

Copyright © 2014 - 2016 Data Science Dojo

DATA SCIENCE DOJO

2205 152ND AVE NE, REDMOND, WA 98052

All rights reserved. No part of this publication may be reproduced, distributed, or transmitted in any form or by any means, including photocopying, recording, or other electronic or mechanical methods, without the prior written permission of the publisher, except in the case of brief quotations embodied in critical reviews and certain other noncommercial uses permitted by copyright law.

The Library of Congress has catalogued this edition as follows:

Data Science Dojo LLC

Data Science and Data Engineering Handbook : 2.3

Registration Number: TX0008030133

2015-05-22

Printed in the United States of America

Copyright © 2014-2016



1. GitHub Tutorial

2. Introduction to R

2.1 Introducing R

2.1.1 What is R?

R is a language and environment for statistical computing and graphics. It provides a wide variety of statistical (linear and nonlinear modeling, classical statistical tests, time-series analysis, classification, clustering) and graphical techniques, and is highly extensible. The S language is often the vehicle of choice for research in statistical methodology and R provides an Open Source route to participation in that activity.¹

R is open source and highly extensible. There are currently 6438 packages in the Comprehensive R Archive Network ("CRAN") package repository.

In this chapter we will go over R data types, basic operations, reading and writing data, statistical simulation, and basic plotting systems.

2.1.2 R is Vectorized

```
1 f = 1 * 11 + 2 * 22 + 3 * 33
2
3 > A = c(1, 2, 3)
4 > B = c(11, 22, 33)
5
6 > f = A * B
7 > f
8 [1] 11 44 99
```

The vectorization of R makes it good for statistics and data.

¹The R Project for Statistical Computing: <http://www.r-project.org/>

2.1.3 R for Data Science

There are various advantages and disadvantages to using R for data science. One of its advantages is that it was designed for statistical analysis, making it more straightforward. For example, linear regression can be performed with a single line of R. Contrastingly in Python, this requires the use of several third-party libraries to represent the data (NumPy), perform the analysis (SciPy), and visualize the results (matplotlib).

One of the disadvantages of R is that it does not scale well with large data. Big companies like Google use R to play with data and experiment with new machine learning methods. If it works well and they decide to move forward, they use something like C. However, this doesn't mean that R can't be used for big data!

2.2 R Data Types

2.2.1 Hello World

```
1 C <- "Hello World" # type Enter
2 print(c) # print some text
3 # anything after a hash (#) is a comment
```

2.2.2 Basic Data Types

There are five basic classes (atomic classes) of R. They are character, numeric, integer, complex, and logical.

Among these classes the assignment operators are as follows:

```
1 n <- 10 # assign value 10 to variable 'n'
2 0.3 -> s # the arrow can go both ways
3 m = TRUE # can also use equal (=) operator for
  assignment
```

There are some differences between `<-` and `=` mainly involving the scope. To simplify this, Google's R style guide simplifies the issue by prohibiting the `=` as an assignment operator. You can read more about this at:

<http://google-styleguide.googlecode.com/svn/trunk/google-r-style.html#assignment>

2.2.3 Numbers

Integers in R are generally treated as numeric objects (e.g. double precision real numbers like 3.1415926...). If you explicitly want an integer, you need to specify using the `L` suffix. For example, entering `1` gives you a numeric object while entering `1L` explicitly gives you an integer.

There is also a special number **Inf** which represents infinity (e.g. `1/0`). **Inf** can be used in ordinary calculations (e.g. `1/Inf` is `0`).

Keep in mind that, in general, all objects in R are case sensitive.

2.2.4 Dates and Times

Dates are represented by the date class. Time is represented by the POSIXct class, a very large integer. POSIXlt class is a list which stores a lot of useful metadata.

```
1 > x <- as.Date("2015-03-26")
2 > x

1 > x <- Sys.time()
2 > x
3 [1] "2015-03-23 21:28:10 PDT"
4 > p <- as.POSIXlt(x)
5 > names(unclass(p)) # unclass(p) is a list object
6 [1] "sec"      "min"      "hour"     "mday"     "mon"      "year"
7      "wday"    "yday"
8 [9] "isdst"    "zone"     "gmtoff"
9 > p$sec
10 [1] 10.89086
```

Utilize the **strptime** function in case your times are written in a different format as characters. For the **formatting strings**, check **?strptime** for details

```
1 > timeString <- "March 26, 2015 12:30"
2 > x <- strptime(timeString, "%B %d, %Y %H:%M")
3 > class(x)
4 [1] "POSIXlt" "POSIXt"
5 > x
6 [1] "2015-03-26 12:30:00 PDT"
```

2.2.5 Compound Objects

R has various compound objects. Some of these are vector, list, factor, matrix, and data frame.

2.2.6 Vector

The most basic of the compound objects is a vector. A vector can only contain objects of the same class.

Empty vectors can be created with the **vector()** function.

```
1 > x <- vector("numeric", length = 10)
2 > x
3 [1] 0 0 0 0 0 0 0 0 0 0
```

The `c()` function can be used to create vectors of objects.

```
1 > x <- c(0.5, 0.6) #number
2 > x <- c(TRUE, FALSE) # logical
3 > x <- c(T, F) #logical
4 > x <- c("a", "b", "c") #character
5 > x <- c(1+0i, 2+4i) #complex
```

The `:` operator is used to create integer sequences.

```
1 : operator
2 > x = 1:20
3 > x
4 [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
   19 20
```

2.2.7 List

List is a special type of vector. There are two characteristics of this type of vector.

1. It can contain elements of different classes (basic or compound).
2. Each element of a list can have a name.

Lists are a very important data type of R, so it is important to familiarize yourself with them.

```
1 > x <- list(1, "a", TRUE, 1 + 4i)
2 > x
3 [[1]]
4 [1] 1
5 [[2]]
6 [1] "a"
7 [[3]]
8 [1] TRUE
9 [[4]]
10 [1] 1+4i
```

```
1 > x <- list(a=c(T,T,F,F), b=2)
2 > x
3 $a
4 [1] TRUE TRUE FALSE FALSE
5
6 $b
7 [1] 2
```

2.2.8 Factor

Factors are used to represent categorical data and can be unordered or ordered. One can think of a factor as a numerical vector, where each integer has a label. Using factors with labels is better

than using integers. This is because factors are self-describing; having a variable with the values of "Male" and "Female" is better than a variable with values of 1 and 2.

```
1 > x <- factor(c("yes", "yes", "no", "yes", "no"))
2 > x
3 [1] yes yes no yes no
4 Levels: no yes
5 > table(x)
6 x
7 no yes
8 2 3
```

