DATA SOCIETY®

*"If you can't explain it simply, you don't understand it well enough."*

- Albert Einstein

# Course syllabus

1. What is data science?

2. Manipulating data in R

3. Visualization in R



DATA SOCIETY © 2017                    1

# Setting expectations

Data science takes dedication! You will need to:

1. Take this course ☺
2. Practice
3. Review class material on your own
4. Practice
5. Complete exercises outside of class
6. Practice
7. Share and read latest news

# Outline

- What is data science?
- A data scientist's approach
- Introduction to R
  - Calculations in R
  - Reading data into R
  - Manipulating data in R
- Visualization in R
  - Basic plotting

# How is data being used?

| Retail | Finance | Marketing | Real estate | Cool |
|--------|---------|-----------|-------------|------|
| Target: | Kabbage: | Netflix: | Zillow: | Andrew Ng: |
| The store knows you're pregnant based on what you buy | Makes lending decisions based on Amazon product reviews, etc. | What movie should you watch? | Calculates Zestimate (value of your home) | Machine learning techniques recognize cat faces online using pictures and videos |

# What is "Big Data"?

- Big Data is **large** volumes of information
  - Moving
  - Storing
  - Manipulating
  - Accessing

- It is <u>not</u>:
  - Analysis or insights

*That's why you're in this class!*

# What is data science?

- Data science applies the scientific method to analyzing data

- It lies at the intersection of several disciplines

- It draws on domain specific knowledge that makes the analysis of Big Data possible

# Who is a data scientist?

- An analyst who can:
  1. **Pose** the right question

  2. **Wrangle** the data (gather, clean, and sample data to get a suitable data set)

  3. **Manage** the data for easy access by the organization

  4. **Explore** the data to generate a hypothesis

  5. **Make predictions** using statistical methods such as regression and classification

  6. **Communicate** the results using visualizations, presentations, and products

# Levels of expertise

**Data analyst**

**Data modeler**

**Data scientist**



- **Wrangles** the data

- **Manages** the data

- **Creates** basic analyses and visualizations

- **Models** to answer specific questions

- **Understands the data**, its source and structure

- **Asks** the right questions

- **Looks** for patterns in data

- **Interprets** results critically

# Data science job market

Somewhat important ✔    Very important ✔

| | A non-data-driven company | The business is just starting to collect data | Data is the product of the company | Company uses data to make decisions |
|---|---|---|---|---|
| Basic tools | ✔ | ✔ | ✔ | ✔ |
| Software engineering | | ✔ | ✔ | ✔ |
| Statistics | ✔ | ✔ | ✔ | ✔ |
| Machine learning | | | ✔ | ✔ |
| Data processing | | ✔ | ✔ | ✔ |
| Data visualization and communication | ✔ | ✔ | ✔ | ✔ |
| Thinking like a data scientist | ✔ | ✔ | ✔ | ✔ |

# Who hires data scientists?

Source: datasciencecentral.com

# How much do data scientists make?

- According to a Burtch Works 2014 data science job market survey:

  *"Data scientists earn a median salary that can be up to 40% higher than predictive analytics professionals at the same job level"*

- The graphic on the right provides detail on median salaries by experience level



Source: http://www.burtchworks.com/big-data-analyst-salary/big-data-career-tips/

# Outline

- What is data science?
- A data scientist's approach
- Introduction to R
  - Calculations in R
  - Reading data into R
  - Manipulating data in R
- Visualization in R
  - Basic plotting

# Data science control cycle

# Data science control cycle



**1** Ask
What is the problem(s) we need to solve?

**2** Research
What data do we need and how do we get it?

**6** Interpret
How can we use the conclusions in the real world?

Sanity check yourself before you...

**3** Model
Which method(s) is appropriate to use?

**5** Test
How does the model generalize to real world data?

**4** Validate
Do the model and assumptions work as expected?

# For every job there is a tool

## Data storage

- **Hadoop**
- **Spark**
- **SQL**
- **...**

## Data analysis

- **R**
- **SAS**
- **SPSS**
- **Matlab**
- **Python**
- **Google Prediction API**
- **...**

# For every job there is a tool

## Visualization

- R
- Tableau
- iVEDiX



Density plot

*All of these visualizations were created in R*

# Supervised machine learning

**Pattern discovery when inputs (x) and outputs (y) are known**

Input *x:*
Voter

Output *y:*
Political
affiliation

Examples:  Classification and regression are supervised machine learning

# Unsupervised machine learning

**The data inputs _(x)_  have no target outputs _(y)_**

Input _x_:
Voter

**?**

Output _y:_
Not given
(to be discovered)

We want to impose structure on the inputs _(x)_ to say something
meaningful about the data

