

Jeffrey M. Stanton

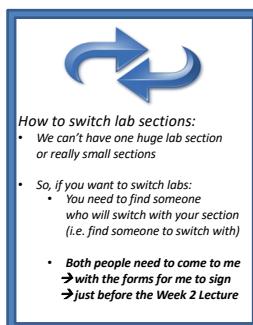
- Computer Science/Organizational Psychology Degrees
- Started work life as a programmer
- Led embedded operating system and signal processing teams
- Worked at venture capital funded startup firms
- Developed text processing and digital audio processing platforms
- Now do research on workers and technology in organizations

Typical Structure of Lecture Classes

- 1. Review & Introduce Key Concepts**
- 2. Data Science Case Study / Discussion**
- 3. Hands-On Activity (Using R)**

Each Week

- Hands on: Goal is to enable you to **practice** data science
- Bring your laptop to every lab session with access to R and R-Studio
- Read the chapter before the class!!!



Labs/Homework

- Bring your laptop to every lab
- Lab ==> Homework: The lab focuses on developing the techniques and code needed for successful completion of the homework
- Your lab instructor is your coach for helping you learn R, master the homework, and execute the term project

Grades

Lab is designed for you to practice the necessary skills in carrying out data processing, analysis, and management tasks

Homework is designed for you to practice the necessary skills in carrying out data processing, analysis, and management tasks

Participation includes professionalism in class, and actually attending class...

Quiz are weekly quick multiple choice tests

Test is designed to evaluate your mastery of concepts, methods, and tools in data analysis and management.

Project: For your project you work on a data set provided to you, transform the data as needed, and provide a written analysis with visualizations.

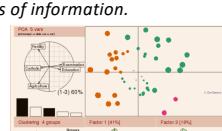
Data Science: Many Skills



School of Information Studies
SYRACUSE UNIVERSITY

Data Science: Definition

Data science: *An emerging area of work concerned with the collection, preparation, analysis, visualization, management, and preservation of large collections of information.*



Connects strongly with areas such as databases, statistics, and computer science, but many different skills—including nonmathematical skills—are needed.

Core Goal of Data Science: Creating Actionable Insight

What Is Data Science?

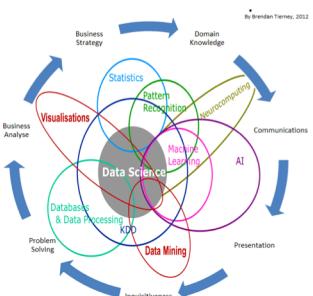
Identify which article belongs to which field:

1. Zhang, N. F. (2006). The uncertainty associated with the weighted mean of measurement data. *Metrologia*, 43(3), 195.
2. Janssen, M., Konopnicki, D., Snowdon, J. L., & Ojo, A. (2017). Driving public sector innovation using big and open linked data (BOLD). *Information Systems Frontiers*, 19(2), 189-195.
3. Hillar, C. J., & Lim, L. H. (2013). Most tensor problems are NP-hard. *Journal of the ACM (JACM)*, 60(6), 45.

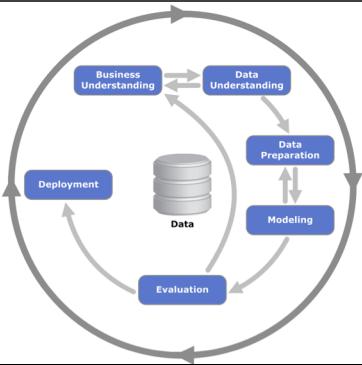
The diagram consists of four overlapping circles arranged in a diamond shape. The top-left circle is yellow and labeled 'Computer science'. The bottom-left circle is blue and labeled 'Domain knowledge'. The bottom-right circle is green and labeled 'Math and statistics'. The top-right circle is light blue and labeled 'Data science'.