The order of the levels in a **factor()** can be set using the **levels** argument. This can be important in linear modeling, as the first level is used as the baseline.

```
1 > x <- factor(c("yes", "yes", "no", "yes", "no"))
2 > x
3 [1] yes yes no yes no
4 Levels: yes no
5 > x <- factor(x, levels=c("no", "yes"))
6 > x
7 [1] yes yes no yes no
8 Levels: no yes
```

2.2.9 Matrix

A matrix is a vector with a dimension attribute. The dimension attribute itself is an integer vector, with a length of 2 (nrow, ncol).

```
1 > m <- matrix(nrow = 2, ncol = 3)
2 > m
3      [,1] [,2] [,3]
4 [1,]   NA   NA   NA
5 [2,]   NA   NA   NA
```

```
1 > dim(m)
2 [1] 2 3
3 > attributes(m)
4 $dim
5 [1] 2 3
```

A matrix is constructed column-wise, meaning entries can be thought of as starting in the upper left corner and running down the columns.

```
1 > m <- matrix(1:6, nrow = 2, ncol = 3)
2 > m
3      [,1] [,2] [,3]
4 [1,]    1    3    5
5 [2,]    2    4    6
```

```
1 > m <- 1:10
2 > m
3 [1]  1  2  3  4  5  6  7  8  9 10
4 > dim(m) = c(2,5)
5 > m
6      [,1] [,2] [,3] [,4] [,5]
7 [1,]    1    3    5    7    9
8 [2,]    2    4    6    8   10
9 > n <- 1:10
10 > dim(n) = c(3,5)
11 Error in dim(n) = c(3, 5) :
12   dims [product 15] do not match the length of object [10]
```

A matrix can be created by column-binding or row-binding with **cbind()** and **rbind()**.

```
1 > dim(m)
2 > x <- 1:3
3 > y <- 10:12
4 > cbind(x,y)
5      x y
6 [1,] 1 10
7 [2,] 2 11
8 [3,] 3 12
```

```
1 > rbind(x,y)
2  [,1] [,2] [,3]
3 x     1     2     3
4 y    10    11    12
```

```
1 > m <- matrix(1:4, nrow = 2, ncol = 2)
2 > dimnames(m)
3 NULL
4 > dimnames(m) <- list(c("a", "b"), c("c", "d"))
5 > m
6   c d
7 A 1 3
8 b 2 4
```

2.2.10 Data Frames

Data frames are used to store tabular data. Unlike matrices which must have each element belonging to the same class, data frames can store different classes of objects in each column. Data frames also have a special attribute called (row.names). Finally, data frames are usually created by calling `read.table()` or `read.csv()`

```
1 > x <- data.frame(foo = 1:4, bar = c(T, T, F, F))
2 > x
3   foo bar
4 1    1 TRUE
5 2    2 TRUE
6 3    3 FALSE
7 4    4 FALSE
8 > nrow(x)
9 [1] 4
10 > ncol(x)
11 [1] 2
```

2.2.11 Coercion

When different objects are mixed in a vector, coercion occurs so that every element in the vector is of the same class.

```
1 > y <- c(1.7, "a") #character
2 > y
3 [1] "1.7" "a"
4 > y <- c(TRUE, 2) #numeric
5 > y
6 [1] 1 2
```

Objects can be explicitly coerced from one class to another using the `as.*` functions, when available.

```
1 > x <- 0:6
2 > class(x)
3 [1] "integer"
4
5 > as.numeric(x)
6 [1] 0 1 2 3 4 5 6
```

```
1 > as.logical(x)
2 [1] FALSE TRUE TRUE TRUE TRUE TRUE
3
4 > as.character(x)
5 [1] "0" "1" "2" "3" "4" "5" "6"
6
7 > as.complex()
8 complex(0)
9
10 > as.complex(x)
11 [1] 0+0i 1+0i 2+0i 3+0i 4+0i 5+0i 6+0i
```

```
1 > x <- c("a", "b", "c")
2 > as.numeric(x)
3 [1] NA NA NA
4 Warning message:
5 NAs introduced by coercion
6 > x
7 [1] "a" "b" "c"
8 > as.logical(x)
9 [1] NA NA NA
```

2.2.12 Missing Values

Missing values are denoted by **NA** or **NaN** for undefined mathematical operations. **NA** means a missing value and has various forms (e.g. `NA_integer`, `NA_character`, etc.). **NaN** means the value is not a number. **NaN** is also **NA** but the converse of this is not true.

```
1 > x <- c(1, 2, NA, 10, 3)
2 > is.na(x)
3 [1] FALSE FALSE TRUE FALSE FALSE
4 > is.nan(x)
5 [1] FALSE FALSE FALSE FALSE FALSE
6
7 > x <- c(1, 2, NaN, NA, 4)
8 > is.na(x)
9 [1] FALSE FALSE TRUE TRUE FALSE
10 > is.nan(x)
11 [1] FALSE FALSE TRUE FALSE FALSE
```

Data creation

c(...) generic function to combine arguments with the default forming a vector; with `recursive=TRUE` descends through lists combining all elements into one vector

from:to generates a sequence; “:” has operator priority; `1:4 + 1` is “2,3,4,5”

seq(from,to) generates a sequence `by=` specifies increment; `length=` specifies desired length

seq(along=x) generates `1, 2, ..., length(along)`; useful for for loops

rep(x,times) replicate `x` times; use `each=` to repeat “each” element of `x` each times; `rep(c(1,2,3),2)` is `1 2 3 1 2 3`; `rep(c(1,2,3),each=2)` is `1 1 2 2 3 3`

data.frame(...) create a data frame of the named or unnamed arguments; `data.frame(v=1:4, ch=c("a", "B", "c", "d"), n=10)`; shorter vectors are recycled to the length of the longest

list(...) create a list of the named or unnamed arguments; `list(a=c(1,2), b="hi", c=3i)`;

array(x,dim=) array with data `x`; specify dimensions like `dim=c(3,4,2)`; elements of `x` recycle if `x` is not long enough

matrix(x,nrow=,ncol=) matrix; elements of `x` recycle

factor(x,levels=) encodes a vector `x` as a factor

gl(n,k,length=n*k,labels=1:n) generate levels (factors) by specifying the pattern of their levels; `k` is the number of levels, and `n` is the number of replications

expand.grid() a data frame from all combinations of the supplied vectors or factors

rbind(...) combine arguments by rows for matrices, data frames, and others

cbind(...) id. by columns

Figure 2.1: Snippet from the R cheatsheet. Source: <http://cran.r-project.org/doc/contrib/Short-refcard.pdf>

