



**USING R FOR
EPIDEMIOLOGICAL
RESEARCH** Kiffer G. Card, PhD
Kirk J. Hepburn, MPP


INTRODUCTION TO R
Simon Fraser University, Faculty of Health
Sciences, Health Sciences (HSCI) 432






OUTLINE


- R Basics
- Exploring and Summarizing Data
- Combining Data Sets
- Recoding Variables
- Descriptive Statistics





1.1. R-BASICS

- Types of statistical software
- About *R*
- Installing *R*
- Installing *RStudio*
- A quick tour of *RStudio*
- Getting help
- R* packages
- R* syntax
- Annotating your *R* code



Types of Statistical Software

- Command-line software (e.g., SAS, R, SPSS scripts)
 - requires knowledge of syntax of commands
 - reproducible results through scripts
 - detailed analyses possible
- GUI-based software (e.g., SPSS point-and-click functions, Excel)
 - does not require knowledge of commands
 - not reproducible actions
- Hybrid types (both command-line and GUI)

R Basics


Exploring Data

Summarizing Data


Merging Data Sets

Rescaling Variables

Plotting Data



Types of Statistical Software



SPSS, STATA, and SAS users are like muggles. They are limited in their ability to change their environment. They have to rely on algorithms that have been developed for them. The way they approach a problem is constrained by how SAS/IBM-employed programmers thought to approach them. And they have to pay money to use these constraining algorithms.

R Basics

Exploring Data


Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Types of Statistical Software



R users are like wizards. They can rely on functions (spells) that have been developed for them by statistical researchers, but they can also create their own. They don't have to pay for the use of them, and once experienced enough (like Dumbledore), they are almost unlimited in their ability to change their environment.

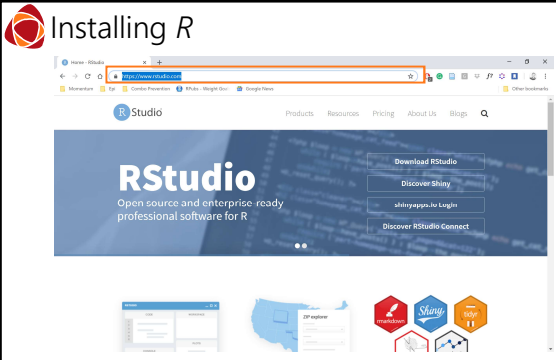
- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescuing Variables
- Plotting Data

About R

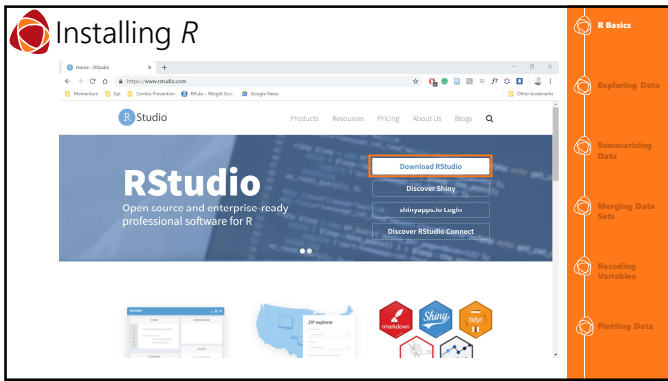
- R is a **statistical language** that was developed in response to Bell Labs (1976) S/S-plus language.
- R is commonly accessed through the **integrated development environment** called **RStudio**.
- Base R** performs a vast number of useful statistical operations.
- Base R can be enhanced with **packages**.
- These packages are **open source**, created by the community of R users, and typically documented in the *Journal of Statistical Software*. These resources are therefore **open for public use free of charge**.

