

# Exploratory Data Mining via Search Strategies Lab #1

Ross Jacobucci

The first lab will go over the basics of programming in R along with covering Smoothing and Stepwise Regression

To get started, I recommend downloading Rstudio as an interface to R. Benefits include:

- ▶ Easier to download packages
- ▶ Easier to view plots
- ▶ Can load datasets without code

# Loading Datasets

Many datasets that we will be using come with installed packages

```
library(MASS) # for boston data  
data(Boston)
```

Now the Boston dataset is in your workplace

```
head(Boston[,1:7],6)
```

```
##      crim zn indus chas   nox    rm  age  
## 1 0.00632 18  2.31    0 0.538 6.575 65.2  
## 2 0.02731  0  7.07    0 0.469 6.421 78.9  
## 3 0.02729  0  7.07    0 0.469 7.185 61.1  
## 4 0.03237  0  2.18    0 0.458 6.998 45.8  
## 5 0.06905  0  2.18    0 0.458 7.147 54.2  
## 6 0.02985  0  2.18    0 0.458 6.430 58.7
```

# Load your own datasets

Data saved as .dat or .txt

```
data = read.table(file.choose(), sep=" ", header=F, na.strings="NA")
```

.csv

```
data = read.csv(file.choose(), header=T, na.strings=".")
```

.sav for SPSS files

```
library(foreign)
data = read.spss(file.choose(), to.data.frame=TRUE)
```

.sas7bdat for SAS data files – note doesn't work very well, best to convert first

```
library(sas7bdat)
data = read.sas7bdat(file.choose())
# note many options
```

# Subsetting datasets

Let's say I have a dataset read in fine. How about subsetting columns and rows: \* specific columns – two ways

```
library(MASS)
data(Boston)
## I want specific columns

Boston.sub <- Boston[,c(1,2,8,9)]
## or
Boston.sub <- Boston[,c("crim","zn","dis","rad")]
```

Have to use c() to combine non-continuous elements in R !!!!! VERY IMPORTANT !!!!!

How about rows – Same process but better shortcuts

```
comps <- complete.cases(Boston)
Boston.comp <- Boston[comps,]

subs <- Boston$zn == 0
Boston.sub <- Boston[subs,]
dim(Boston.sub)
```

```
## [1] 372 14
```

Notice: left of comma = row indicator; right = column indicator

# Dataset Types

For manipulating datasets, almost always best to make sure in data.frame

```
library(MASS)
data(Boston)
```

```
Boston.df <- data.frame(Boston)
```

Can tell type by running str()

Some R packages require X and Y variables to be in separate matrices

```
library(MASS)
data(Boston)

Y <- as.numeric(Boston$medv)
X <- data.matrix(Boston[, -14])

out <- lm.fit(X,Y)
```

This is very specific to data mining packages

# Missing Data

Can be summed up in two words: Not Fun

R requires all missingness to be coded as NA, therefore it is best to deal with when reading in data:

```
data = read.csv(file.choose(),header=T,na.strings=".")
```

This takes all missing values coded with a period in the original dataset and converts them to NA for R

Once you have a dataset in R and missingness coded as NA:

```
# have to use is.na a lot  
is.na(c(1,2,NA,6))
```

```
## [1] FALSE FALSE  TRUE FALSE
```

Returns a logical indicator. Can save this to subset dataset based on number of missings

## Missing Data Continued

```
library(psych)
data(bfi)
dim(bfi)
```

```
## [1] 2800 28
```

Now, subsetting based on number of missing values per person

```
ids <- rowSums(is.na(bfi)) < 3
bfi.sub <- bfi[ids,]
dim(bfi.sub)
```

```
## [1] 2779 28
```



# Variable Types

This is really important as it will change the type of estimator used e.g. logistic vs. linear regression

```
bfi.sub <- bfi[,1:5]
library(rpart)
out1 <- rpart(A1 ~ ., bfi.sub)
out1
```

```
## n=2784 (16 observations deleted due to missingness)
##
## node), split, n, deviance, yval
##      * denotes terminal node
##
## 1) root 2784 5515.1380 2.413434
##    2) A2>=4.5 1907 3382.5430 2.126377
##      4) A2>=5.5 874 1618.1660 1.925629 *
##      5) A2< 5.5 1033 1699.3550 2.296225 *
##    3) A2< 4.5 877 1633.7580 3.037628
##      6) A2>=3.5 554 879.9495 2.819495 *
##      7) A2< 3.5 323 682.2353 3.411765 *
```

## Variable Types Continued

```
out2 <- rpart(as.factor(A1) ~ ., bfi.sub)
out2
```

```
## n=2784 (16 observations deleted due to missingness)
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
## 1) root 2784 1862 1 (0.33 0.29 0.14 0.12 0.08 0.029)
##    2) A2>=5.5 874  386 1 (0.56 0.22 0.072 0.064 0.055 0.03) *
##    3) A2< 5.5 1910 1285 2 (0.23 0.33 0.18 0.15 0.092 0.029) *
```

Changing the variable type, from integer to factor also changes the cost function. This results in very different equations and also very different results.

Most important is the variable type of the Y variable. Also can change it in the actual function by setting “family =”

Also can change variable type in dataset and the data mining function will recognize this. However, I almost always just include it in the actual script to make sure it changes to the correct estimator.

# How to get help

```
library(caret)  
?train
```

However, I usually just google the R package and look at the manual

Once you have run a script and have it save as an object, you can

```
library(MASS)  
data(Boston)  
out <- lm(medv ~ ., data=Boston)  
str(out)
```

Using `str()` will list all of the attributes of the “out” object. This is important to use if `print()` or `summary()` don't give you the information you want. Most applicable if you need nitty gritty details

# Most Important Details for Learning R

1. Read in data
2. How to subset dataset
3. Variable types
4. Understand where the results go and how to access them
5. The intricacies to each package