Image credit:
Quora.com

Data Science is Multidisciplinary



Doing Data Science: CRISP-DM



Doing Data Science: OSEMN

The pipeline:

- O — Obtaining the data
- S — Scrubbing / Cleaning
- E — Exploring / Visualizing (patterns / trends)
- M — Modeling the data (doing predictions)
- N — Interpreting the analysis

→ Understanding the business problem?
→ This is iterative (need to keep looping)

<http://www.dataists.com/2010/09/a-taxonomy-of-data-science/>

The Data Science Process

- Domain analysis/decomposition
- Identification of subject matter expert(s)
- Question/interview/observation process

Via anomalies, we learn:

- Is there an error in the system definition
- Is there an error in the system implementation
- Is there an **interesting exception** that teaches us something

Via stories, we learn:

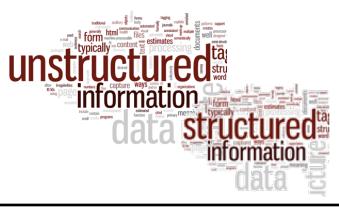
- What people do
- How they do it
- Information produced and consumed
- Process and information touch points
- Decisions made
- Key Challenges

Skills Through the Life Cycle



Structured & Unstructured Data

- What is the difference between structured and unstructured data?



Questions

- Do data scientists mainly deal with structured data or unstructured data?
 - Which kind of data is there more of?



