



# DATA EXPLORATION USING R

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 : 1.

Registration Number: TX0008030133

2015-05-22

Printed in the United States of America

*Copyright © 2014-2015*

# Contents

<b>1</b>	<b>Appendix B: Data Exploration Using R .....</b>	<b>5</b>
1.1	<b>Overview</b>	<b>5</b>
1.2	<b>Getting the Data</b>	<b>5</b>
1.3	<b>Data Exploration</b>	<b>6</b>
1.4	<b>Categorical Preprocessing</b>	<b>8</b>
1.5	<b>Data Exploration Through Visualization</b>	<b>8</b>
1.5.1	Pie Charts .....	8
1.5.2	Box Plot Distributions .....	11
1.5.3	Histogram .....	13
1.5.4	Kernel Density Plot .....	13
1.6	<b>Feature Engineering</b>	<b>15</b>
1.7	<b>Advanced Visualizations in R</b>	<b>17</b>



# 1. Appendix B: Data Exploration Using R

## 1.1 Overview

## 1.2 Getting the Data

- Download the data set “Titanic\_train.csv” at  
[https://github.com/datasciencedojo/bootcamp/tree/master/Datasets/Titanic\\_train.csv](https://github.com/datasciencedojo/bootcamp/tree/master/Datasets/Titanic_train.csv)
- Set working directory of R to the directory of the file using setwd()

```
1 # Reads in the Titanic dataset and sets it to a variable
2 # called 'Titanic'.
3 titanic = read.csv("C:\\\\myFolder\\\\titanic.csv");
4
5 colnames(titanic)
6 # "PassengerId" "Survived" "Pclass" "Name"
7 # "Sex"          "Age"      "SibSp"   "Parch"
8 # "Ticket"       "Fare"     "Cabin"   "Embarked"
9
10 # revise the name of a column
11 colnames(titanic)[5] <- "Gender"
```

### 1.3 Data Exploration

Look at the first few rows.

```

1 head(titanic[,1:5], 4)
2 ##   PassengerId Survived Pclass
3 ## 1             1       0     3
4 ## 2             2       1     1
5 ## 3             3       1     3
6 ## 4             4       1     1
7 ##
8 ##                               Name
9 ## 1                         Braund, Mr. Owen Harris
10 ## 2 Cumings, Mrs. John Bradley (Florence Briggs Thayer)
11 ## 3                         Heikkinen, Miss. Laina
12 ## 4      Futrelle, Mrs. Jacques Heath (Lily May Peel)
13 ## 5                         Allen, Mr. William Henry
14 ## 6                         ...
```

Let's look at a more specific sample of the data.

```

1 # View the last six columns of the data frame, and the
2 # first 4 rows in those columns.
3 head(titanic[,6:ncol(titanic)], 4)
4
5 #Age SibSp Parch          Ticket    Fare Cabin Embarked
6 #1 22    1    0            A/5 21171 7.2500   S
7 #2 38    1    0            PC 17599 71.2833   C85   C
8 #3 26    0    0  STON/O2. 3101282 7.9250   S
9 #4 35    1    0            113803 53.1000  C123   S
```

Summary of the data frame using str() function.

```
1 str(titanic)
2
3 ##  'data.frame':   891 obs. of  12 variables:
4 ##   $ PassengerId: int  1 2 3 4 5 6 7 8 9 10 ...
5 ##   $ Survived    : int  0 1 1 1 0 0 0 0 1 1 ...
6 ##   $ Pclass      : int  3 1 3 1 3 3 1 3 3 2 ...
7 ##   $ Name        : Factor w/ 891 levels "Abbing , Mr.
8 ##                 Anthony",...: 109 191 35..
9 ##   $ Gender      : Factor w/ 2 levels "female","male": 2 1
10 ##                 1 1 2 2 2 2 1 1 ...
11 ##   $ Age         : num  22 38 26 35 35 NA 54 2 27 14 ...
12 ##   $ SibSp       : int  1 1 0 1 0 0 0 3 0 1 ...
13 ##   $ Parch       : int  0 0 0 0 0 0 0 1 2 0 ...
14 ##   $ Ticket      : Factor w/ 681 levels
15 ##                 "110152","110413",...: 524 597 670 50..
```

## 1.4 Categorical Preprocessing

There are some values that must be casted into categorical features so that our visualizations do not treat them as numerical values.

```

1 # Converts survived into a factor
2 titanic$Survived = factor(titanic$Survived)
3
4 # Converts embarked into a factor
5 titanic$Embarked = factor(titanic$Embarked)
6
7 # Summarizes the newly casted data frame.
8 str(titanic$Survived)
9 ##  Factor w/ 2 levels "0","1": 1 2 2 2 1 1 1 1 2 2 ...
10 str(titanic$Embarked)
11 ##  Factor w/ 4 levels "", "C", "Q", "S": 4 2 4 4 4 3 4 4 4 2
    ...

