

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

STAT 408 - R OVERVIEW

January 11, 2018

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

R INTRO

WHY USE R?

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

R is:

- a free public domain implementation of S,
- the standard among (academic) professional statisticians,
- available for Windows, Mac, and Linux,
- an object-oriented and functional programming structure, and
- designed to connect to high-performance programming languages like C and Fortran.

WHY USE R?

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

R has:

- an open-software environment with a large community that makes getting help easy,
- a massive set of packages for statistical modeling, data science, visualization, and importing and manipulating data,
- powerful tools for replicating and communicating your results,
- an interactive development environment (R Studio) tailored for interactive data analysis and statistical programming,
- available for Windows, Mac, and Linux, and
- an object-oriented and functional programming structure.

READING DATA FILES

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

The ability to load datasets into R is an essential skill. For this class, most of the files will be on the course webpage and can be directly downloaded using `read.csv`. Consider a dataset available at: <http://math.montana.edu/ahoegh/teaching/stat408/datasets/SeattleHousing.csv>

```
Seattle <- read.csv(  
  'http://math.montana.edu/ahoegh/teaching/stat408/datasets/SeattleHousing.csv',  
  stringsAsFactors = F)
```

VIEWING DATA FILES

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

A common function that we will use is `head`, which shows the first few rows of a data frame.

```
head(Seattle)
```

```
##      price bedrooms bathrooms sqft_living sqft_lot floors waterfront
## 1 1350000         3         2.50         2753   65005    1.0          1
## 2  228000         3         1.00         1190    9199    1.0          0
## 3  289000         3         1.75         1260   8400    1.0          0
## 4  720000         4         2.50         3450  39683    2.0          0
## 5  247500         3         1.75         1960  15681    1.0          0
## 6  850830         3         2.50         2070  13241    1.5          0
##      sqft_above sqft_basement zipcode      lat      long yr_sold mn_sold
## 1          2165              588  98070 47.4041 -122.451   2015     3
## 2          1190              0  98148 47.4258 -122.322   2014     9
## 3          1260              0  98148 47.4366 -122.335   2014     8
## 4          3450              0  98010 47.3420 -122.025   2015     3
## 5          1960              0  98032 47.3576 -122.277   2015     3
## 6          1270             800  98102 47.6415 -122.315   2014     6
```

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

R DATA STRUCTURES

DATA STRUCTURE OVERVIEW

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

R has four common types of data structures:

- Vectors
- Matrices (and Arrays)
- Lists
- Data Frames

DATA STRUCTURE OVERVIEW

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

The base data structures in R can be organized by dimensionality and whether they are homogenous.

Dimension	Homogenous	Heterogenous
1d	Vector	List
2d	Matrix	Data Frame
no d	Array	

VECTOR TYPES

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

There are four common types of vectors: logical, integer, double (or numeric), and character. The `c()` function is used for combining elements into a vector

```
dbl <- c(1,2.5,pi)
int <- c(1L,4L,10L)
log <- c(TRUE,FALSE,F,T)
char <- c('this is','a character string')
```

VECTOR TYPES

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

The type of vector can be identified using the `typeof()` function. Note that only a single data type is allowed.

```
typeof(dbl)
```

```
## [1] "double"
```

```
comb <- c(char,dbl)  
typeof(comb)
```

```
## [1] "character"
```

```
comb
```

```
## [1] "this is"
```

```
## [4] "2.5"
```

```
"a character string" "1"
```

```
"3.14159265358979"
```

EXERCISE: VECTORS

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

Create a vector with your first, middle, and last names.

SOLUTION: VECTORS

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

- 1 Create a vector with your first, middle, and last names.

```
andy.names <- c("Andrew", "Blake", "Hoegh")  
andy.names
```

```
## [1] "Andrew" "Blake" "Hoegh"
```

DATA FRAME OVERVIEW

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

A data frame:

- is the most common way of storing data in R
- is like a matrix with rows-and-column structure; however, unlike a matrix each column may have a different mode
- in a technical sense, a data frame is a list of equal-length vectors.

```
df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))  
kable(df)
```

x	y
1	a
2	b
3	c

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

SUBSETTING

VECTOR SUBSETTING: I

Subsetting allows you to extract certain elements from a data frame or vector (or matrix, array, list).

```
num.vec <- seq(from = 1, to = 9, by = 1); num.vec
```

```
## [1] 1 2 3 4 5 6 7 8 9
```

```
num.vec[1:3]
```

```
## [1] 1 2 3
```

```
num.vec[c(1,5,8)]
```

```
## [1] 1 5 8
```

```
num.vec[-5]
```

```
## [1] 1 2 3 4 6 7 8 9
```


VECTOR SUBSETTING: II

Subsetting also works with logical values or expressions.

```
num.vec[num.vec > 5]
```

```
## [1] 6 7 8 9
```

```
num.vec[num.vec != 6]
```

```
## [1] 1 2 3 4 5 7 8 9
```

```
num.vec[rep(c(TRUE, FALSE, TRUE), each=3)]
```

```
## [1] 1 2 3 7 8 9
```