# Machine vs. human

| | Machine | Human |
|---|---|---|
| **Understanding context** | | ✔ |
| **Thinking through the problem** | | ✔ |
| **Asking the right questions** | | ✔ |
| **Selecting the right tools** | | ✔ |
| **Performing calculations quickly** | ✔ | |
| **Performing repetitive tasks** | ✔ | |
| **Following pre-defined rules** | ✔ | |
| **Interpreting results** | | ✔ |

# Setting up R: overview

1. What is R?
2. Download R from the CRAN website (http://cran.us.r-project.org/)
   - R for Windows
   - R for Mac
3. Install R Studio (http://www.rstudio.com/products/rstudio/download/)
   - RStudio a brief tour
4. Running a script
   - Variables
5. Reading in a data
   - Manually
   - Through the script

# What is R?

- R is a statistical programming software
  - Has many similar features to SAS, Excel and SPSS
  - Scripting language

- It's free and open source

- Has lots of helpful pre-built functions
  - You can build your models quicker

- Easy to learn

# Install R

# R for Windows

# R for Mac

# What is RStudio?

The user interface in the R terminal
is bulky and non-intuitive

RStudio provides a better interface
for a more intuitive user experience

# Install R Studio

# RStudio

# Script and Console

- Script
  - We use scripts to generate re-usable code
  - They are like macros
  - Code is not executed here unless you press run (see next slide)
  - Good coding practice:
    - When saving a script use a name that describes what the script does i.e. counter
    - Comment your code using #

- Console
  - This is where the code executes
  - It is the actual R program (backend)
  - Anything that you type in here will be executed

# Running code from script



Type code in here

# Running code from a script



Highlight the code
you want to run

# Running code from a script



To run the code in the console
Press: Run
or
Hit: ctrl/command + enter

# Running code from a script



Code is executed in the console

# Running code from a script



Notice #a=10
is a comment and not
"executed", otherwise
a+b+d = 16 rather than 9

# Variables



Once the code is run the variables are "stored" in memory. Here we have a,b,d

# Reading in data

- Data comes in many forms
  - R can read in most forms
- The most common to work with now days are
  - `csv` files – comma separated value files
  - `tsv` files – tab separated value files
- R can also read in others (excel and json) but that is not covered here

# Manually



Tools > Import Dataset > From Text File...

# Manually

- R will automatically select options for you such as heading, sometimes it makes mistakes so be sure to double check

# Through the Script



Step 1:
Select your working directory
(This is where your data files are located and where your documents will save)

# Through the script



Step 2:
Use the read.csv command
to read in
"white_wine.csv"

# Through the script



Step 3:
Highlight and press run

# Through the script

# Outline

- What is data science?
- A data scientist's approach
- **Introduction to R**
  - Calculations in R
  - Reading data into R
  - Manipulating data in R
- Visualization in R
  - Basic plotting

# Why use R

1. De facto standard among professional statisticians

2. Comparable and often superior in power to commercial products (SAS, SPSS)

3. Available for the Windows, Mac, and Linux operating systems

4. R is a general-purpose programming language, so you can use it to automate analyses

5. Create dynamic graphics and visualization

6. Large community of users, many are prominent scientists: www.r-bloggers.com

7. Pre-made packages to run data analyses contributed by user base (over 6,500 packages)

Source: http://cran.r-project.org/web/packages/

# Uses of R

1. Can be used to analyze and visualize data

2. Can be used to write software

3. Can be used to create data products and applications

*In this course, we will focus on how to analyze and visualize data*

# Data formats R can read

- Can work with many types of data

# Companies that use R

# R vs. Excel

| | R | Excel |
|---|---|---|
| **Data capacity** | R can read files as big as several gigabytes and trillions of data points; only limitation is your RAM | Excel can't read more than 1,048,576 rows and 16,384 columns (2011 version), files over ~300 megabytes can be very slow to work with |
| **Customization** | Can create custom visualizations through code, very flexible | Drop down menus limit ability to manipulate charts and graphs |
| **Analyzing data** | Powerful, pre-built packages that speed up work flow | Less flexible built-in analytic abilities that can be augmented by macros |
| **Modeling** | Data analysis and statistical models | Complex financial and accounting models |
| **Seeing data** | Built-in spreadsheet viewer | Easy to use spreadsheet interface |
| **Usability** | Direct commands similar to Excel "if-statements" | Keyboard shortcuts and slower point-and-click functionality |

# Visualizations in R

## R

Simple customizable code: flexible



## Excel

Drag and drop: rigid

# R vs. Python

- R has more convenient statistical packages to analyze data than Python
  – More than any other software tool, over 6,500 as of April 2015

- R is easier to learn for non-programmers than Python, less code is required to perform tasks

- Python is used by many data scientists to build data products (they also tend to be computer scientists)