Variable conversion

as.array(x), **as.data.frame(x)**, **as.numeric(x)**,
as.logical(x), **as.complex(x)**, **as.character(x)**,
 ... convert type; for a complete list, use **methods(as)**

Variable information

is.na(x), **is.null(x)**, **is.array(x)**, **is.data.frame(x)**,
is.numeric(x), **is.complex(x)**, **is.character(x)**,
 ... test for type; for a complete list, use **methods(is)**

length(x) number of elements in x

dim(x) Retrieve or set the dimension of an object; **dim(x) <- c(3,2)**

dimnames(x) Retrieve or set the dimension names of an object

nrow(x) number of rows; **NROW(x)** is the same but treats a vector as a one-row matrix

ncol(x) and **NCOL(x)** id. for columns

class(x) get or set the class of x; **class(x) <- "myclass"**

unclass(x) remove the class attribute of x

attr(x, which) get or set the attribute *which* of x

attributes(obj) get or set the list of attributes of obj

Figure 2.2: Snippet from the R cheatsheet. Source: <http://cran.r-project.org/doc/contrib/Short-refcard.pdf>

2.2.13 Exercise

To sum up what we learned in this section. R has basic data types and compound objects. It is also important to keep in mind the attributes of coercion and missing values in R.

In the following exercise we will declare two vectors with three elements each. Call these two vectors **x** and **y** and do the following:

1. Use **rbind()** and **cbind()** to create a matrix from these two variables. Observe the difference in structure resulting from these two functions.
2. Name the columns of the matrix.
3. Create a data frame from **x** and **y** and give the columns appropriate names.

```

1 > x <- c(1, 2, 3)
2 > y <- c(4, 5, 6)
3 > rbind(x, y)
4   [,1] [,2] [,3]
5 x    1    2    3
6 y    4    5    6

```

```

1 > cbind(x, y)
2       x y
3 [1,]  1 4
4 [2,]  2 5
5 [3,]  3 6
6 > a = cbind(x, y)
7 > colnames(a) <- c("x", "y")

```

2.3 Basic Operations

In this section we will be covering subsetting, control structures, build-in functions, user written functions, apply, packages, and the working directory and R script.

2.3.1 Subsetting

There are a number of operators that can be used to extract subsets of R objects. For example, `[]` always returns an object of the same class as the original object. With one exception, it can be used to select more than one element. `[[` is used to extract elements of a list or a data frame. It can only be used to extract a single element and the class of the returned object will not necessarily be a list or data frame. `$` is used to extract elements of a list or data frame by name. The semantics of `$` are similar to `[[`.

```

1 > x <- c("a", "b", "c", "c", "d", "a")
2 > x[2]
3 [1] "b"
4 > x[1:4]
5 [1] "a" "b" "c" "c"
6 > x > "a"
7 [1] FALSE TRUE  TRUE  TRUE  TRUE FALSE
8 > x[x>"a"]
9 [1] "b" "c" "c" "d"

```

```

1 > x <- matrix(1:6, 2, 3)
2 > x[1,2]
3 [1] 3
4 > x[1,] # Entire first row.
5 [1] 1 3 5
6 > x[,2] # Entire second column.
7 [1] 3 4

```


By default, when a single element of a matrix is retrieved, it is returned as a vector with a length of 1 rather than a 1 x 1 matrix. This behavior can be turned off by setting `drop` equal to `FALSE`.

```
1 > x <- matrix(1:6, 2, 3)
2 > x[1,2]
3 [1] 3
4 > x[1,2,drop=FALSE]
5      [,1]
6 [1,]    3
7 > x[1,]
8 [1] 1 3 5
```

Similarly, subsetting a single column or a single row will give you a vector, not a matrix by default.

```
1 > x[1,,drop=FALSE]
2      [,1] [,2] [,3]
3 [1,]    1    3    5
```

```
1 > x <- list(foo = 1:4, bar = 0.6)
2 > x[1]
3 $foo
4 [1] 1 2 3 4
5 > x$foo
6 [1] 1 2 3 4
7 > x$bar
8 [1] 0.6
9 > x["bar"]
10 $bar
11 [1] 0.6
12 > x[["bar"]]
13 [1] 0.6
```

Below is an example of extracting multiple elements of a list:

```
1 > x <- list(foo = 1:4, bar = 0.6, baz = "hello")
2 > x[c(1, 3)]
3 $foo
4 [1] 1 2 3 4
5
6 $baz
7 [1] "hello"
```

Partial matching of names is allowed with `[]` and `$`

```
1 > x <- list(addedName = 1:5)
2 > x$a
3 [1] 1 2 3 4 5
4 > x[["a"]]
5 NULL
6 > x[["a", exact=FALSE]]
7 [1] 1 2 3 4 5
```

2.3.2 Control Structures

Conditional:

```
1 if (logical.expression) {
2   statements
3 }
4 else if (another.logical.expression) {
5   statements
6 }
7 else {
8   alternative.statements
9 }
10 else if, else branch is optional
```

for-loop:

```
1 for(i in 1:10) {
2   print(i*i)
3 }
```

repeat-loop:

```
1 i=1
2 repeat {
3   i= i+1
4   if (i>=10) break # the only
5   # way to get out of the loop
6 }
```

while-loop:

```
1 i=1
2 while(i<=10) {
3   print(i*i)
4   i=i+sqrt(i)
5 }
```

R can support looping however, keep in mind:

1. Looping is not recommended for compound objects like lists, vectors, or data.frames, etc. For example, a loop could be used to multiply all elements by 2, but the loop is implicit.

```
1 > x <- c(1, 2, 3)
2 > x * 2
3 [1] 2 4 6
```

2. Control structures mentioned here are primarily useful for writing programs. For command-line interactive work, the ***apply** functions (lapply, sapply, etc.) are more useful.