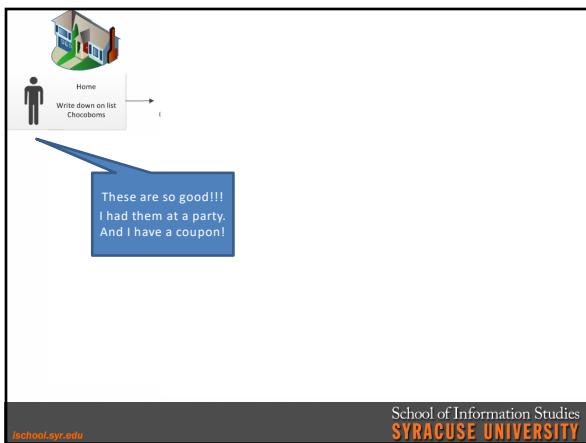
Data Science: First Example

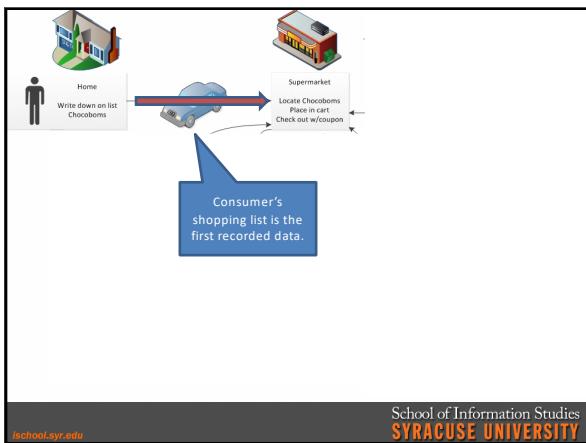


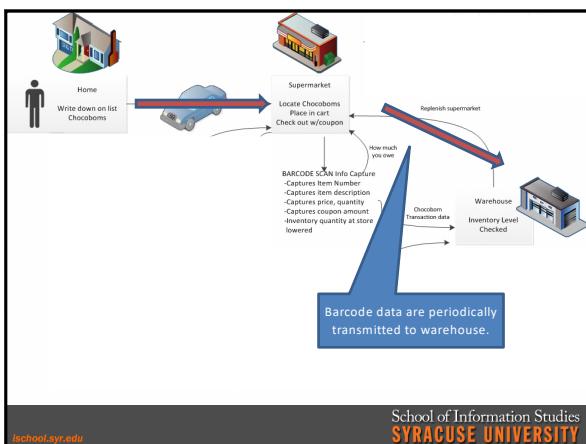
Data Science Example

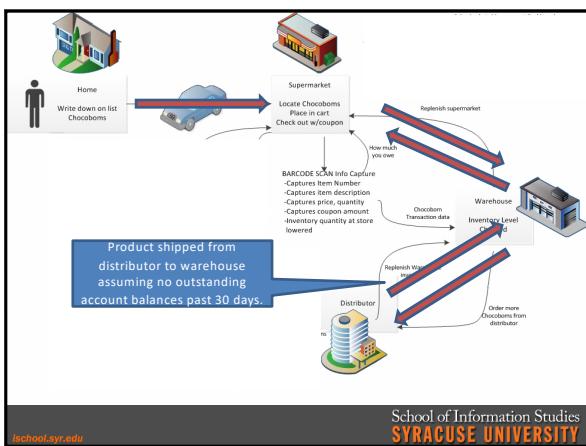
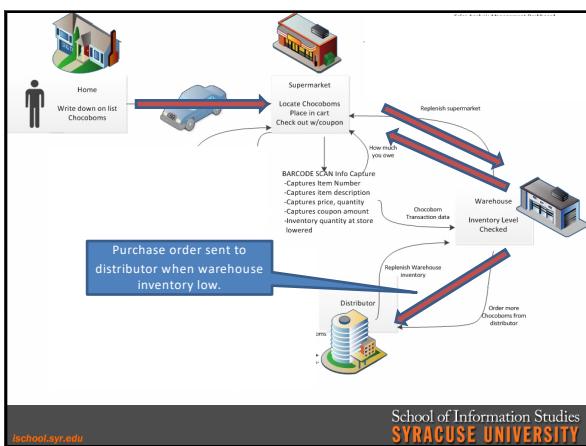
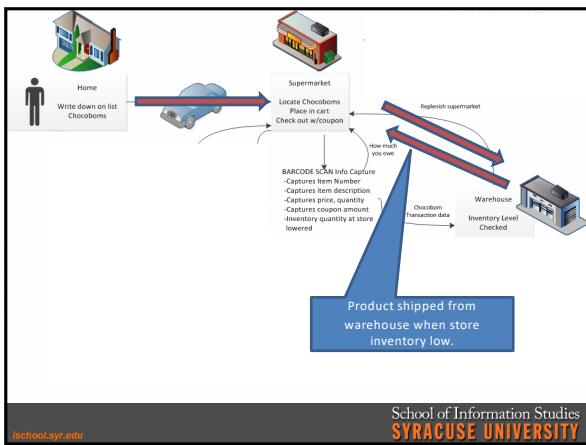


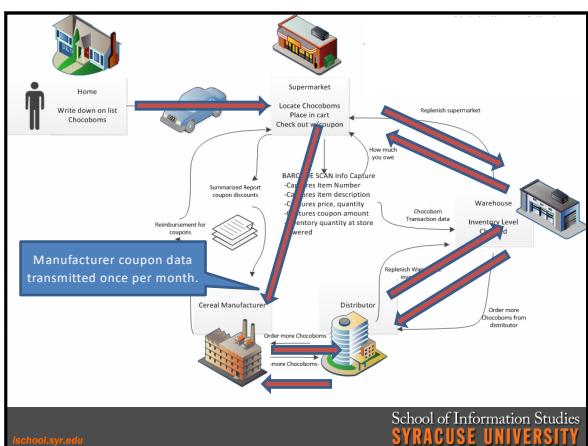
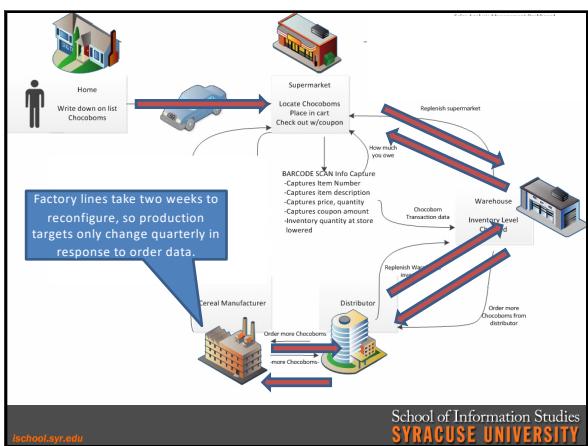
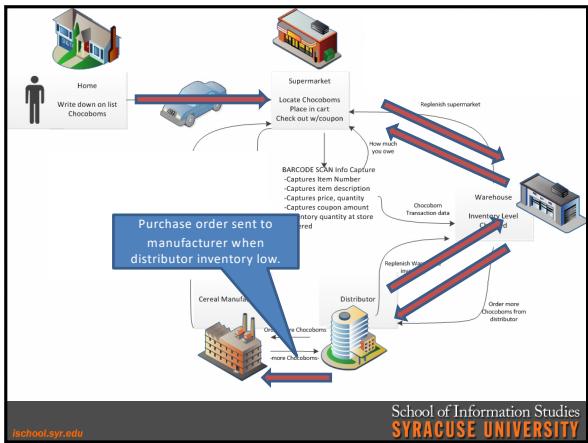
The problem as described by client:
Many variably priced products, like this one,
have unanticipated demand spikes that create severe restocking delays.