- Python can be easier to integrate into web applications

Source: http://cran.r-project.org/web/packages/

# RStudio overview



- Top left runs commands, called *Script*

- Top right shows summary of data and variables currently loaded, called *Environment*

- Bottom left shows results, called *Console*

- Bottom right shows help files, graphical outputs, packages, etc

# Working with R: Comments

- Hashmarks are used to add comments and annotate your code

```
# Comments need to start with a hashmark, but don't need to end with one


# Hashmarks show up in green and are included to explain your code
```

- It's good practice to annotate your code
  - You can go back later and understand what you were doing

# Executing commands in R

- Code is executed when you press "Run" in the top right hand corner of the script window

- R runs the line of code where your cursor is located

- You can also highlight multiple lines to run at once

*Note: R is case sensitive*

# Outline

- What is data science?
- A data scientist's approach
- Introduction to R
  - Calculations in R
  - Reading data into R
  - Manipulating data in R
- Visualization in R
  - Basic plotting

# Working with R: variables

- A series of numbers (think columns in Excel) can be defined using the arrow (**<-**) or equals (**=**) sign

```
# Define variables with arrow
A <- c(5.5, -6.5, 7.5, 8.5)
B <- c(1, 2, 3, 4)
```
Script

or

```
# Define variables with equals sign
A = c(5.5, -6.5, 7.5, 8.5)
B = c(1, 2, 3, 4)
```
Script

- The command `c( )` stands for "concatenate" (join) a series of numbers

# Basic operations in R

## Adding

- Just use + sign

```
# Add variables
A = c(5.5, -6.5, 7.5, 8.5)
B = c(1, 2, 3, 4)
D = A + B
D
```
Script

```
> D
 [1] 6.5  -4.5  10.5  12.5
```
Console

## Multiplying

- Just use * sign

```
# Multiply variables
E = D*33
E
```
Script

```
> E
 [1] 214.5  -148.5  346.5  412.5
```
Console

*Enter formulas in top left window (script)*
*Output is shown in bottom left window (console)*

# Working with R: variables

- When a variable is named (instantiated), R stores it in its "environment" and can use it for subsequent operations

# R can run several lines of code

- You can highlight several lines of code and press "Run" to execute all of them

- Highlighting can be done either with the mouse or by holding "Shift" and using the arrow keys

- You can execute a command by pressing "Ctrl" + "Enter" for PCs or "Command" + "Enter" for Macs

*Troubleshooting: if you have trouble with this, try restarting R, restarting your computer, or reinstalling R*

# Executing operations

- You can run several operations in 1 line of code

- Or you can separate steps and instantiate new variables to check your code more easily

# Outline

- What is data science?
- A data scientist's approach
- Introduction to R
  - Calculations in R
  - **Reading data into R**
  - Manipulating data in R
- Visualization in R
  - Basic plotting

# A note about data

- Data can be found on a variety of sites on the internet

- Processing data stored in different formats is covered in a separate course

- For the purposes of this course, we will provide all the data sets already cleaned

# Loading data in R

- Loading data from your computer
  - Point and click

# Loading data in R

- Loading data from your computer
  - Enter code into script window
  - `crime_incidents_2013` is instantiated as the label of the data set

# Loading data in R

- Loading data from the internet
  - Enter code into script window
  - `crime_incidents_2013` is instantiated as the label of the data set

# Visualizing data

- Once data is loaded, you can see it as a spreadsheet by either:
  - Pressing the "spreadsheet" button in the top right window
  - Using the `View()` function in the script window

# `dir()` function

- Lists all the files in a particular directory



Number of the file in the list

# R can read many types of files

```
read.csv("filename.csv")          # read Excel files converted to csv format

read.table("filename")            # reads a table from a text file

read.spss("filename.spss")        # reads SPSS files

read.dta("filename.dta")          # reads Stata files

read.ssd("filename.ssd")          # reads SAS files

read.octave("filename.octave")    # read Octave files

read.mtp("filename.mtp")          # read Minitab files

read.systat("filename.systat")    # read Systat files

read.JPEG("filename.jpg")         # read JPEG image files
```

Note: this requires us to install package called 'jpeg', we will cover packages later

# Types of data used in R

## Basic (units)

- Integers ( -1, 5, 100 )
- Numerics ( 2.54 )
- Characters ( "Hello" )
- Logicals ( TRUE )
- Factors ( A B C )

## Composites

- Vectors
- Matrices (arrays)
- Lists
- Data frames

*Basics make up composites*

*Note: R works only with types of data it "understands." When loading new data, you need to tell R what kind of data you are loading*

# Why do data types matter?

- Different functions and commands in R can only work with certain types of data

- You may need to tell R that data from a spreadsheet is either a vector a list or another type of data before you can perform analysis on it

- Note that a spreadsheet that only includes numbers can be called either a list or a matrix