2.3.3 Built-in Functions

Numeric Functions

```
1 abs(x) absolute value
2 sqrt(x) square root
3 ceiling(x) ceiling(3.475) is 4
4 floor(x) floor(3.475) is 3
5 trunc(x) trunc(5.99) is 5
6 round(x, digits=n) round(3.475, digits=2) is 3.48
7 ...
```

Character Functions

```
1 paste(..., sep="") Concatenate strings after using sep
  string to separate them.
2 paste("x",1:3,sep="") returns c("x1","x2" "x3")
3 paste("x",1:3,sep="M") returns c("xM1","xM2" "xM3")
4
5 grep(pattern, x , ignore.case=FALSE, fixed=FALSE)
6
7 Search for pattern in x. If fixed =FALSE then pattern is a
  regular expression. If fixed=TRUE then pattern is a
  text string. Returns matching indices.
8 grep("A", c("b","A","c"), fixed=TRUE) returns 2
9 ...
```

Statistical Functions

```
1 rnorm(), dunif(), mean(), sum()
```

You can find a more complete list of other built-in functions at: <http://www.statmethods.net/management/functions.html>

2.3.4 User-written Functions

```
1 foo <- function(x,y=1,...) {  
2   cat("extra args:",..., "\n")  
3   sqrt(x)+sin(y)  
4 }  
5 foo(1,2,3, "bar")  
6 foo(1)  
7 foo(y=3,x=7) # named arguments in any order
```

Type the function name to see the code. This works for any function. `foo`



The three dots used in the code allow for an arbitrary number and variety of arguments. They also allow for the passing of arguments on to other functions.

2.3.5 str Function

str compactly displays the internal structure of an R object (data or function). It is a diagnostic function that acts as an alternative to **summary()** for a data object. **str** is especially well suited to compactly display the abbreviated contents of a possibly nested list. It also tells you what is in the object and the arguments of the function.

Try out: **str(lm)**

2.3.6 apply

lapply loops over a list and evaluates a function for each element. **sapply** is the same as **lapply** and tries to simplify the result.

```
1 > lapply(c(-1, -5), abs)  
2 [[1]]  
3 [1] 1  
4  
5 [[2]]  
6 [1] 5
```

```
1 > sapply(c(-1, -5), abs)  
2 [1] 1 5
```

apply applies a function over the margins of an array.

```
1 > x <- matrix(c(1, 2, 3, 4), 2, 2)
2 > x
3      [,1] [,2]
4 [1,]    1    3
5 [2,]    2    4
6 > apply(x, 1, sum)
7 [1] 4 6
8 > apply(x, 2, sum)
9 [1] 3 7
```

tapply applies a function over the subsets of a vector.

```
1 > score <- c(90, 79, 94, 85)
2 > gender <- factor(c("Male", "Female", "Female", "Male"))
3 > tapply(score, gender, mean)
4 Female      Male
5   86.5     87.5
```

split is an auxiliary function. A common idiom is split followed by lapply

```
1 > score <- c(90, 79, 94, 85)
2 > gender <- factor(c("Male", "Female", "Female", "Male"))
3 > splitted <- split(score, gender)
4 > splitted
5 $Female
6 [1] 79 94
7 $Male
8 [1] 90 85
9 > lapply(splitted, mean)
10 $Female
11 [1] 86.5
12 $Male
13 [1] 87.5
```

mapply is a multi-variable input. It applies a function in parallel over a set of arguments.

```
1 > mapply(rep, c(0,2), c(3,5)) # input certain combinations
2 ## <- list(rep(0,3), rep(2,5)), it vectorizes a function
3 [[1]]
4 [1] 0 0 0
5
6 [[2]]
7 [1] 2 2 2 2 2
```

Don't worry if you forget the logics of **apply** functions. You can always use **lapply**, **str(lapply)**, or Google/Bing to check.

```
1 > str(lapply)
2 function (X, INDEX, FUN = NULL, ..., simplify = TRUE)
```

2.3.7 Packages

For almost everything you want to do with R, there is probably a package written to do just that. A list of packages can be found in the official repository CRAN: <http://cran.fhcrc.org/web/packages/>. If you need a package, it can be installed very easily from within R using the command:

```
1 install.packages("packagename") # if package already
   installed, it'll bypass
```

Libraries in Github can be installed using the devtools library.

2.3.8 Working directory and R script

You can set the working directory from the menu if using the R-gui (Change dir...) or from the R command line:

```
1 setwd("C:\\MyWorkingDirectory")
2 setwd("C:/MyWorkingDirectory") # can use forward slash
3 setwd(choose.dir()) # opens a file browser
4
5 getwd() # returns a string with
6         # the current working directory
```

To see a list of the files in the current directory:

```
1 dir() # returns a list of strings of file names
2 dir(pattern=".R$") # list of files ending in ".R"
3 dir("C:\\Users") # show files in directory C:\\Users
```

Run a script:

```
1 source("helloworld.R") # execute a script
```

2.4 Reading and Writing Data

In this section we will be covering reading and writing local flat files, reading and writing local Excel files, connection interfaces, reading XML/HTML files, reading and writing to JSON, connecting to a database, and the textual format.

2.4.1 Reading and Writing Local Flat Files

You can read local flat files as **read.table**, **read.csv**, **readLines**. CSV stands for “Comma Separated Values.”

For this lab we will be working with the Titanic dataset. To procure this dataset, please visit the following link:

<https://raw.githubusercontent.com/datasciencedojo/datasets/master/titanic.csv>

Right click and select "save as" to save it to your local directory. The code samples in the remainder of this chapter assumes that you have set your working directory to the location of the titanic.csv file.

```
1 > titanic_data <- read.csv("titanic.csv")
2 > head(titanic_data, 3)
3   X Class  Sex   Age Survived Freq
4 1 1   1st Male Child       No    0
5 2 2   2nd Male Child       No    0
6 3 3   3rd Male Child       No   35
```

```
1 > line1 <- readLines("titanic.csv", 1)
2 > line1
3 [1] "\"\"\", \"Class\", \"Sex\", \"Age\", \"Survived\", \"Freq\""
```

You can write local flat files as **write.table**, **write.csv**, **writeLines**.

```
1 > write.table(titanic_data, "new_Titanic.csv")
```


Before reading or writing files, make sure to take note of the parameters. For example the **read.table** function is one of the most commonly used functions for reading data. It has a few important arguments to take into account:

- **file** is the name of a file, or a connection.
- **header** a logical value indicating the variables of the first line.
- **sep** is a string indicating how the columns are separated.
- **colClasses** is a character vector indicating the class of each column in the dataset.
- **nrows** are the number of rows in the dataset
- **comment.char** is character string indicating the comment character .
- **skip** is the number of lines to skip from the beginning.
- **stringsAsFactors** determines whether character variables be coded as factors.