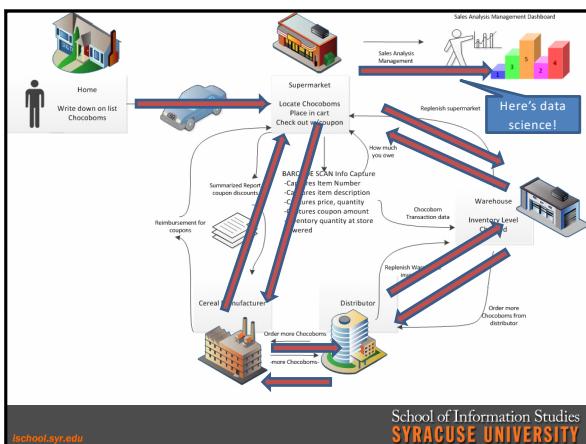
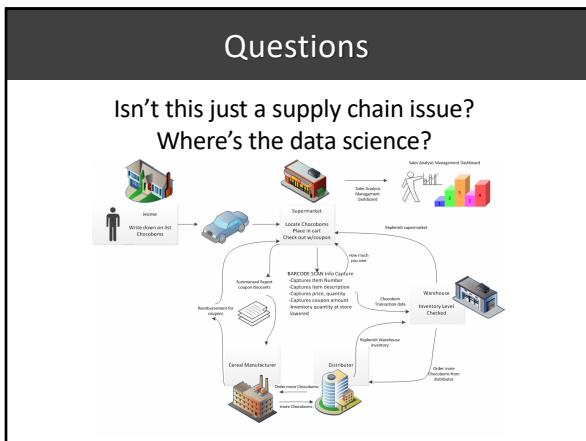
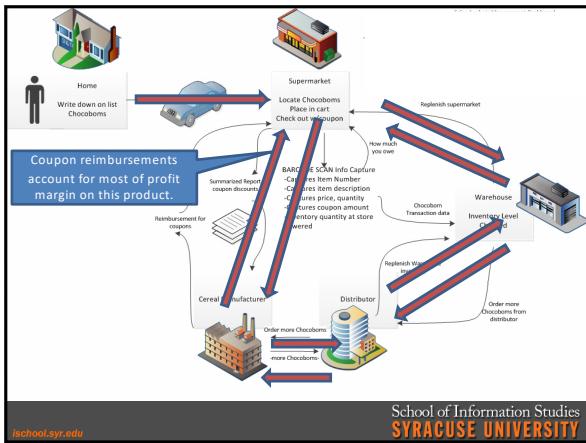






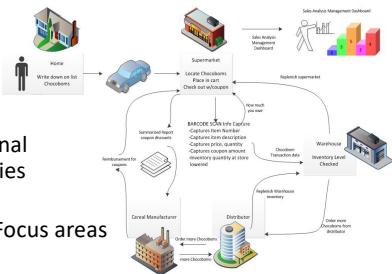






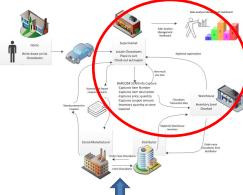
Thinking About the Problem Domain

- Define Scope
- Set Internal Boundaries
- Choose Focus areas
- Remain Aware of Larger Project Context



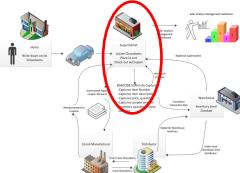
Set Overall Scope

- Problem domain
 - Define scope
 - Set internal boundaries
 - Choose focus areas
 - Context
- Avoid “boiling the ocean”



Define Internal Boundaries

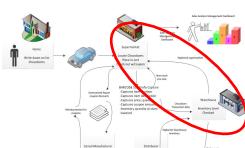
- Problem domain
 - Define Scope
 - **Set Internal Boundaries**
 - Focus area
 - Context
- Avoid “boiling the ocean”
- Decompose domain into manageable modules



Set Focus Area(s)

- Problem domain
 - Scope
 - Internal Boundaries
 - Focus area(s)**
 - Context

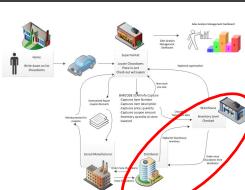
Avoid “boiling the ocean”
 Decompose domain into manageable modules
 Solve (at least) one important problem well



Enhance Contextual Awareness

- Problem domain
 - Scope
 - Boundaries
 - Focus area
 - Context**

Avoid “boiling the ocean”
 Decompose domain into manageable modules
 Solve (at least) one important problem well
 All (good) data science includes change processes



The Data Science Process

- Domain analysis/decomposition
- Identification of subject matter expert(s)
- Question/interview/observation process
 - Risks and uncertainty
 - Stories
 - Anomalies

Via anomalies, we learn:

- Is there an error in the system definition
- Is there an error in the system implementation
- Is there an **interesting exception** that teaches us something

Via stories, we learn:

- What people do
- How they do it
- Information produced and consumed
- Who owns what risks
- Process and information touch points
- Decisions made
- Challenges associated with all of the above

Key Takeaways

- In contrast to other fields, data science is at its heart concerned with **communicating**.
- **Domain identification** and **understanding** is essential such that we focus on the most pressing *and* tractable problem or opportunity.
- **Identify SMEs** and engage them in “storytelling” about their activities such that you can get a sense of who, what, why they do what they do to also include surfacing anomalies and risk situations.