- You should tell R how to interpret the data based on what analysis you'd like to perform

# Vectors in R

**Vectors allow you to automatically sort text documents, images, etc.**

- A vector is a collection of elements of the same type (a column with either all numbers or all letters)

- R reads data as vectors

- Vectors allow you to manipulate a lot of data with a single command

- Operation between two vectors requires that vectors contain the same number of entries (same length)

Console

```
> # Create a vector
> x = c(1,2,3,4); x
[1] 1 2 3 4

> # Create a vector of years
> year = 2005:2010; year
[1] 2005 2006 2007 2008 2009 2010

> # Multiply constant by a vector
> x*3
[1]  3  6  9 12

> # Add two vectors
> y = c(1,2,2,1); x + y
[1] 2 4 5 5
```

**";"** separates commands that can be on separate lines

# Vectors in R

A vector



| | CCN | REPORTEDDATE | REPORTEDTIME | SHIFT | OFFENSE |
|---|---|---|---|---|---|
| 1 | 4104147 | 4/16/13 | 12:00:00 AM | MIDNIG... | HOMICIDE |
| 2 | 5047867 | 6/5/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 3 | 7083463 | 7/8/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 4 | 9172197 | 4/8/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 5 | 9251354 | 2/27/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 6 | 100289... | 2/27/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 7 | 100335... | 10/10/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 8 | 101249... | 4/9/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |
| 9 | 110101... | 7/31/13 | 12:00:00 AM | MIDNIG... | HOMICIDE |
| 10 | 110455... | 1/31/13 | 12:00:00 AM | MIDNIG... | HOMICIDE |
| 11 | 112502... | 7/8/13 | 12:00:00 AM | MIDNIG... | SEX ABUSE |

Showing 1 to 12 of 35,826 entries

# Matrices in R

## Matrices allow us to work with several columns of data at once

- A matrix is one or more vectors stacked next to each other

- Matrices can have row and column names, which can be determined and/or assigned (i.e. a customer list) by `dimnames, rownames` or `colnames` functions

- A matrix is a table where all the data is of the same type (i.e. numbers or letters)

```
> # Create a matrix
> mat = matrix(1:6, nrow = 3, ncol = 2)
> mat
     [,1] [,2]
[1,]    1    4
[2,]    2    5
[3,]    3    6
```

```
# Get matrix information
class(mat)         # "matrix"
is.vector(mat)     # FALSE
is.matrix(mat)     # TRUE
length(mat)        # 6
dim(mat)           # 3 2
```

# Customizing matrices in R

## Add names and adjust data order

- If you want to put in your data by row instead of by column, type `byrow = TRUE`

- You can add names for rows and columns by using the function `dimnames`

```
> mat1 = matrix(1:6, nrow = 3, ncol = 2,
byrow = TRUE)
> mat1
     [,1] [,2]
[1,]    1    2
[2,]    3    4
[3,]    5    6

> mat2 = matrix(1:6, nrow = 3, ncol = 2,
byrow = TRUE, dimnames = list(c("Row1",
"Row2", "Row3"), c("Col1", "Col2")))
> mat2
     Col1 Col2
Row1    1    2
Row2    3    4
Row3    5    6
```

# Lists in R

## Lists allow you to work with different types of data mixed together

- A list is a vector with different types of elements (data)

- The elements of a list can be numeric vectors, character vectors, matrices and other lists

- List components are determined by $ signs

```
> # Create a list
> Johnsons = list(husband = "Bill", wife =
"Joanna", children = TRUE, child.ages = c(3,
13, 18))
> Johnsons
$husband
[1] "Bill"

$wife
[1] "Joanna"

$children
[1] TRUE

$child.ages
[1]  3 13 18
```

# Lists in R

## You can pull out detailed information from lists

- `length("list")`: gives you the number of components

- `class("list")`: identifies the data type of the information

- `names("list")`: identifies the name of each component

- `"list"[1:2]`: identifies the specific data in those components

- `"list"$"component_name"`: identifies the specific data in that component

```
> # Get list information
> length(Johnsons)
[1] 4
> class(Johnsons)
[1] "list"
> names(Johnsons)
[1] "husband"    "wife"       "children"
[4] "child.ages"
> Johnsons[1:2]
$husband
[1] "Bill"

$wife
[1] "Joanna"


> Johnsons$child.ages
[1]  3 13 18
```

R automatically numbers the components

# Lists in R

A list includes several data types and unstructured data

# Data frames in R

## Data frames allow you to work with alpha-numeric and other data types

- Data frames are matrices with different data types

- Each column is a vector of the same length, types may differ

- The elements of a data frame can be numeric vectors, factor vectors, and logical vectors

- Consist of observations (rows), variables (columns)