2.4.2 Reading and Writing Excel Files

read.xlsx, write.xlsx, or read.xlsx2, write.xlsx2 (faster, but unstable)

```
1 # Install the xlsx library
2 > install.packages('xlsx')
3 # Load the library
4 > library(xlsx)
5 # Now you can Read the Excel file
6 > titanic_data <- read.xlsx("titanic3.xls", sheetIndex=1)
```

2.4.3 Connection Interfaces

In practice, we don't need to deal with the connection interface directly. Connections can be made to files or to other channels:

- **file** opens a connection to a file
- **gzfile** opens a connection to a file compressed with gzip
- **bzfile** opens a connection to a file compressed with bzip2
- **url** opens a connection to a webpage

Here is a simple example:

```
1 > con <- file("Titanic.csv", "r")
2 > titanic_data <- read.csv(con)
3 > close(con)
4
5 > con <- url("http://vincentarelbundock.github.io/
  Rdatasets/csv/datasets/Titanic.csv")
6 > another_data <- read.csv(con)
```

2.4.4 Reading XML/HTML Files

```

1 > library(XML) # you need to install this library
2 > url <- "http://www.w3schools.com/xml/simple.xml"
3 > doc <- xmlTreeParse(url, useInternal=TRUE) # also works
  for html
4 > rootNode <- xmlRoot(doc)
5 > rootNode[[1]]
6 <food>
7   <name>Belgian Waffles</name>
8   <price>$5.95</price>
9   <description>Two of our famous Belgian Waffles with
  plenty of real maple syrup</description>
10  <calories>650</calories>
11 </food>

```

```

1 > rootNode[[1]][[1]]
2 <name>Belgian Waffles</name>
3 > xpathSApply(rootNode, "//name", xmlValue)
4 [1] "Belgian Waffles" "Strawberry Belgian Waffles" "Berry-
  Berry Belgian Waffles" "French Toast" "Homestyle
  Breakfast"

```

- **xpath: /node:** top level node; **//node:** node at any level;
- **node[@attr-name]:** node with an attribute name;
- **node[@attr-name = 'bob']:** node with an attribute name = 'bob'

2.4.5 Reading and Writing JSON

JSON or JavaScript Object Notation or "JSON" is a common format for data from application programming interfaces ("APIs"). It is an alternative to XML.

```

1 > library(jsonlite)
2 > jsonData <- fromJSON("http://citibikenyc.com/stations/
  json")
3 > names(jsonData)
4 [1] "executionTime" "stationBeanList"
5 > jsonData$stationBeanList[1,1:3]
6   id      stationName  availableDocks
7 1 72 W 52 St & 11 Ave      31

```

```
1 > head(iris, 3)
2   Sepal.Length Sepal.Width Petal.Length Petal.Width
3   Species
4 1          5.1          3.5          1.4          0.2
5   setosa
6 2          4.9          3.0          1.4          0.2
7   setosa
8 3          4.7          3.2          1.3          0.2
9   setosa
10
11 > iris2 <- toJSON(iris, pretty=TRUE)
```

2.4.6 Connect to a Database

JSONPackages: RMySQL, RpostresSQL, RODBC, RMONGO

```
1 > library(RMySQL)# load the library
2 > ucscDb <- dbConnect(MySQL(), user="genome", host="genome
3   -mysql.cse.ucsc.edu")
4 > data <- dbGetQuery(ucscDb, "show databases;") # get the
5   output of SQL query as data frame in R
6 > head(data)
7       Database
8 1 information_schema
9 2          ailMel1
10 3          allMis1
11 4          anoCar1
12 5          anoCar2
13 6          anoGam1
14 > dbDisconnect(ucscDb) # don't forget to close the
15   connection
```

2.5 Statistical Simulation

Probability distribution functions usually have four functions associated with them. The functions are prefixed with:

- **d** for density, like `dnorm()`
- **r** for random number generation, like `rnorm()`
- **p** for cumulative distribution, like `pnorm()`
- **q** for quantile function, like `qnorm()`

Below is an example of normal distribution:

```
1 > dnorm(0, mean = 0, sd = 1)
2 [1] 0.3989423
3 > dnorm(10, mean = 0, sd = 1)
4 [1] 7.694599e-23
5
6 > pnorm(0, mean = 0, sd = 1)
7 [1] 0.5
8 > pnorm(1, mean = 0, sd = 1)
9 [1] 0.8413447
10 > pnorm(100, mean = 0, sd = 1)
11 [1] 1
```

```
1 > qnorm(0.5, mean = 0, sd = 1)
2 [1] 0
3 > qnorm(0, mean = 0, sd = 1)
4 [1] -Inf
5 > qnorm(0.2, mean = 0, sd = 1)
6 [1] -0.8416212
7
8 > rnorm(4, mean = 0, sd = 1)
9 [1] -0.80938783 -0.07203091 0.99059330 0.69783570
```

Built-in standard distributions: Normal (**norm**), Poisson (**pois**), Binomial (**binom**), Exponential (**exp**), Gamma (**gamma**), Uniform (**unif**)

Reproducible random number: **set.seed()** in advance

```
1 > rnorm(3, mean = 0, sd = 1)
2 [1] -1.1991719 -0.3227133  0.4802463
3 > set.seed(1)
4 > rnorm(3, mean = 0, sd = 1)
5 [1] -0.6264538  0.1836433 -0.8356286
6 > set.seed(1)
7 > rnorm(3, mean = 0, sd = 1)
8 [1] -0.6264538  0.1836433 -0.8356286
```

2.6 Plotting Systems (optional)

In this chapter we will be covering a few of the different kinds of plotting systems.

- **Base graphics** are constructed piecemeal. This system is conceptually simpler and allows plotting to mirror the thought process.
- **Lattice graphics** are entire plots created through a simple function call.
- **ggplot2 graphics** are an implementation of the Grammar of Graphics by Leland Wilkinson. They combine concepts from both base and lattice graphics. To use this system you must install the ggplot2 library.