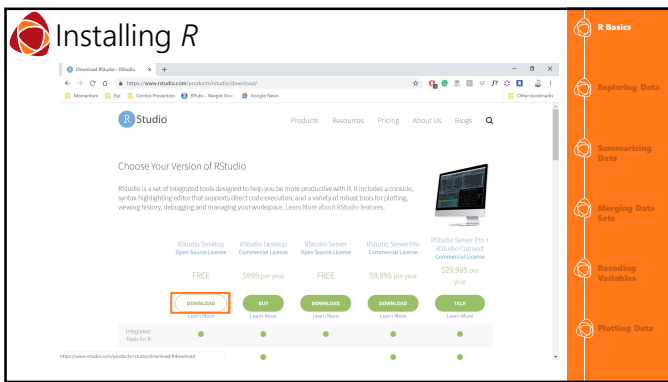
- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescuing Variables
- Plotting Data

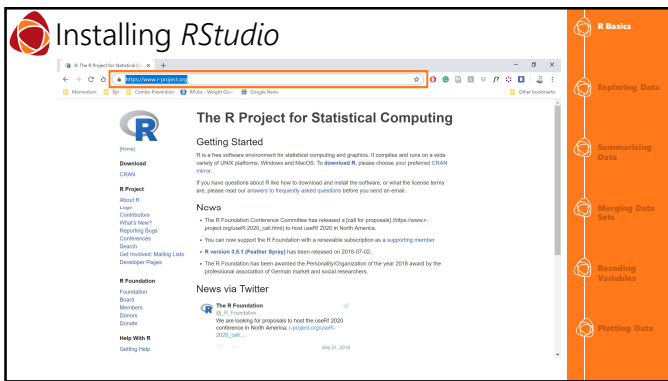
Installing R

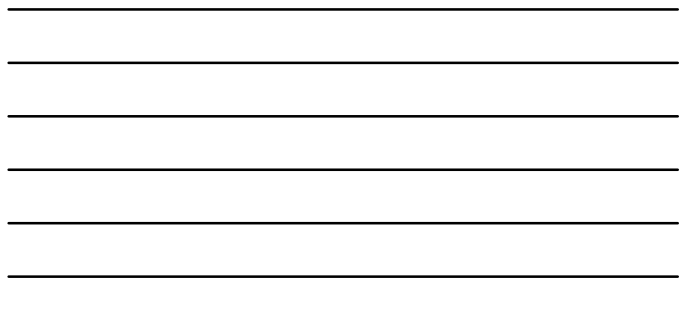
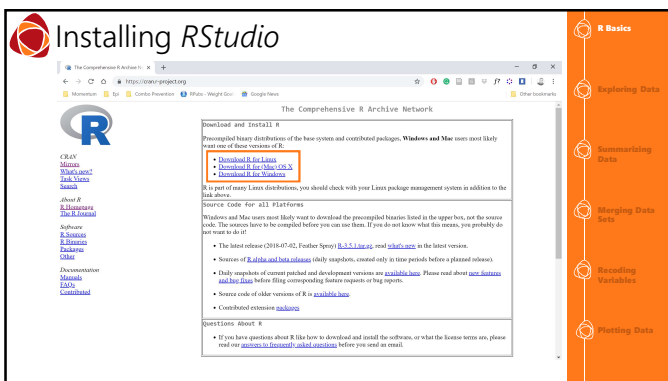
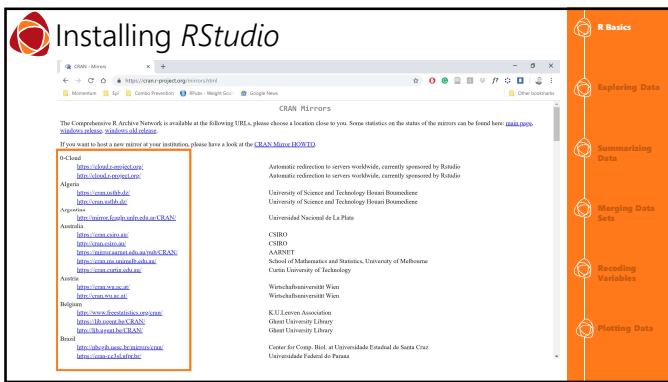
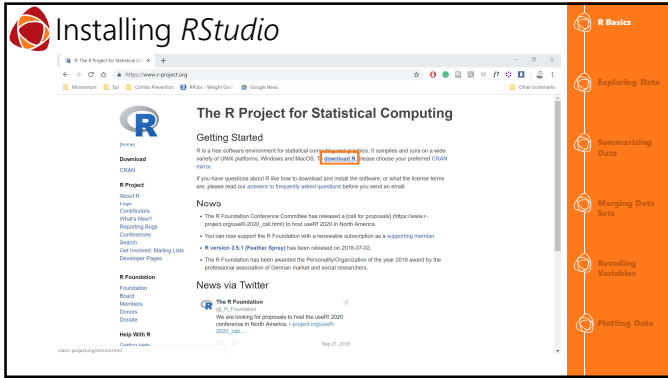