Console

```
> # Setting up variables for data.frame
> first.name = c('Joe', 'Bob', 'Jill')
> last.name = c('Li', 'Dan', 'Smith')
> age = c(45, 20, 37)

> # Setting up data.frame
> df = data.frame(first.name,
> last.name, age)


> # Output
> df
  first.name last.name age
1        Joe        Li  45
2        Bob       Dan  20
3       Jill     Smith  37
```

# Combining lists into data frames

You can combine lists and create data frames

- Note that these lists must have the same number of components

- Use function `rbind` (row bind) to match the categories and combine lists
  - When using `rbind`, the names of the columns must be the same

*Note: the `rbind` function won't work with multiple data points within an element of a list*

```
# Combining lists into a data frame

Simpsons = list(husband = "Homer", wife =
                "Marge", children = TRUE,
                num.kids = 3)
Belchers = list(husband = "Bob", wife =
                "Linda", children = TRUE,
                num.kids = 3)
rbind(Simpsons, Johnsons)
```

Console

```
> rbind(Simpsons, Belchers)
         husband wife    children num.kids
Simpsons "Homer" "Marge" TRUE     3
Belchers "Bob"   "Linda" TRUE     3
```

# Data frames in R

A data frame includes several data types

# Creating data types in R

| Composite data types | Instantiating (creating) | Appropriate use |
|---|---|---|
| Vector | `c( . . . )` | Store simple data types in a row or column |
| Matrix | `matrix(data, nrows, ncols)` | Store multiple rows and columns of a single data type |
| Data Frame | `data.frame(data matrix)` | Work with alphanumeric and other data types |
| List | `list(element_1,...,element_k)` | Work with different data types simultaneously |

# Outline

- What is data science?
- A data scientist's approach
- Introduction to R
  - Calculations in R
  - Reading data into R
  - **Manipulating data in R**
- Visualization in R
  - Basic plotting

# DC crime data 2013

## What does the data look like?

- Almost 36,000 rows, each row represents a crime

- Information includes date, time, offense, method, location, precinct, etc.

- Always take the time to look over your data – this habit will help you understand the information you are working with

# Subsetting data

- You can separate out certain columns of the data by using a $ sign

- The format is
  `"data_set"$"name_of_column"`

- If the list you are creating has columns of different lengths, the `View()` function will not work
  - Use the `head()` function or just enter the new variable you created instead

Script

```
crime_list = list(a = crime_incidents_2013$REPORTEDDATE,
                  b = crime_incidents_2013$SHIFT)
```

Tells R that we are using a piece of the data set

```
View(crime_list)
head(crime_list)
crime_list
```



| | a | b |
|---|---|---|
| 1 | 4/16/13 | MIDNIGHT |
| 2 | 6/5/13 | MIDNIGHT |

Console ~/Dropbox/
[9993] EVENING  EVENING  EVENING  EVENING  EVENING  EVENING  EVENING  EVENING
[ reached getOption("max.print") -- omitted 25826 entries ]
Levels: DAY EVENING MIDNIGHT

# Factors in R

- A factor is a unique value in a data set

- For example, if five crimes are labeled as "arson, burglary, burglary, burglary, robbery", there would be three factors: arson, burglary, and robbery

- Use the `as.factor()` function to tell R to read data as a factor

```
# Set a category of our crime data set as a factor using as.factor()
crime_categories = as.factor(crime_incidents_2013$OFFENSE)

# Visualize the unique values in the OFFENSE category of the crime data set
levels(crime_categories)
View(levels(crime_categories))
```

# Factors in R

# Searching for terms in your data

- The function `grep()` can help you check if your data set includes a given term

```
# grep() only works on a data type called a character vector
# Recall that a vector is a single column
crime = as.character(crime_incidents_2013$METHOD)    Select the "METHOD" column

# Create a new variable that includes the output of grep()
gun_rows = grep("GUN", crime)

# View gun_rows, this includes all the row numbers with the word "GUN"
# Recall that the View() function displays data frames so tell R that "gun_rows"
# is a data frame
View(as.data.frame(gun_rows))
```

- You can search the entire data set by combining all the columns into one or searching each column individually

# Searching for terms in your data

- How many rows include "GUN" as the method?

```
> # Find the number of rows
> length(gun_rows)
[1] 2156
```

# Summarizing your data

- The `table()` function is a convenient way to see unique categories in your data and determine the frequency of occurrence

```
crime_summary = table(crime_incidents_2013$OFFENSE)

View(crime_summary)
```

Column to summarize

# Summarizing your data

- What if I want to see how crime is distributed over 24 hours?

```
crime_time = table(crime_incidents_2013$OFFENSE,
                   crime_incidents_2013$SHIFT)


View(crime_time)
```

# Sorting your data

- The `order()` function is a convenient way to sort your data
  - For more details of what this function can do, run the `?order` command