Fancier and more telling ones will be elaborated on further during the bootcamp! In the meantime, check out a list of different interactive visualizations in R: <http://ouzor.github.io/blog/2014/11/21/interactive-visualizations.html>

2.6.1 Basic Plotting Systems

```
1 > library(datasets)
2 > names(airquality)
3 [1] "Ozone" "Solar.R" "Wind" "Temp" "Month" "Day"
4 > plot(x = airquality$Temp, y = airquality$Ozone)
```

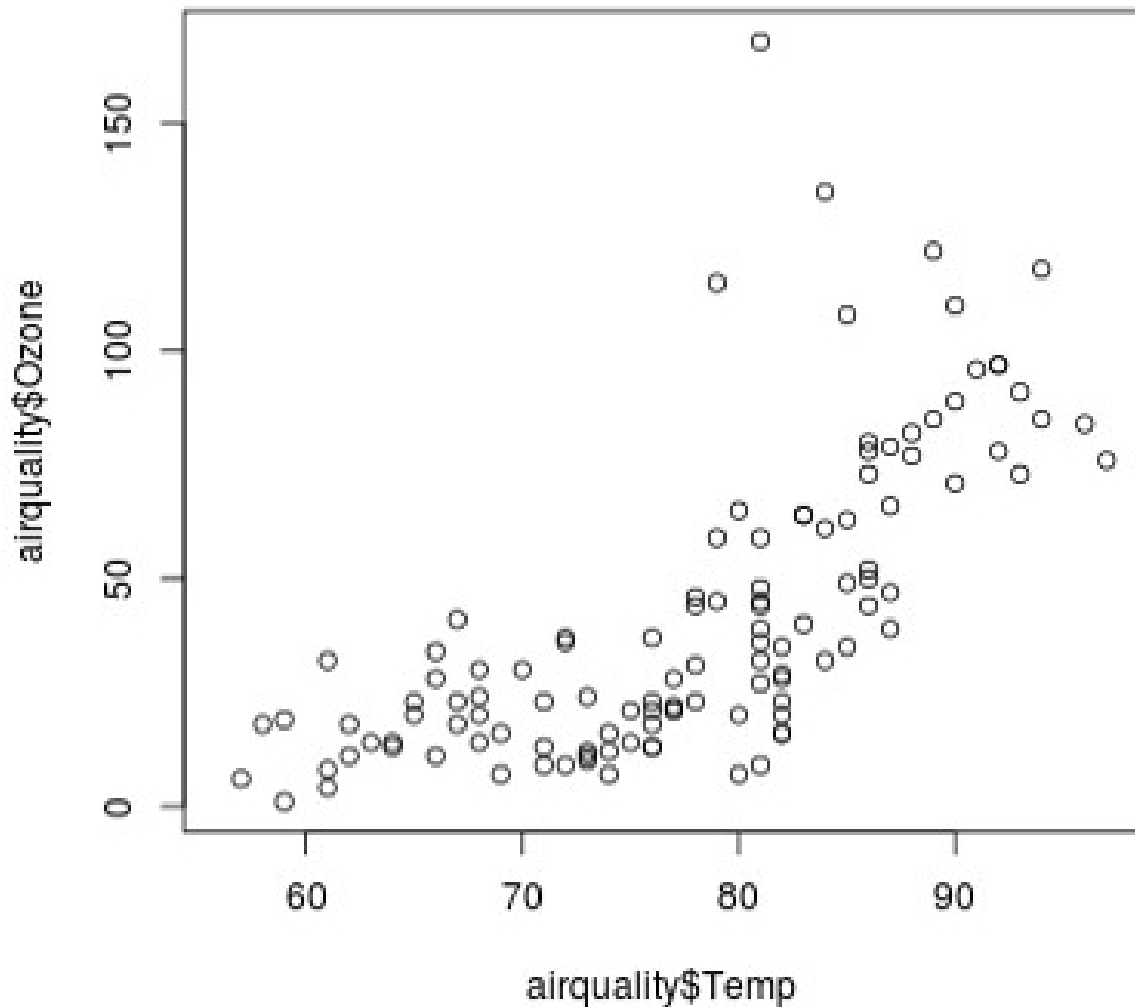


Figure 2.3: Scatter plot of temperature and ozone level

```

1 # par() function is used to specify global graphics
  parameters
2 # that affect all plots in an R session. Type ?par to see
  all
3 # parameters
4 > par(mfrow = c(1, 2), mar = c(4, 4, 2, 1), oma = c(0, 0,
  2, 0))
5 > with(airquality, {
6 + plot(Wind, Ozone, main="Ozone and Wind")
7 + plot(Temp, Ozone, main="Ozone and Temperature")
8 + mtext("Ozone and Weather in New York City", outer=TRUE)
9 + })

```

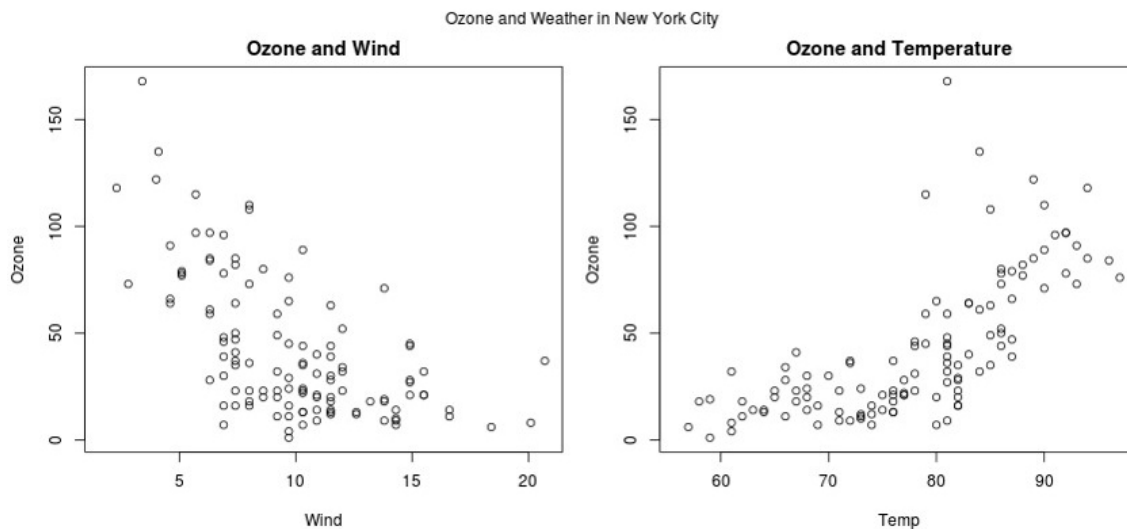


Figure 2.4: Correlation between ozone and weather conditions in New York City

There are various plotting functions that can be utilized.