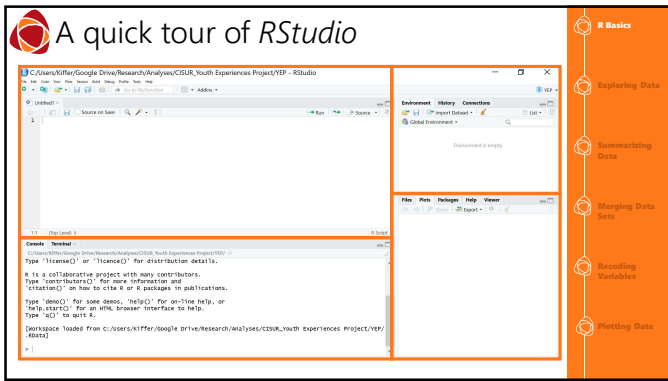
- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescuing Variables
- Plotting Data

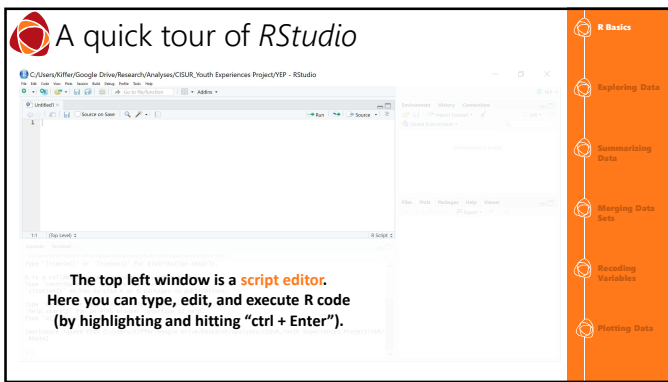


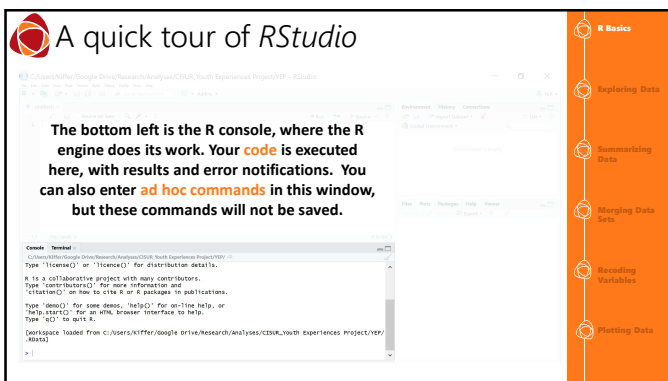












A quick tour of RStudio

Retains Code You Type Here

Does NOT Retain Code You Type Here

R Basics

- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescuing Variables
- Plotting Data

A quick tour of RStudio

In the top right, R objects are displayed and summarized. Objects are named sets of data, whether a set of data from a survey, a list of names, a set of randomly generated numbers, or even functions you have defined (created) yourself.

Note: R is case sensitive. DATA, data, and Data would be three different objects.

R Basics

- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescuing Variables
- Plotting Data

A quick tour of RStudio

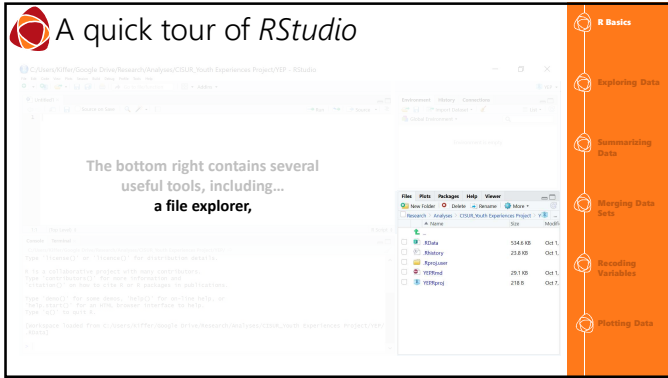
The bottom right contains several useful tools, including...