DATA FRAME SUBSETTING: I

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

The same ideas apply to data frames, but the indices now constitute rows and columns of the data frame.

```
df <- data.frame(x=1:3, y=3:1, z=c('a','b','c'))  
df[,1]
```

```
## [1] 1 2 3
```

```
df[-1,c(2:3)]
```

```
##   y z  
## 2 2 b  
## 3 1 c
```

DATA FRAME SUBSETTING: II

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

There are also a couple built in functions in R for subsetting data frames.

```
df$x
```

```
## [1] 1 2 3
```

```
new.df <- subset(df, x >1); new.df
```

```
##   x y z
```

```
## 2 2 2 b
```

```
## 3 3 1 c
```

EXERCISE: SUBSETTING

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

- 1 Create a new data frame that only includes houses worth more than \$1,000,000.
- 2 (bonus) From this new data frame what is the average living square footage of houses. Hint columns in a `data.frame` can be indexed by `Seattle$sqft_living`

EXERCISE: SUBSETTING - SOLUTIONS

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

- 1 Create a new data frame that only includes houses worth more than \$1,000,000.

```
expensive.houses <- subset(Seattle, price > 1000000)
```

- 2 (bonus) From this new data frame what is the average living square footage of houses. Hint columns in a data.frame can be indexed by `Seattle$sqft_living`

```
mean(expensive.houses$sqft_living)
```

```
## [1] 3890.065
```

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

GRAPHICS

BASIC PLOTTING IN R: SCATTERPLOT

STAT 408 -
R OVERVIEW

R INTRO

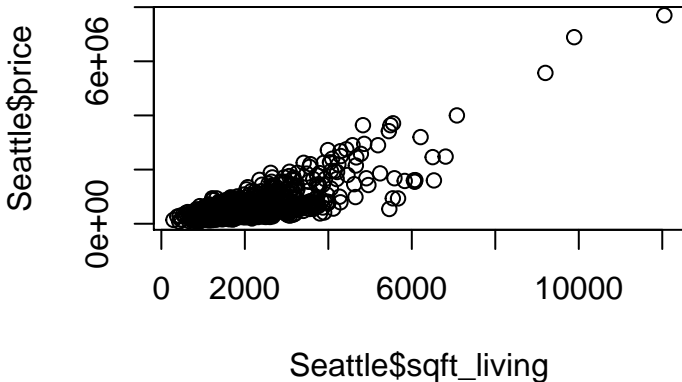
R DATA
STRUCTURES

SUBSETTING

GRAPHICS

Later in the course, we will spend considerable time on graphics. For now, let's consider some of the basic functionality in R.

```
plot(Seattle$price~Seattle$sqft_living)
```



BASIC PLOTTING IN R: LABELS

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
plot(Seattle$price~Seattle$sqft_living,  
     ylab='Price',xlab='Living Sqft')
```



BASIC PLOTTING IN R: PCH

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
plot(Seattle$price~Seattle$sqft_living,  
     ylab='Price',xlab='Living Sqft', pch=16)
```



BASIC PLOTTING IN R: COLOR

STAT 408 -
R OVERVIEW

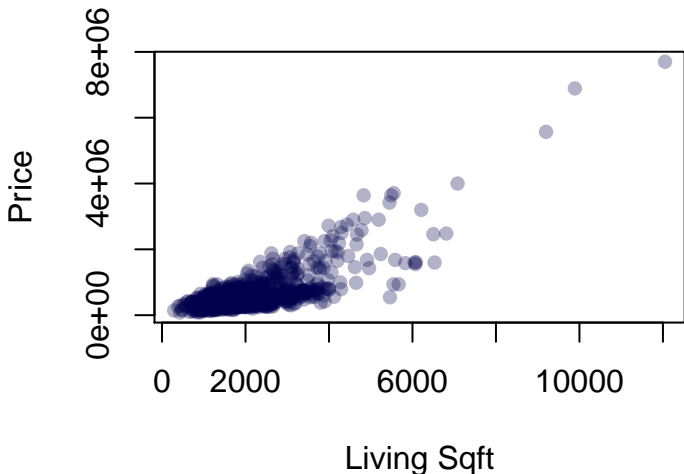
R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
plot(Seattle$price~Seattle$sqft_living, pch=16,  
     col=rgb(0,0,.3,.3),ylab='Price',xlab='Living Sqft')
```



BASIC PLOTTING IN R: TITLE

STAT 408 -
R OVERVIEW

R INTRO

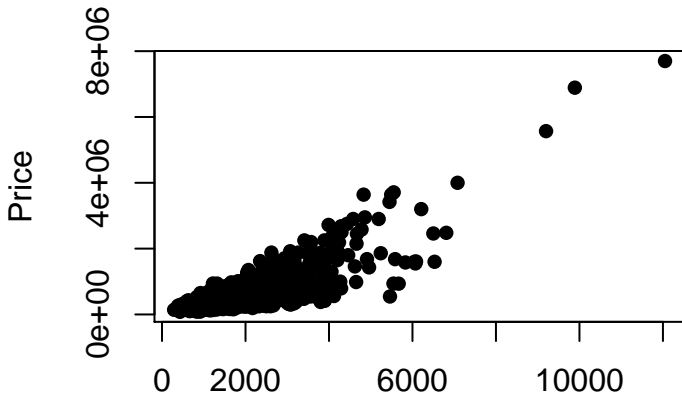
R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
plot(Seattle$price~Seattle$sqft_living, pch=16, ylab='Price',  
     xlab='Living Sqft',main='Price vs. Living Sqft')
```