- **lines** add lines to a plot, given a vector of x values and corresponding vector of y values.
- **points** adds a point to the plot.
- **text** adds text labels to a plot using specified x,y coordinates.
- **title** adds annotations to x,y axis labels, title, subtitles, and the outer margin.
- **mtext** adds arbitrary text to the inner or outer margins of the plot.
- **axis** specifies the axis ticks

2.6.2 Lattice Plotting Systems

```
1 > library(lattice) # need to load the lattice library
2 > set.seed(10) # set the seed so our plots are the same
3 > x <- rnorm(100)
4 > f <- rep(1:4, each = 25) # first 25 elements are 1,
   # second 25 elements are 2, ...
5 > y <- x + f - f * x + rnorm(100, sd = 0.5)
6 > f <- factor(f, labels = c("Group 1", "Group 2", "Group 3",
   # "Group 4"))
7 > xyplot(y ~ x | f)
```

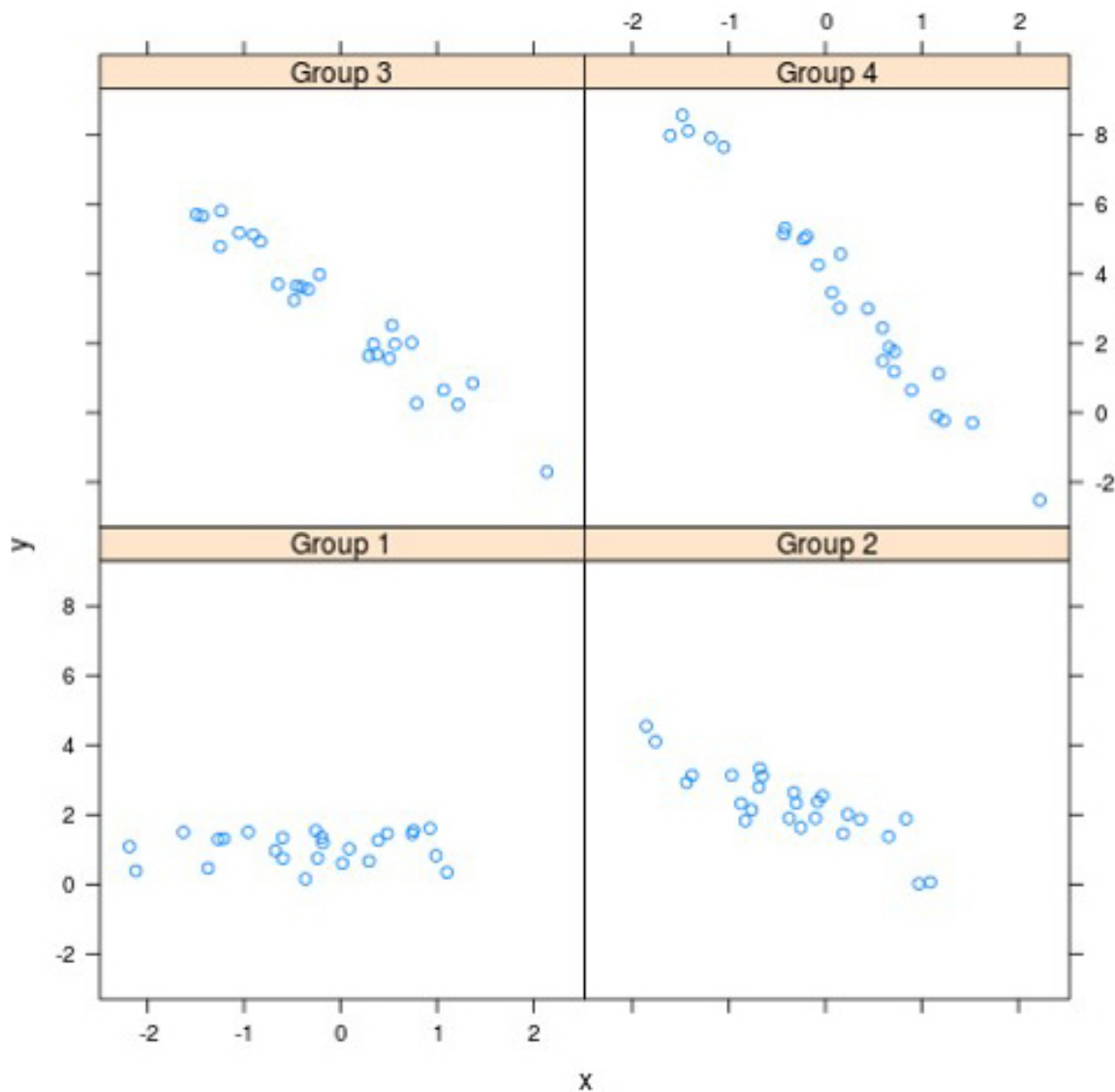


Figure 2.5: Example of simulated data by factor level

What more can you do with the plots?

```

1 > xyplot(y ~ x | f, panel = function(x, y, ...) {
2   # call the default panel function for xyplot
3   panel.xyplot(x, y, ...)
4   # adds a horizontal line at the median
5   panel.abline(h = median(y), lty = 2)
6   # overlays a simple linear regression line
7   panel.lmline(x, y, col = 2)
8 })