Getting Started With R



school.syr.edu

School of Information Studies
SYRACUSE UNIVERSITY

Data Science: Getting Started with R

- R is an **open-source** software program
 - Developed by volunteers as a service to the community of scientists, researchers, and data analysts who use it.
- R is **free** to download and use.
- **Lots of advice is available** online to help users learn R, which is good because it is a powerful and complex program.
- R is a full-featured programming language dedicated to data.

Data Science: Getting Started with R

```
R version 3.0.2 (2013-09-25) -- "Frisbee Sailing"
Copyright (C) 2013 The R Foundation for Statistical Computing
Platform: i386-w64-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

Attempting to load the environment 'package:stats'
[Previously saved workspace restored]
> |
```

Data Science: Getting Started with R

```
R version 3.5.1 (2018-07-02) -- "Feather Spray"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
> |
```

Getting Started with R

- Download R installer from:
<https://cran.r-project.org>
- Available for Windows, Mac, and Linux
- Run installer
- Start up R
- Type commands
- End R q() at the command prompt
- Respond ‘yes’ to save your work space

→ Or use [rstudio.cloud](#)

An Introduction to Data in R

"Atomic" Types in R:

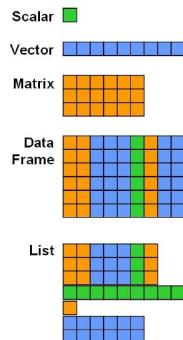
Logical – Either TRUE or FALSE; TRUE generally coded as 1 and FALSE as 0; TRUE and FALSE are reserved keywords

Integer – A number with nothing after the decimal point; sometimes displayed with an "L" after it as in 95L

Double – A number that supports digits right of the decimal point

Character – Text data encoded using any of a variety of encoding systems

Rarely used types: Complex and Raw



An Introduction to Data: Type / Mode

List of integers

43, 42, 12, 8, 5

Each integer represents the age of a family member.

Integer list is all the same "type/mode."

R refers to an atomic list as a "vector."

R code to create this vector looks like

c(43, 42, 12, 8, 5)

myFamilyAges <- c(43, 42, 12, 8, 5)

Named vector

The Initial Console View

```
R version 3.0.2 (2013-09-25) -- "Frisbee Sailing"
Copyright (C) 2013 The R Foundation for Statistical Computing
Platform: i386-w64-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

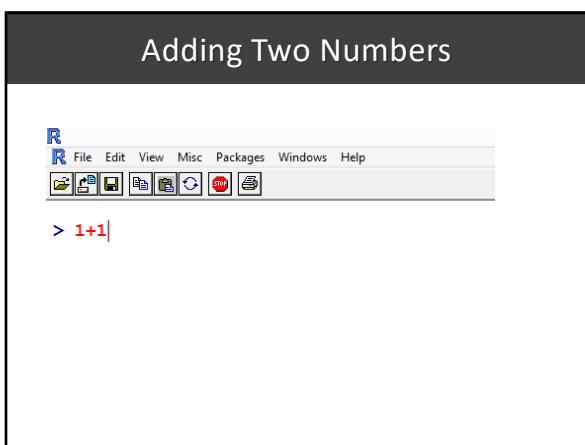
Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

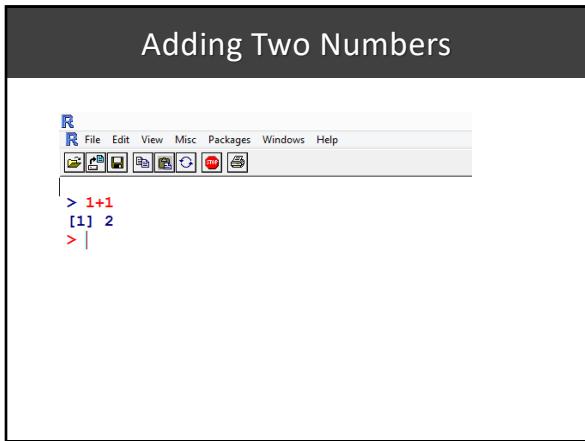
Attempting to load the environment 'package:stats'
[Previously saved workspace restored]
```