```

To remove confusion, rename the categorical values themselves from numeric zeros and ones into their abstracted values.

```

1 levels(titanic$Survived) <- c("died", "survived")
2 levels(titanic$Embarked) <- c("Unknown", "Cherbourg",
  "Queenstown", "Southampton")
3 str(titanic$Survived)
4 ##  Factor w/ 2 levels "died","survived": 1 2 2 2 1 1 1 1
  2 2 ...
5 str(titanic$Embarked)
6 ##  Factor w/ 4 levels "Unknown","Cherbourg",...: 4 2 4 4 4
  3 4 4 4 2 ...

```

## 1.5 Data Exploration Through Visualization

### 1.5.1 Pie Charts

What is the rate of survival overall? Let's visualize the distribution of survival in a pie chart. (Figure: 1.1)

```

1 survivedTable = table(titanic$Survived); survivedTable
2 ##
3 ##      died   survived
4 ##      549       342
5 # mar, and oma mean margin and outer margin of the plot,
6 # cex gives the size of texts
7 par(mar = c(0, 0, 0, 0), oma = c(0, 0, 0, 0), cex=1.5)
8 pie(survivedTable, labels=c("Died", "Survived"))

```

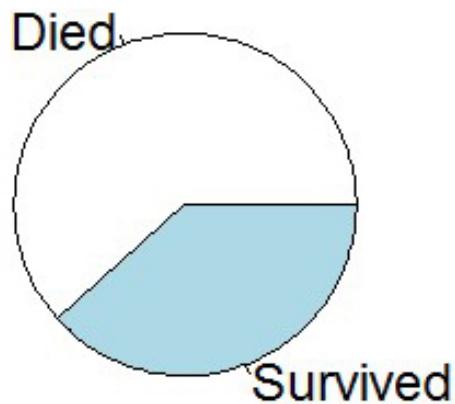


Figure 1.1

Was gender a good predictor? Let's visualize pie chart distributions on survival rate based upon gender. (Figure: 1.2)

```
1 par(mfrow = c(1, 2), mar = c(0, 0, 2, 0), oma = c(0, 1, 0,
   1), cex=1.5)
2 pie(table(titanic[titanic$Gender=="male", "Survived"]),
3   labels=c("Dead","Survived"), main="Survival Portion
   of Men")
4 pie(table(titanic[titanic$Gender=="female", "Survived"]),
5   labels=c("Dead","Survived"), main="Survival Portion of
   Women")
```

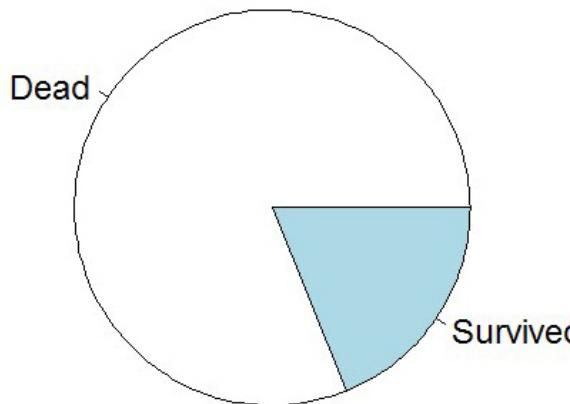
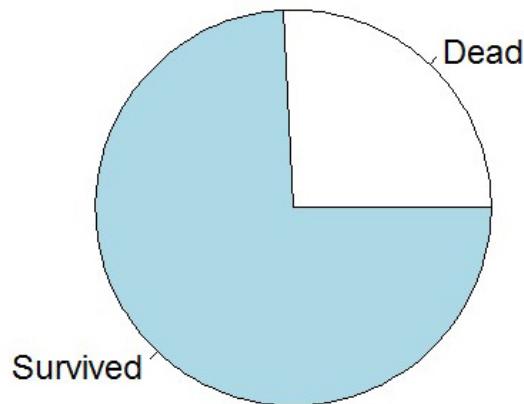
**Survival Portion of Men****Survival Portion of Women**

Figure 1.2

What is the distribution of age among titanic passengers?

```
1 # isolate the Age column
2 Age <- titanic$Age
3
4 #summary stats for the age coulmn
5 summary(Age)
6 ##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.    NA's
7 ##      0.42   20.12  28.00    29.70  38.00    80.00    177
```