```

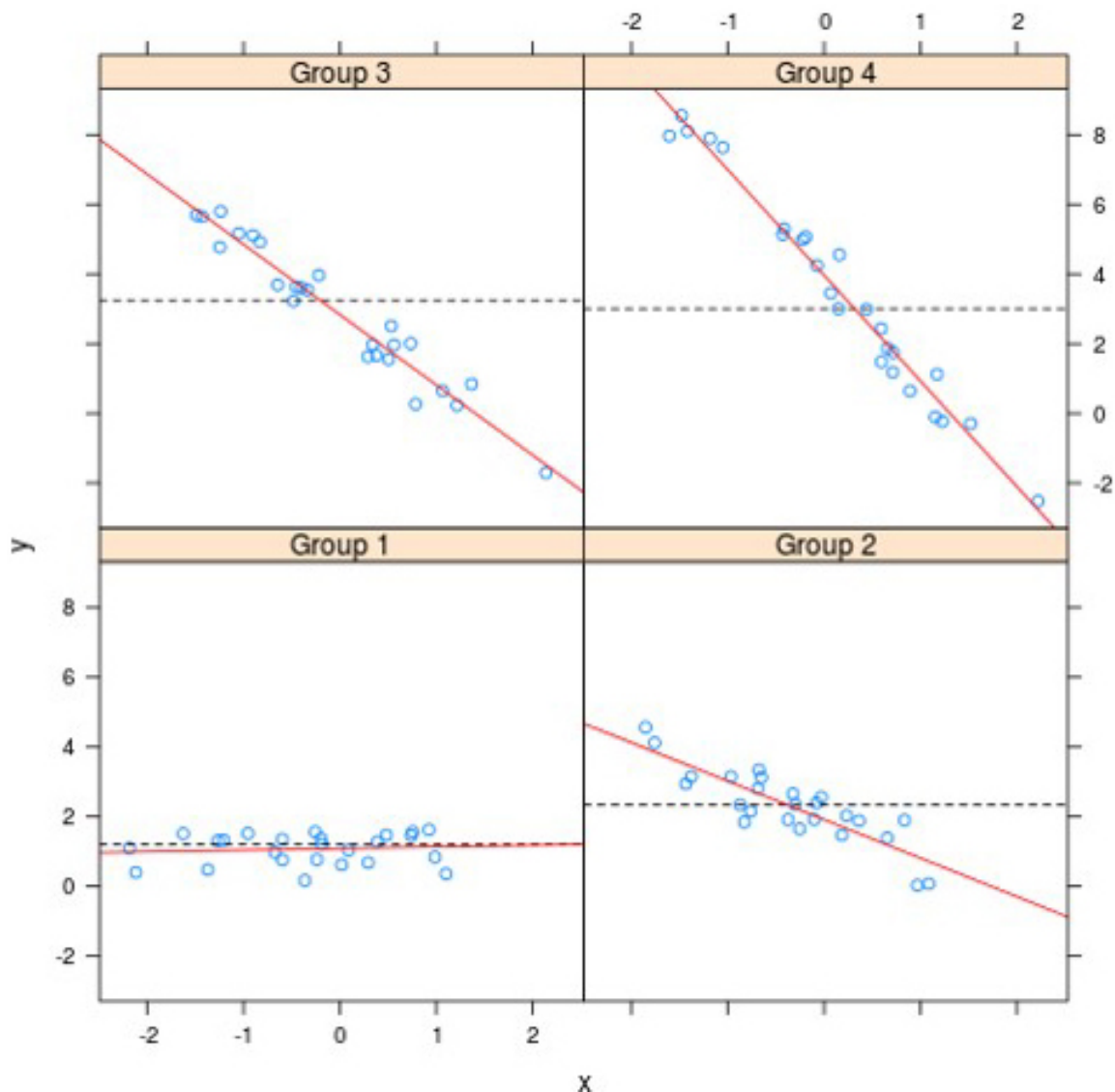


Figure 2.6: Adding linear trend lines to scatter plot

Plotting Functions:

xyplot() is the main function for creating scatterplots **bwplot()** is a box and whiskers plot or box plot **histogram()** is a histogram **stripplot()** is a box plot with actual points **dotplot()** is a plot with dots on “violin strings” **splo**m() is scatterplot matrix, similar to **pairs()** in the base plotting system **levelplot()**/**contourplot()** plots image data

2.6.3 ggplot2 Plotting systems

Before using ggplots, you must install the ggplot2 library. Ggplots mix the elements of base and lattice systems. A good tutorial on the three basic plotting systems can be found at:

https://sux13.github.io/DataScienceSpCourseNotes/4_EXDATA/Exploratory_Data_Analysis_Course_Notes.html.

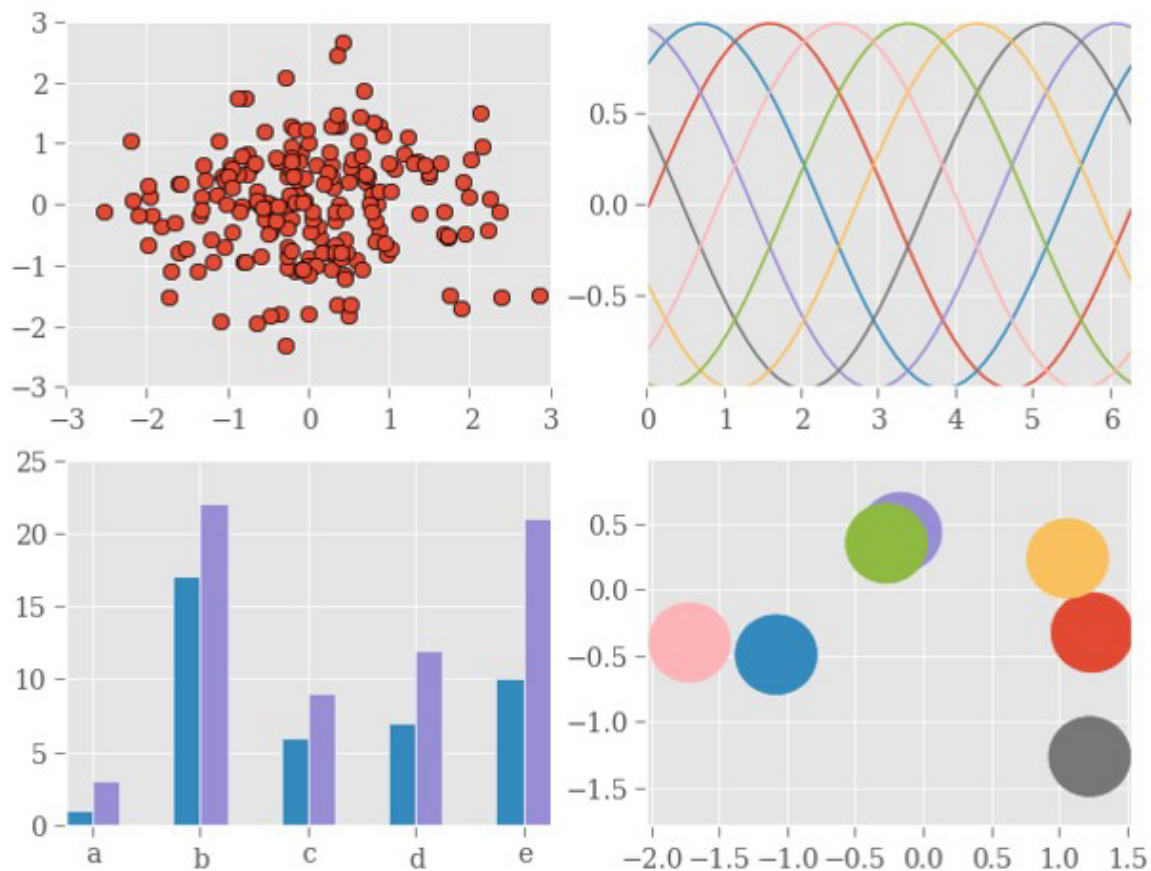


Figure 2.7: Example demonstrations of ggplots. Source: http://matplotlib.org/examples/style_sheets/plot_ggplot.html