Adding Two Numbers



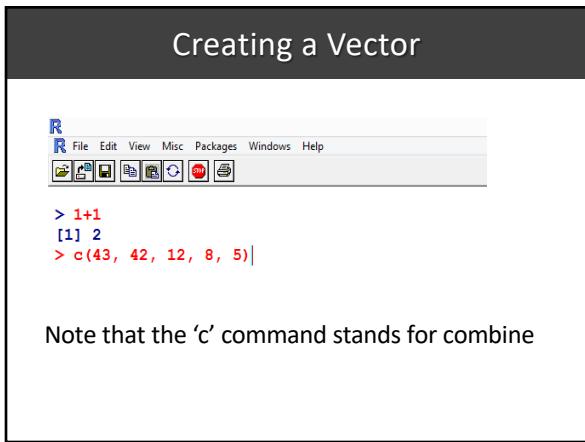
The screenshot shows the RStudio interface with a dark header bar. The title bar says "Adding Two Numbers". Below it is a menu bar with "File", "Edit", "View", "Misc", "Packages", "Windows", and "Help". A toolbar with various icons follows. The main workspace shows the command `> 1+1` in red, with the result `[1] 2` displayed below it.

Adding Two Numbers



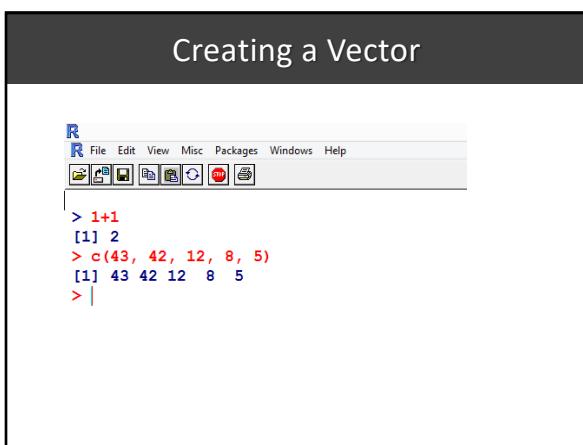
The screenshot shows the RStudio interface with a dark header bar. The title bar says "Adding Two Numbers". Below it is a menu bar with "File", "Edit", "View", "Misc", "Packages", "Windows", and "Help". A toolbar with various icons follows. The main workspace shows the command `> 1+1` in red, followed by the result `[1] 2`, and then an empty line prompt `> |`.

Creating a Vector



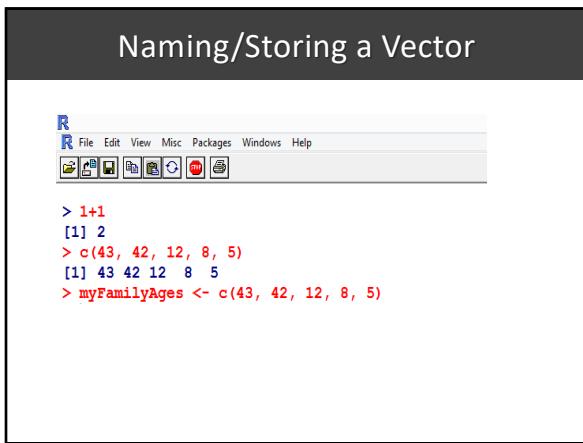
The screenshot shows the RStudio interface with a dark header bar. The title bar says "Creating a Vector". Below it is a menu bar with "File", "Edit", "View", "Misc", "Packages", "Windows", and "Help". A toolbar with various icons follows. The main workspace shows the command `> 1+1` in red, followed by the result `[1] 2`, then `> c(43, 42, 12, 8, 5)` in red, and finally an empty line prompt `> |`. Below the workspace, a note reads: "Note that the 'c' command stands for combine".