R Basics

- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescuing Variables
- Plotting Data

A quick tour of RStudio

The bottom right contains several useful tools, including...
a file explorer,

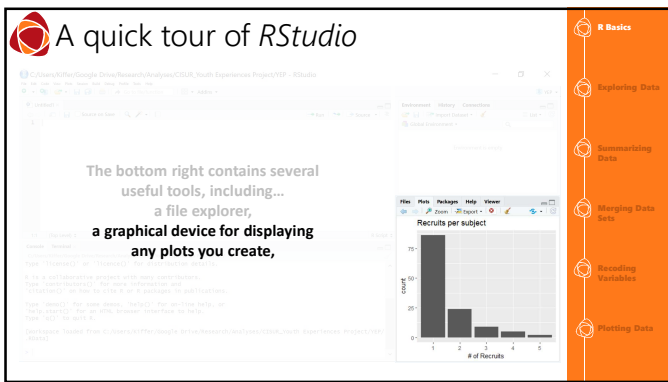


R Basics

- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

A quick tour of RStudio

The bottom right contains several useful tools, including...
a file explorer,
a graphical device for displaying any plots you create,

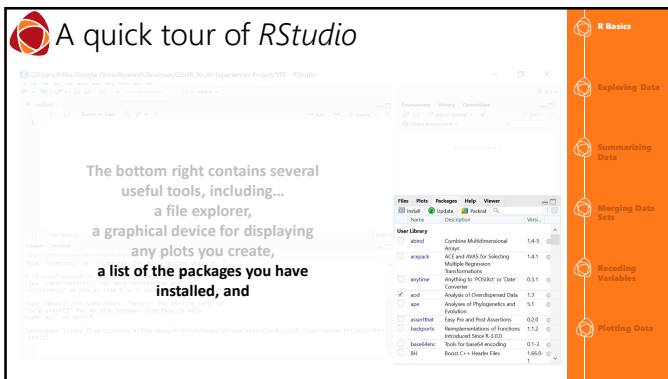


R Basics

- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

A quick tour of RStudio

The bottom right contains several useful tools, including...
a file explorer,
a graphical device for displaying any plots you create,
a list of the packages you have installed, and



R Basics

- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

A quick tour of RStudio

The bottom right contains several useful tools, including...

- a file explorer,
- a graphical device for displaying any plots you create,
- a list of the packages you have installed, and
- a browser for help documentation.

The screenshot shows the RStudio interface with the help viewer open for the `as_rds` function. The help text includes the title "Coerces a data frame object into an rds data frame object.", a description stating "This function converts a regular R data frame into an rds::rds object.", and usage instructions: "## Defaults: 03 methods available: rds::as_rds, rds::as_rds, ...".

Getting Help

Specific functions are described fully in R and can be found in the "help" viewer.

The screenshot shows the RStudio interface with the help viewer open for the `as_rds` function. The help text includes the title "Coerces a data frame object into an rds data frame object.", a description stating "This function converts a regular R data frame into an rds::rds object.", and usage instructions: "## Defaults: 03 methods available: rds::as_rds, rds::as_rds, ...".

Getting Help

The screenshot shows the RStudio interface with the help viewer open for the `as_rds` function. Two red callout boxes are present: box 1 points to the prompt `?<function name>` in the console, and box 2 points to the help text in the help viewer. The help text includes the title "Arithmetic Mean", a description "Generic function for the (binned) arithmetic mean.", and usage instructions: "## Defaults: 03 methods available: rds::as_rds, rds::as_rds, ...".

Getting Help

<function name>

mean Base

Arithmetic Mean

Description
Generic function for the (trimmed) arithmetic mean.

Usage
mean(x, ...)