Create new data set　　　　Data set used　　　　Column to sort by

```
                                                                    Script
crime_order = crime_incidents_2013[order(crime_incidents_2013$METHOD,
                        decreasing = FALSE),]

View(crime_order)
```

- Note that the `decreasing` operator can be set to `TRUE` or `FALSE` depending on the order you are looking for (i.e. ascending or descending)

# Sorting your data

# Eliminating duplicate records

- You can eliminate repetitive instances of your data using `unique()`
  - Eliminates rows where every entry is the same as another row

```
crime_no_dups = unique(crime_incidents_2013)

# Check how many rows are identical by looking at the dimensions of each data set
dim(crime_incidents_2013)
dim(crime_no_dups)

# You can just compute the difference in the dimensions of the data set
dim(crime_incidents_2013) - dim(crime_no_dups)
```

**9 rows were eliminated**

```
Console ~/Desktop/Course Revisions/Week 1/
> dim(crime_incidents_2013)
[1] 35826    22
> dim(crime_no_dups)
[1] 35817    22
> dim(crime_incidents_2013) - dim(crime_no_dups)
[1] 9  0
```

# Subsetting: real data example

# Subsetting: selects only values you want

- When working with large data sets, you may want to use or visualize only a portion of your data at any given time

```r
# Vector subsetting
v = c(10, 20, 30, 40, 50, 60)    # define your vector with 6 numbers
v[3]                              # select the 3rd term
v[1:4]                            # select the first 4 terms
v[3:6]                            # select the 3rd through the 6th term
v[c(1, 3, 6)]                     # select the 1st, 3rd and 6th term
2*v[3:6]                          # multiply the 3rd through the 6th term by 2
```

```
[1] 30
[1] 10 20 30 40
[1] 30 40 50 60
[1] 10 30 60
[1] 60  80 100 120
```

# Subsetting: selects only values you want

- When working with large data sets, you may want to use or visualize only a portion of your data at any given time

```
# Matrix subsetting
m = matrix(c(5, 6, 7, 8), 2, 2) # define your matrix with 4 numbers in
                                          2 rows and 2 columns

m                          # view your matrix
m[2,]                      # select the numbers in the 2nd row
m[,1]                      # select the numbers in the 1st column
m[1,2]                     # select the number in the 1st row and 2nd column
m[1:3]                     # select the first 3 terms in the matrix
m[1:2,1:2]                 # select the #s in 1st and 2nd rows and 1st and 2nd cols
2*m[1:3]                   # multiply the first 3 terms in the matrix by 2
```

# Subsetting: selects only values you want

```
> m                      # view your matrix
       [,1]  [,2]
[1,]     5     7
[2,]     6     8
> m[2,]                  # select the numbers in the 2nd row
[1] 6 8
> m[,1]                  # select the numbers in the 1st column
[1] 5 6
> m[1,2]                 # select the number in the 1st row and 2nd column
[1] 7
> m[1:3]                 # select the first 3 terms in the matrix
[1] 5 6 7
> m[1:2,1:2]             # select the #s in 1st and 2nd rows and 1st and 2nd cols
       [,1]  [,2]
[1,]     5     7
[2,]     6     8
> 2*m[1:3]               # multiply the first 3 terms in the matrix by 2
[1] 10 12 14
```

Console

# Subsetting: real data example

# Subsetting: real data example

```
# Select data from the crime data set based on time of day the crime occurred
crime_time = subset(crime_incidents_2013, SHIFT == "DAY")

# Or you can also use
crime_time = crime_incidents_2013[crime_incidents_2013$SHIFT == "DAY",]
View(crime_time)

# Select 1 column from the data set
crime_date = subset(crime_incidents_2013, select = REPORTEDDATE)
View(crime_date)
```

# Subsetting: real data example

# Subsetting: real data example

- What if we wanted to select only robberies that happened during the day?
  - When performing several operations, keep the steps separate so it's easier to check your code

```
# First, only select the crimes that happened during the day
crime_time = subset(crime_incidents_2013, SHIFT == "DAY")

# Second, select only robberies out of the new data set called "crime_time"
crime_time_robbery = subset(crime_time, OFFENSE == "ROBBERY")
View(crime_time_robbery)
```

# Outline

- What is data science?
- A data scientist's approach
- Introduction to R
  - Calculations in R
  - Reading data into R
  - Manipulating data in R
- Visualization in R
  - Basic plotting

# Why build a visualization?

- To communicate your ideas
- To better understand your data
- To discover a new insight
- Visualization is a great way to do exploratory data analysis ("EDA")

# Exploratory data analysis

- R is a powerful tool for EDA because the graphics tie in with the functions used to analyze data

- You can create graphs without breaking your train of thought as you explore your data

- Visualization is an iterative process
  1. Analyze
  2. Manipulate
  3. Graph
  4. Repeat