Price vs. Living Sqft



BASIC PLOTTING IN R: HISTOGRAM

STAT 408 -
R OVERVIEW

R INTRO

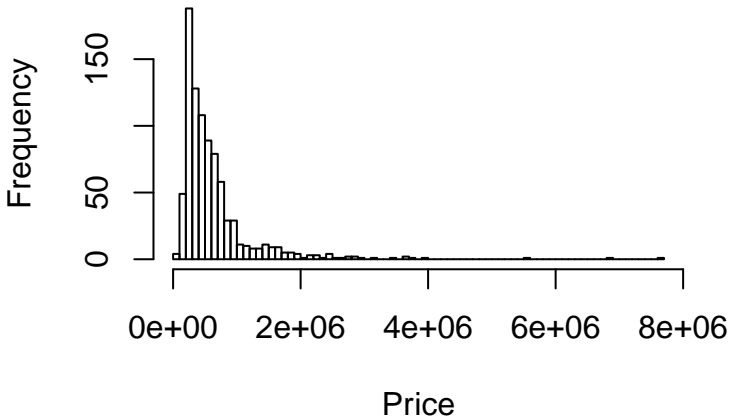
R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
hist(Seattle$price,xlab='Price', breaks='FD')
```

Histogram of Seattle\$price



BASIC PLOTTING IN R: HISTOGRAM

STAT 408 -
R OVERVIEW

R INTRO

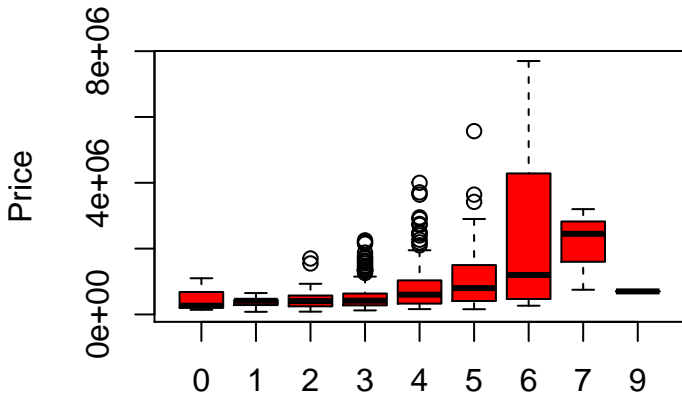
R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
boxplot(Seattle$price~Seattle$bedrooms,ylab='Price', col='red',  
        xlab='bedrooms',main='Price by Bedrooms for Seattle')
```

Price by Bedrooms for Seattle



EXERCISE: BASIC PLOT

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

- Using only the subset of homes worth more than a million dollars, create a graphic.

SOLUTION: BASIC PLOT

STAT 408 -
R OVERVIEW

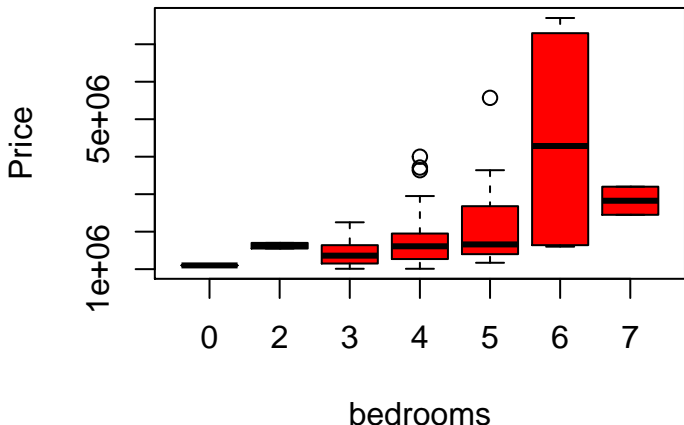
R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

Price by Bedrooms for Seattle



For homes worth more than \$1,000,000

SOLUTION: BASIC PLOT - WITH CODE

STAT 408 -
R OVERVIEW

R INTRO

R DATA
STRUCTURES

SUBSETTING

GRAPHICS

```
boxplot(expensive.houses$price ~
        expensive.houses$bedrooms,
        ylab='Price', col='red', xlab='bedrooms',
        main='Price by Bedrooms for Seattle',
        sub='For homes worth more than $1,000,000')
```