Arguments
x: numeric (or methods)

Getting Help

- Furthermore, if you have questions about a specific function or error, the R community is extremely active.
 - Your question has likely already been asked on *StackExchange* or similar websites:
 - <https://stackoverflow.com/>
 - <https://www.r-bloggers.com/>
 - <https://www.statmethods.net/>

R Packages

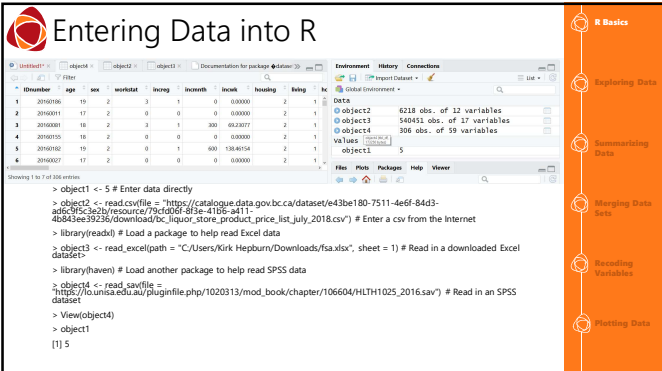
dplyr **ggplot2** **forcats** **stringr**



1.2. EXPLORING DATA

- Entering your data
- Using R to look at your data
- Types of data in R
- Data dictionaries
- Skip Patterns

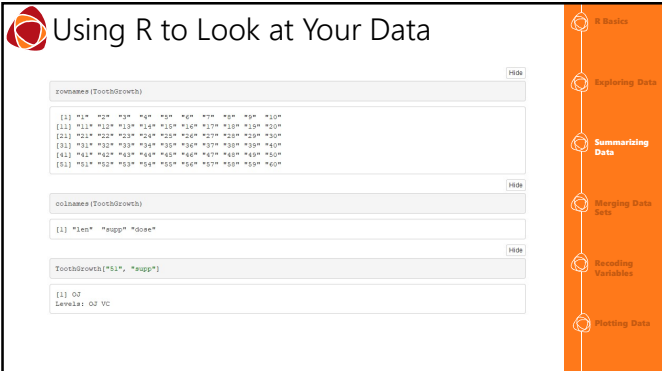
Entering Data into R



```

> object1 <- 5 # Enter data directly
> object2 <- read.csvfile = "https://catalogue.data.gov.bc.ca/dataset/643be180-7511-4e6f-84d3-
ad69fc3e2b/resource/79cd06f-83e-4195-a411-
4b44ee39236/download/cv_liquor_store_product_price_list_july_2018.csv" # Enter a csv from the Internet
> library(readxl) # Load a package to help read Excel data
> object3 <- read_excel(path = "C:/Users/Kirk Hepburn/Downloads/fsa.xlsx", sheet = 1) # Read in a downloaded Excel
dataset
> library(haven) # Load another package to help read SPSS data
> object4 <- read_savfile =
"https://o.unisa.edu.au/pluginfile.php/1020313/mod_book/chapter/106604/HLTH1025_2016.sav" # Read in an SPSS
dataset
> View(object4)
> object1
[1] 5
    
```

Using R to Look at Your Data



```

> viewname(ToothGrowth)
[1] "1" "2" "3" "4" "5" "6" "7" "8" "9" "10"
[1] "11" "12" "13" "14" "15" "16" "17" "18" "19" "20"
[1] "21" "22" "23" "24" "25" "26" "27" "28" "29" "30"
[1] "31" "32" "33" "34" "35" "36" "37" "38" "39" "40"
[1] "41" "42" "43" "44" "45" "46" "47" "48" "49" "50"
[1] "51" "52" "53" "54" "55" "56" "57" "58" "59" "60"

> colnames(ToothGrowth)
[1] "len" "supp" "dose"

> ToothGrowth[1:1, "supp"]
[1] 0
Levels: 02 VC
    
```

R Data Types

- R has a wide variety of classes of objects.
- The classes that you should be concerned with:
 - Vectors
 - Numeric - 1, 2, 5.3, 6, -2, 4, NA
 - Factor - 1, 2, 3, 2, 3, 3, 2, 1, NA
 - Character/String - "one", "two", "three", NA
 - Logical - TRUE, TRUE, FALSE, NA
 - Data frames – objects that can contain multiple objects (as columns) with different classes
 - You can check an object's class using the "class()" function.