1 **Analyze**

2 **Manipulate**

3 **Graph**

# What can you ultimately do

- The key to visualization in R is understanding how your data maps to the image you are creating
  - How R interprets your data
- Visualization types
  - Maps
  - Dynamic visualizations
  - Web based & interactive
  - 3D
  - Composite graphs made of many different elements
- Repeatable (no need to ever re-create a graph, just re-use the commands you already typed out)
- Use `??visualization` to see other examples

# Visualization in R

- R comes with basic plotting functionality

- More advanced visualizations are done through packages
  - `ggplot2`
  - `ggpairs`
  - `ggmaps`
  - `rgl`
  - `googleVis`
  - `lattice`

# Basic scatter plot

- `x, y` – coordinates
- `xlab, ylab, main` – graph labels
- `lwd` – size of dot / line
- `col` – colors, "red", "blue", "green",…

Script

```
# Suppose we want to plot y = x^2 + 5
x = seq(1:10)
y = x^2 + 5

plot(x, y,
     lwd = 10,
     col = "red",
     main = "y = x^2 + 5")
```

# Basic line plot

- `x, y` – coordinates
- `xlab, ylab, main` – graph labels
- `lwd` – size of dot / line
- `type` – line, points or both
  - `"l"` – line
  - `"b"` – line and dots
- `col` – colors, "red", "blue", "green",...

Script

```
# Suppose we want to plot y = x^2 + 5
x = seq(1:10)
y = x^2 + 5
plot(x, y,
    xlab = "x",
    ylab = "y",
    main = "y = x^2 + 5",
    lwd = 3,
    type = "l",
    col = "red")
```



y = x^2 + 5

# Basic bar plot

- `height` – value of each bar
- `xlab, ylab, main` – graph labels
- `names.arg` – labels to add to the x axis
- `col` – colors, "red", "blue", "green",…

Script

```
# Suppose we want to plot 10 random
# numbers using the formula y = x^2 + 5
x = seq(1:10)
y = x^2 + 5
barplot(height = y,
        width = 1,
        xlab = "x",
        ylab = "y",
        main = "Random numbers",
        names.arg = x,
        col = "red")
```



Random numbers

# Basic histogram

- `col` – colors, "red", "blue", "green",…
- `xlab, ylab, main` – graph labels
- `labels` – labels on top of the columns
- `breaks` – allows you to specify a custom number of groupings

Script

```
# Let's plot a histogram of a sequence
of numbers
Count_Data = c(1, 1, 1, 1, 1,
               2,
               3, 3, 3, 3,
               4, 4, 4,
               5, 5, 5, 5, 5, 5, 5)

hist(Count_Data,
     col = "blue",
     labels = TRUE,
     breaks = 2)
```



Histogram of Count_Data

# Functions for basic graphs

| Chart type | Command |
|:---:|:---:|
| Line chart | `plot()` |
| Bar chart | `barplot()` |
| Histogram | `hist()` |
| Boxplot | `boxplot()` |
| Density plot | `density()` |
| Pie chart | `pie()` |
| Scatter plot | `plot()` |
| Heat map | `image()` |

# Grammar of graphics: `ggplot2`

- R comes with basic plotting functionality

- More advanced visualizations done through packages: `ggplot2`

- `ggplot2` allows the user to create flexible visualizations by combining many elements into a single image



ALPHABET

GRAMMAR

CREATIVITY

# Why `ggplot2`

- Explore your data efficiently

- Communicate a visual story flexibly and efficiently

- Layer raw, summarized and contextual data
  - Demonstrate relationships

- Reproduce and extend your work easily



Age and Marriage Rates by Country

# Combine basic plots for enhanced effect


Bar plot


Histogram


Scatterplot


Bubble plot


Area plot


Time series


Heat map

# ggplot2: **building a graph**



**1. Specify data**

**2. Link data to visuals**

**3. Assign shapes**

**4. Adjust vis. effects**

**5. Adjust axes**

**6. Adjust legend**

**7. Customize theme**

**8. Layer statistics**

**9. Overlay text**

# Reading in data and loading library

```
# Install ggplot2
install.packages("ggplot2")

# Loading the library ggplot2
library(ggplot2)

# Learning what ggplot2 can do
library(help = ggplot2)

# Reading in files, windows uses the C:/ directory, Mac uses the ~/ directory
setwd("path to files")
US_GDP = read.csv("US GDP.csv")

# View the file
View(US_GDP)
```