What was the age statistics for those who survived versus those that died?

```
1 summary(titanic[titanic$Survived=="died", "Age"])
2 ##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.    NA's
3 ##      1.00   21.00  28.00    30.63  39.00    74.00    125
4
5 summary(titanic[titanic$Survived=="survived", "Age"])
6 ##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.    NA's
7 ##      0.42   19.00  28.00    28.34  36.00    80.00     52
```

### 1.5.2 Box Plot Distributions

When you want to compare the distribution or find patterns between two attributes where one is a numeric and the other is a categorical, a comparison box plot is a good option.

```
1 #Comparing age (numeric), with gender (categorical)
2 boxplot(
3   titanic$Age~titanic$Sex, #comparison attributes
4   data=titanic, #data source
5   main="Age Distribution by Gender", #title
6   xlab="Gender", #x-axis label
7   ylab="Age", #y-axis label
8   col=c("red", "green") #color
9 )
```

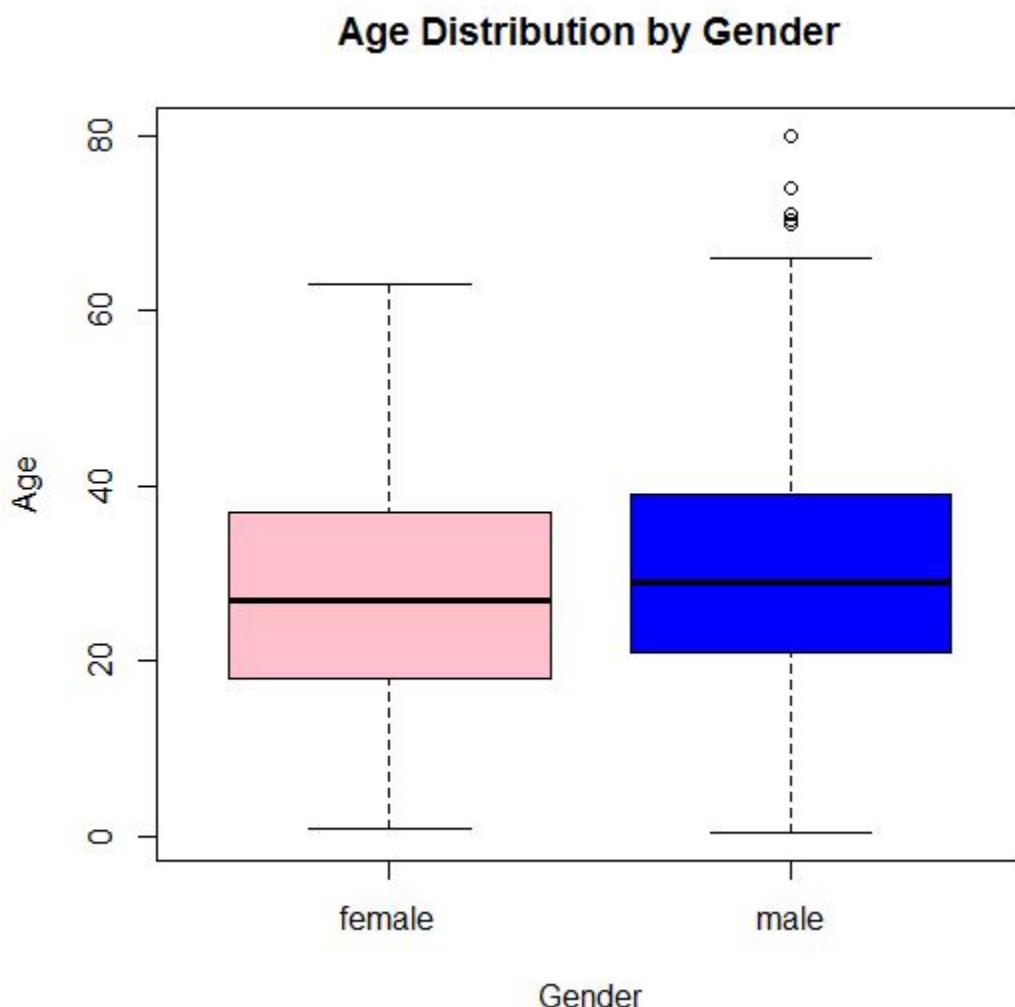


Figure 1.3: Comparison box plot between age and gender

```
1 # Comparing age (numeric) with survived (categorical)
2 boxplot(
3   titanic$Age~titanic$Survived, #comparison attributes
4   data=titanic, #data source
5   main="Age Distribution by Survival", #title
6   xlab="Survived", #x-axis label
7   ylab="Age", #y-axis label
8   col=c("blue", "red") #color
9 )
```