Data Dictionaries

- A data dictionary, or metadata repository, as defined in the IBM Dictionary of Computing, is a "centralized repository of information about data such as meaning, relationships to other data, origin, usage, and format."
 - Provides question & response text
 - Provides variable name & response codes
 - Provides descriptive frequency counts & proportions
 - Provides information about skip/display logic
 - Provides codes for missing-ness
 - Provides information about the data type for each variable

1.3. SUMMARIZING DATA

- Frequencies
- Proportions
- Cross-Tabulated Frequencies
- Cross-Tabulated Proportions
- Measures of Central Location
- Measures of Spread
- Measures of Kurtosis
- Measures of Skew
- Calculating Rates

Frequencies

- Count of the occurrences of a value
- "Cross-tabulate": to count the occurrences of a value across different sets of data

Frequency

```
table(Education = inferEducation) # Frequency of education levels in infertility data
```

Education	0-5yrs	6-11yrs	12+ yrs
	12	120	116

Cross-tabulation

```
table(Education = inferEducation, NumberOfBirths = inferParity) # Education level by Number of Births
```

Education	1	2	3	4	5	6
0-5yrs	3	0	0	2	0	6
6-11yrs	42	42	21	12	3	0
12+ yrs	54	28	15	3	3	2

Proportions

- Proportion of total occurrences matching each value

Proportion Table

```
prop.table(table(Education = inferEducation)) # Frequency of education levels in infertility data
```

Education	0-5yrs	6-11yrs	12+ yrs
	0.0988871	0.4633710	0.4377419

Cross-tabulation of Proportions

```
prop.table(table(Education = inferEducation, NumberOfBirths = inferParity, margin = 1)) # Education level b y Number of Births
```

Education	1	2	3	4	5	6
0-5yrs	0.25000000	0.00000000	0.00000000	0.25000000	0.00000000	0.50000000
6-11yrs	0.35000000	0.35000000	0.17500000	0.10000000	0.02500000	0.00000000
12+ yrs	0.46881724	0.39620490	0.12910394	0.02586207	0.02586207	0.01724138

Measures of Central Location

The **mean**, μ , is the average value of the data set.
 The **median** is the middle value of the data set.
 The **mode** is the most common value of the data set.


Measures of Central Location | Mean

- "Average", with equal weight on either side
- Best for normal data
- Highly susceptible to outliers
- Good for continuous and discrete data
- Not good for ordinal or nominal data

Mean

```
mean(infert$age) # Calculate mean age of sample
[1] 31.80409
```

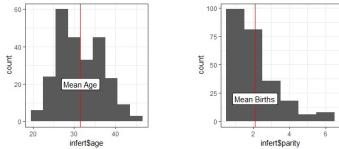

```
mean(infert$parity) # Calculate mean number of births in sample
[1] 2.092742
```



- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Measures of Central Location | Mean

Which mean is a better representation of its data?

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Measures of Central Location | Median

- The "middle point", or 50th percentile of the data
- Not sensitive to outliers, good for non-normal data
- Good for continuous, discrete, and ordinal data
- Not good for nominal data

Median

```
median(infert$age) # Calculate median age of sample
[1] 31
```

```
median(infert$parity) # Calculate median number of births in sample
[1] 2
```

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Measures of Central Location | Mode

- The most frequent value in the data
- Useful for all kinds of data
- No "easy" R function, so just look at a frequency table

Mode

```
table(infert$age) # Frequency of ages in sample to find mode
```

21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	44
6	6	3	15	15	15	15	12	12	21	15	18	15	12	9	9	6	3	6	3			

The histogram shows the distribution of ages. The x-axis is labeled 'intert\$age' and ranges from 20 to 40. The y-axis is labeled 'count' and ranges from 0 to 60. The highest bar is at age 28, with a count of 60. A vertical line points to this bar with the label 'Mode of Age'.