Creating a Vector



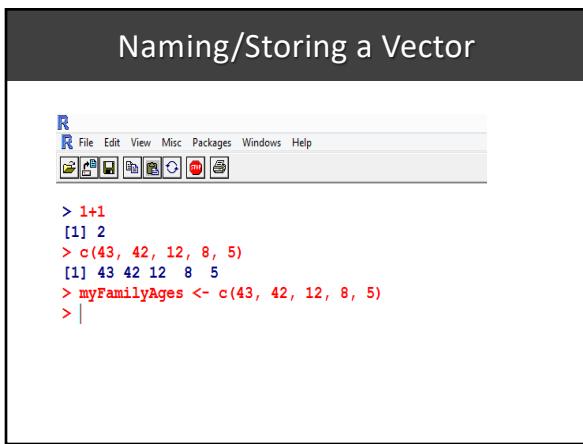
```
R
File Edit View Misc Packages Windows Help
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> |
```

Naming/Storing a Vector



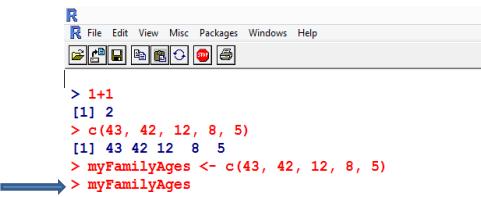
```
R
File Edit View Misc Packages Windows Help
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
```

Naming/Storing a Vector



```
R
File Edit View Misc Packages Windows Help
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> |
```

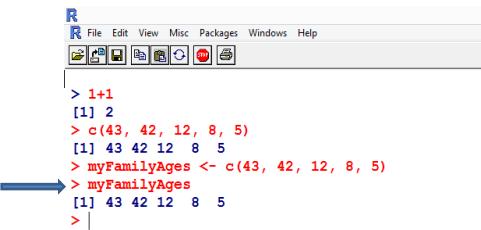
Reporting a Stored Vector



```
R
File Edit View Misc Packages Windows Help
> 1+1
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> myFamilyAges
```

A blue arrow points from the line 'myFamilyAges' to the right margin of the code window.

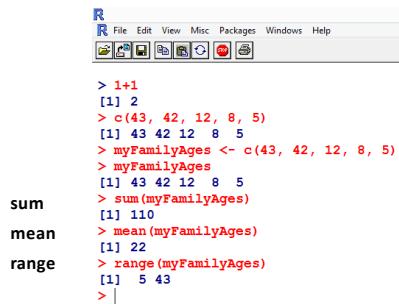
Reporting a Stored Vector



```
R
File Edit View Misc Packages Windows Help
> 1+1
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> myFamilyAges
[1] 43 42 12 8 5
> |
```

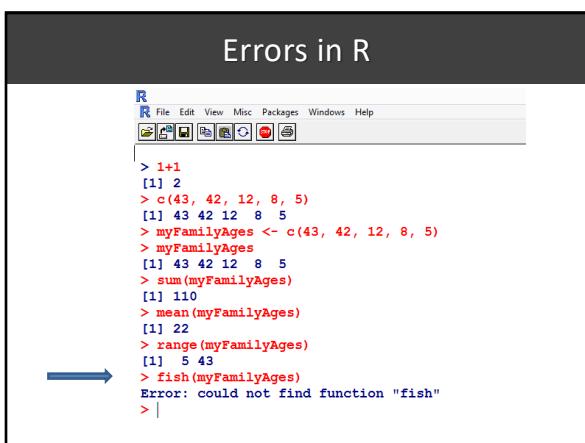
A blue arrow points from the line 'myFamilyAges' to the right margin of the code window.

Functions Operating on a Vector



```
R
File Edit View Misc Packages Windows Help
> 1+1
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> myFamilyAges
[1] 43 42 12 8 5
> sum(myFamilyAges)
[1] 110
> mean(myFamilyAges)
[1] 22
> range(myFamilyAges)
[1] 5 43
> |
```

Errors in R



The screenshot shows the RStudio interface with the title bar "Errors in R". The console window displays the following R session:

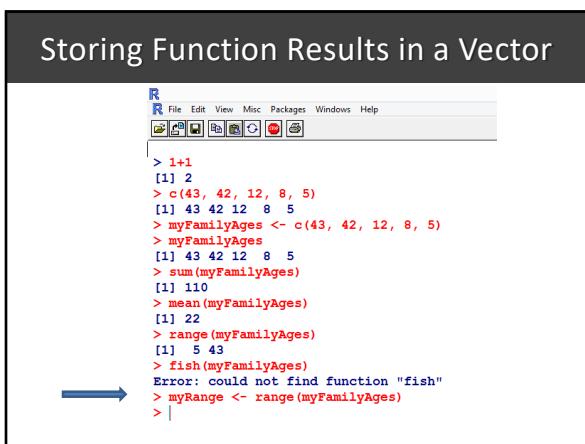
```

> 1+1
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> myFamilyAges
[1] 43 42 12 8 5
> sum(myFamilyAges)
[1] 110
> mean(myFamilyAges)
[1] 22
> range(myFamilyAges)
[1] 5 43
> fish(myFamilyAges)
Error: could not find function "fish"
> |

```

A blue arrow points to the last line of the session where the function "fish" is called.