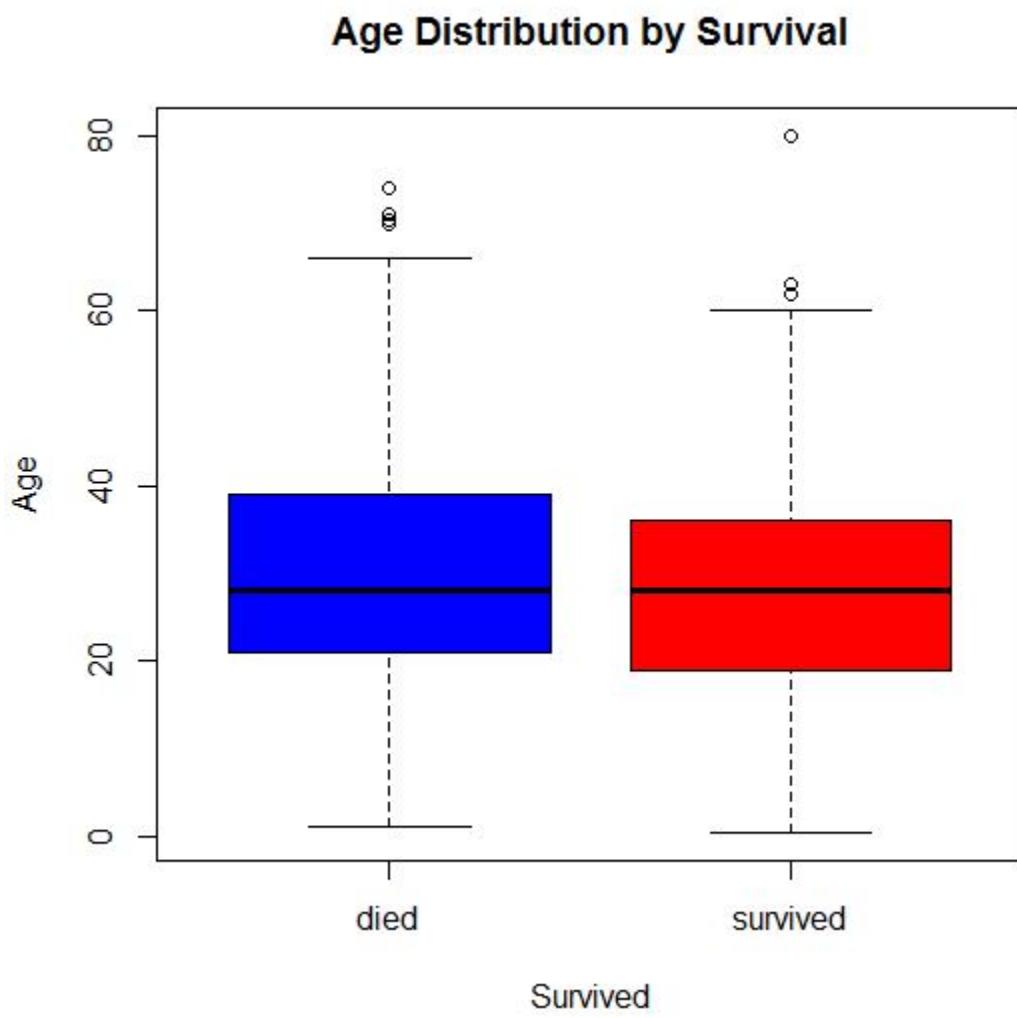


Figure 1.4: Comparison box plot between age and survived

### 1.5.3 Histogram

Let's visualize the distribution of age in a histogram. (Figure: 1.5)

```
1 hist(Age, col="blue", xlab="Age", ylab="Frequency",
2       main = "Distribution of Passenger Ages on Titanic",
3       cex.lab=1.6, cex.axis=1.4, cex.main=1.6)
```