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Measures of Spread

- Range
- Interquartile Range
- Variance
- Standard Deviation
- Coefficient of Variation

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Measures of Spread | Range

- Minimum and maximum values
- Not terribly useful when there are outliers, similar to the mean

Range

```
range(intert$age) # Find the range of ages in the data
```

```
[1] 21 44
```

```
range(interc$parity) # Find the range of births in the data
```

```
[1] 1 4
```

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Measures of Spread | Interquartile Range

- How far apart are the 25th and 75th percentiles?
- Not as easily influenced by outliers

Interquartile Range

```
IQR(x = infer$age) # Find the interquartile range of ages in the data
```

```
[1] 7.25
```

```
IQR(x = infer$parity) # Find the interquartile range of births in the data
```

```
[1] 2
```

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Measures of Spread | Summary

- Multiple measures of spread provide a good description of the data

Summary

```
summary(infer$age) # Summarize the age data
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
21.00	28.00	31.00	32.50	35.25	44.00

```
summary(infer$parity) # Summarize the number of births
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	1.000	2.000	2.093	3.000	6.000

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Measures of Spread | Variance

- The amount of spread around the mean
- Useful for normal data

$$s^2 = \frac{\sum_{i=1}^n (X_i - X_{avg})^2}{n-1}$$

Variance

```
var(x = infer$age) # Calculate the variance of age in the data
```

```
[1] 27.57893
```

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Measures of Spread | Standard Deviation

- The square root of the variance
- Useful for normal data
- Interpreted in units of the data, so more descriptive than the variance

Standard Deviation

```
sd(x = infert$age) # Calculate the standard deviation of age in the data
```

```
[1] 5.261565
```

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Measures of Spread | Coefficient of Variation

- A way to compare the variation of two or more sets of data representing different quantities with different units
- Not frequently used

$$CV (\%) = \left(\frac{\text{Standard deviation}}{\text{Mean}} \right) \times 100$$

Coefficient of Variation

```
sd(x = infert$age) / mean(x = infert$age) * 100 # Calculate the coefficient of variation for age
```

```
[1] 14.4695
```

R Basics


Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data



1.4. MERGING DATASETS

Adding Variables

Adding Observations

Wide vs. Long Data

Dropping Variables

Dropping Observations

Adding Observations

- Continuation of a survey, or multiple administrations of a survey

trial1 # First dataset		
Year	Drug	Sex
2010	A	M
2010	B	F
2010	A	M

3 rows

rbind(trial1, trial2) # Combined datasets		
Year	Drug	Sex
2010	A	M
2010	B	F
2010	A	M
2011	B	F
2011	B	F
2011	A	M

6 rows

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Recoding Variables

Plotting Data

Dropping Variables or Observations

- When you might need to do it
- How to do it

trial1[, -3] # Remove 'Sex', the third variable		
Year	Drug	
2010	A	
2010	B	
2010	A	

3 rows

trial1[-3,] # Remove the third observation		
Year	Drug	Sex
2010	A	M
2010	B	F

2 rows

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Recoding Variables

Plotting Data

1.5. RECODING VARIABLES

- Recoding Categorical Variables
- Categorizing Continuous Data
- Combining Variables
- Missing Data

Data Visualizations

- Most plots can be generated using either
 - Base R**
 - or
 - ggplot**
- Because the R community is quite active, you can usually figure out how to do plots by googling "How to make a _____ in R."
 - For example, I found the following two guides – both of which are quite useful:
 - Base R - <https://www.statmethods.net/graphs/index.html>
 - Ggplot - <http://www.sthda.com/english/wiki/ggplot2-barplots-quick-start-guide-r-software-and-data-visualization>

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Bar Charts

A bar chart or bar graph is a chart or graph that presents **categorical data** with rectangular bars with heights or lengths proportional to the values that they represent. The bars can be plotted vertically or horizontally.

```
counts<-table(infert$education)
barplot(counts,
        main = "Participants by Years of Education",
        xlab = "number of participants", ylab = "Years of Education")
```

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

Plotting Data

Histograms

A histogram is a plot that lets you discover, and show, the underlying frequency distribution (shape) of a set of **continuous data**.

```
hist(infert$age,
     main = "Histogram of Age",
     xlab = "Age (in Years)", ylab = "number of participants")
```