Storing Function Results in a Vector



The screenshot shows the RStudio interface with the title bar "Storing Function Results in a Vector". The console window displays the same R session as the first block, but includes an additional line at the end:

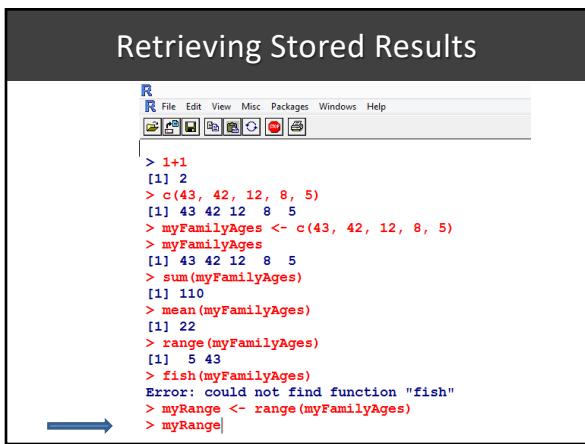
```

> myRange <- range(myFamilyAges)
> |

```

A blue arrow points to the last line of the session where the function "range" is stored in the variable "myRange".

Retrieving Stored Results



The screenshot shows the RStudio interface with the title bar "Retrieving Stored Results". The console window displays the same R session as the previous blocks, but includes an additional line at the end:

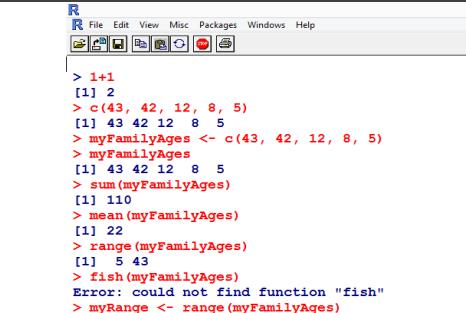
```

> myRange
[1] 5 43
> |

```

A blue arrow points to the last line of the session where the variable "myRange" is retrieved.

Results



```
R
R File Edit View Misc Packages Windows Help
[1] 1+1
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> myFamilyAges
[1] 43 42 12 8 5
> sum(myFamilyAges)
[1] 110
> mean(myFamilyAges)
[1] 22
> range(myFamilyAges)
[1] 5 43
> fish(myFamilyAges)
Error: could not find function "fish"
> myRange <- range(myFamilyAges)
> myRange
[1] 5 43
> |
```

The screenshot shows the RStudio interface with the following details:

- Header:** Capitalization is Important
- Toolbar:** File, Edit, View, Misc, Packages, Windows, Help
- Environment pane:** RG
- Code pane:**

```
> > 1+  
[1] 2  
> c(43, 42, 12, 8, 5)  
[1] 43 42 12 8 5  
> myFamilyAges <- c(43, 42, 12, 8, 5)  
> myFamilyAges  
[1] 43 42 12 8 5  
> sum(myFamilyAges)  
[1] 110  
> mean(myFamilyAges)  
[1] 22  
> range(myFamilyAges)  
[1] 5 43  
> fish(myFamilyAges)  
Error: could not find function "fish"  
> myRange <- range(myFamilyAges)  
> myRange  
[1] 5 43  
> myrange|
```

The screenshot shows the RStudio interface with a warning message in the console:

```
> 1+1
[1] 2
> c(43, 42, 12, 8, 5)
[1] 43 42 12 8 5
> myFamilyAges <- c(43, 42, 12, 8, 5)
> myFamilyAges
[1] 43 42 12 8 5
> sum(myFamilyAges)
[1] 110
> mean(myFamilyAges)
[1] 22
> range(myFamilyAges)
[1] 8 43
> fish(myFamilyAges)
Error: could not find function "fish"
> myRange <- range(myFamilyAges)
> myRange
[1] 5 43
> myrange
Error: object 'myrange' not found
> |
```

A blue arrow points from the left towards the warning message.