- Help: http://docs.ggplot2.org/current/

# U.S. GDP data set

## What data is available?

- U.S. GDP by year in $USD billions

- U.S. GDP growth rate by year in percentages

- Years: 1980 - 2014

- Source: U.S. Federal Reserve

| | Year | US_GDP_BN | GDP_Growth_PC |
|---|---|---|---|
| 1 | 1980 | 2863 | 0.0 |
| 2 | 1981 | 3211 | 12.2 |
| 3 | 1982 | 3345 | 4.2 |
| 4 | 1983 | 3638 | 8.8 |
| 5 | 1984 | 4041 | 11.1 |
| 6 | 1985 | 4347 | 7.6 |
| 7 | 1986 | 4590 | 5.6 |
| 8 | 1987 | 4870 | 6.1 |
| 9 | 1988 | 5253 | 7.9 |
| 10 | 1989 | 5658 | 7.7 |
| 11 | 1990 | 5980 | 5.7 |
| 12 | 1991 | 6174 | 3.3 |

# Basic scatter plot

```
ggplot(US_GDP,
       aes(x = Year,
           y = US_GDP_BN)) +
  geom_point()
```

- `aes` – tells R how to map / assign the data
- **`geom_point()`** – tells R to use points to display the data

# Basic line plot

```
ggplot(US_GDP,
        aes(x = Year,
            y = US_GDP_BN)) +
  geom_line()
```

- aes – tells R how to map / assign the data
- **geom_line()** – tells R to use a line to display the data

# Scatter plot + line plot

```
ggplot(US_GDP,
       aes(x = Year,
           y = US_GDP_BN)) +
  geom_line() +
  geom_point()
```
Script

- **geom_line()** – tells R to use a line to display the data
- **geom_point()** – tells R to use points to display the data



*There are now 2 layers on the graph!*

# Basic histogram

```
ggplot(US_GDP, aes(x = GDP_Growth_PC)) +    Script
  geom_bar()
```

- `aes` – note that we're only including x-axis values
- **`geom_bar()`** – tells R to use bars to represent the data
- When given numeric data, the default setting of `geom_bar()` is to bin (put data into groups) the data and plot the number of instances each bin occurs -- a histogram



*A histogram shows the frequency of occurrence of each value on the x-axis.*
*This is a great way to see how often U.S. GDP grows at various rates.*

# Basic bar plot

```
ggplot(US_GDP,
       aes(x = Year,
           y = GDP_Growth_PC)) +
  geom_bar(stat = "identity")
```

- **geom_bar()** – tells R to use bars to represent the data
- The default setting of `geom_bar()` plots the number of instances of a particular value on the x-axis
- To use the data to determine bar height, we need to tell R to use the `"identity"` statistic (`stat = "identity"`)



*Percentage changes are usually best expressed with bars or area graphs, which are less likely to mislead your audience compared to a line graph*

# Bubble plot

```
ggplot(US_GDP, aes(x = Year,
                   y = US_GDP_BN,
                   size = GDP_Growth_PC)) +

  geom_point()
```

- The **size** argument tells R to interpret the `GDP_Growth_PC` column as a third variable related to the size of the data markers

- **geom_point()** – tells R to use points to represent the data



*A bubble plot allows you to plot a 3rd variable on a 2D graph*

# Colored bubble plot

Script

```
ggplot(US_GDP, aes(x = Year,
                   y = US_GDP_BN,
                   size = GDP_Growth_PC,
                   color = GDP_Growth_PC)) +

  geom_point()
```

- The **color** argument tells R to interpret the `GDP_Growth_PC` column as a third variable distinguished by color

- `geom_point()` – tells R to use points to represent the data



*The prior graph was harder to read because the points were so close in size.*
*Adjusting the color makes the chart more legible.*

# Additional R resources

1. http://www.statmethods.net

2. http://ggplot2.org

3. https://plot.ly/r/

4. www.r-bloggers.com

5. cran.r-project.org/

6. *The Art of R Programming* by Norman Matloff

7. *R for Everyone* by Jared P. Lander

8. *R Graphics Cookbook* by Winston Chang

# Additional data science resources

1. *Doing Data Science* by Cathy O'Neil & Rachel Schutt

2. *Data Science for Business* by Foster Provost & Tom Fawcett

3. *Data Smart* by John W. Foreman

4. *Mining the Social Web* by Matthew A. Russell

5. *Predictive Analytics* by Eric Siegel

6. *Analyzing the Analyzers* by Harlan Harris, Sean Murphy and Marck Vaisman

7. www.datasciencecentral.com

8. www.kdnuggets.com

# Congratulations!!

✔ You have learned how to code in R

✔ You have learned how to think about data

✔ You have learned how to format data

✔ You have learned how to build basic visualizations

# Don't stop here!

Continue your learning with us at datasociety.co! Learn how to:

Automate data cleaning processes
Build beautiful interactive visualizations
Find new patterns and groups in your data
Predict trends from your data
Analyze networks and find influencers
Prioritize and organize thousands of text documents
Classify new data points to predict behavior

And much more!