R Basics

Exploring Data

Summarizing Data

Merging Data Sets

Rescaling Variables

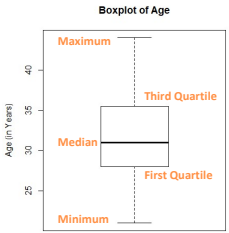
Plotting Data

Box Plots

- The box plot (a.k.a. box and whisker diagram) is a standardized way of displaying the **distribution of continuous data** based on the five number summary: minimum, first quartile, median, third quartile, and maximum.

```

boxplot(infert$age,
  main = "boxplot of Age",
  ylab = "Age (in Years)")
    
```



Boxplot of Age

Age (in Years)

Maximum

Third Quartile

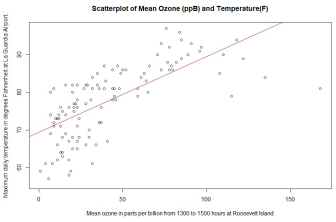
Median

First Quartile

Minimum

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Scatter Plots



Scatterplot of Mean Ozone (ppb) and Temperature (F)

- A scatter plot is a type of plot or mathematical diagram using X & Y coordinates to display values from **two continuous variables** within a dataset – often as a means to highlight a **correlation**.

```

plot(airquality$ozone, airquality$temp,
  main = "Scatterplot of Mean ozone (ppb) and Temperature (F)",
  xlab = "Mean ozone in parts per billion from 1300 to 1500 hours at Roosevelt Island",
  ylab = "Maximum daily temperature in degrees Fahrenheit at La Guardia Airport.",
  abline(lm(airquality$temp~airquality$ozone), col="red") # Linear regression line
    
```

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Line Graphs

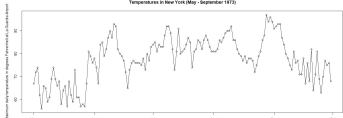
- A line graph, also known as a line chart, is a type of chart used to visualize the value of a **continuous variable over time**.

```

# Create a Year variable because the Air Quality dataset does not have this
airquality$year <- 1973 ## All observations were in 1973

# combine the month, day and year variables using the "paste" and "as.date" functions.
airquality$date <- as.date(paste(month, abbr, airquality$month),
  airquality$day =
  airquality$year ~ sep="-",
  format = "%b-%d-%Y")

# create a plot
plot(airquality$date, airquality$temp, type = "o",
  ylab = "Maximum daily temperature in degrees Fahrenheit at La Guardia Airport",
  xlab = "Date",
  main = "Temperatures in New York (May - September 1973)")
    
```



Temperatures in New York (May - September 1973)

- R Basics
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data

Essential Elements of a Good Graph

- Every graph should have a
 - A descriptive title
 - A key or legend
 - Units for each axis
 - Titles for each axis
 - A citation for the data source
- Often times it is helpful to highlight key values in a graph.
- Error bars are often helpful to highlight uncertainty in a graph.
- Fit lines are often helpful when displaying data with individual observations.
- Generally graphs should have high visual contrast.
- While use of color can enhance a graph – it may also reduce accessibility.
- Footnotes or plot descriptions should be used to clarify any misleading visual factors or to describe limitations in the data.

Percent of adults aged 18-75 years who had a high school diploma in the past 5 years or had a certificate in the past 10 years, by highest level of education attained, 2000-2015

Year	Less than High School	High School	Greater than High School
2000	20	30	35
2003	22	32	38
2006	24	34	40
2009	26	36	42
2012	28	38	44
2015	30	40	46

Source: Centers for Disease Control and Prevention, National Center for Health Statistics, National Health Interview Survey. Estimates for 2006 are projected for 2006-2009 because respondents in the 2000-2009 NHIS were asked about their most recent postsecondary, certificate, or degree attainment, while 2010 respondents were asked about their most recent high school or GED attainment. Percentages adjust for non-response. Data are age-adjusted to the 2000 US standard population using age groups: 18-44, 45-75.

- Basic
- Exploring Data
- Summarizing Data
- Merging Data Sets
- Rescaling Variables
- Plotting Data