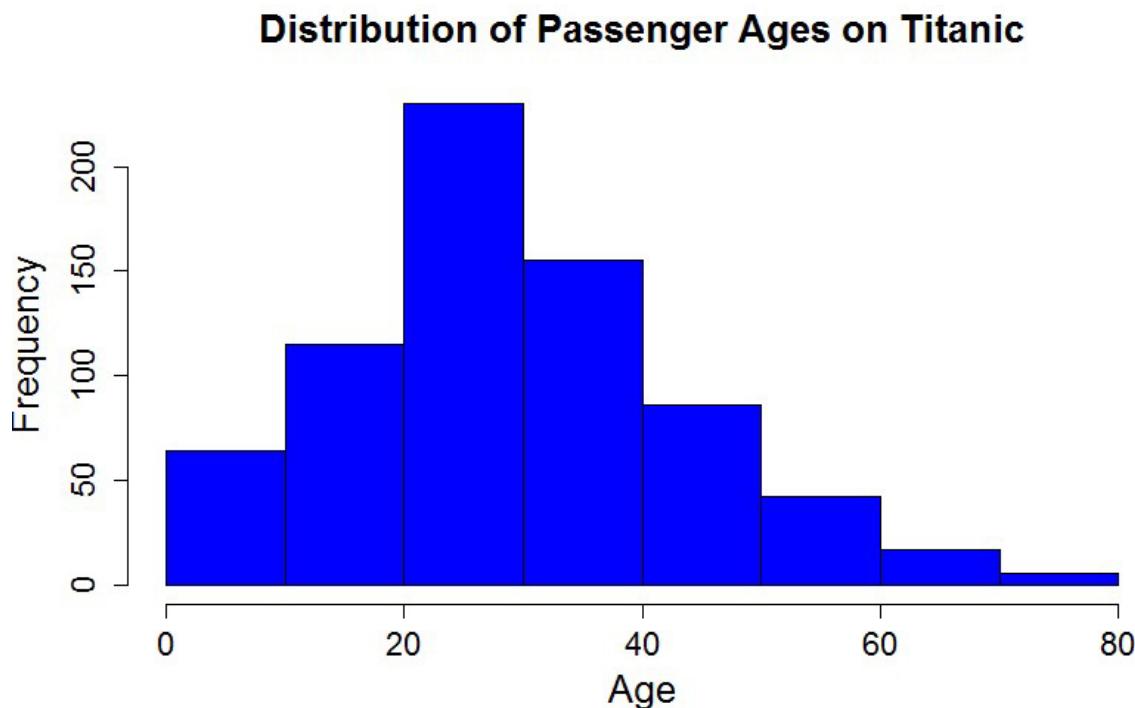


Figure 1.5

### 1.5.4 Kernel Density Plot

Let's smoothen the histogram into a kernel density plot. (Figure: 1.6)

```
1 # density() requires all NAs to be removed
2 d = density(na.omit(Age))
3 plot(d, main = "kernel density of Ages of Passengers",
4       xlab="Age", cex.lab=1.6, cex.axis=1.4)
5 polygon(d, col="red", border="blue")
```

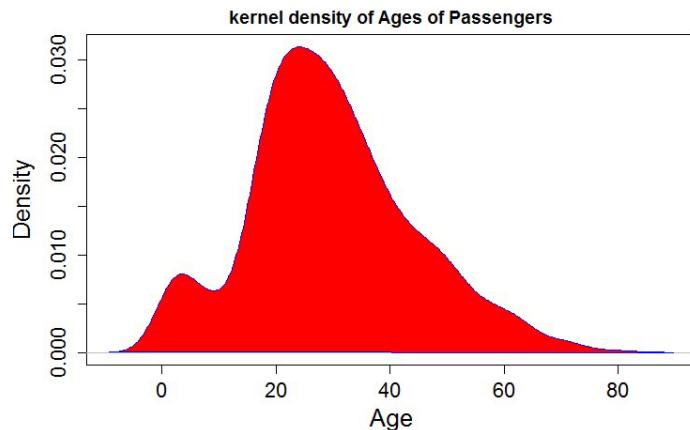


Figure 1.6

Comparison of density plots of Age with different Gender (Figure: 1.7)

### Kernel Density Plot of Ages By Gender

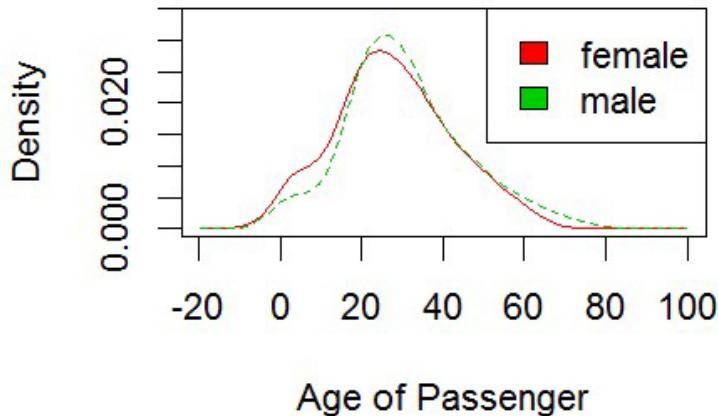


Figure 1.7

Did age have an impact on survival? (Figure: 1.8)

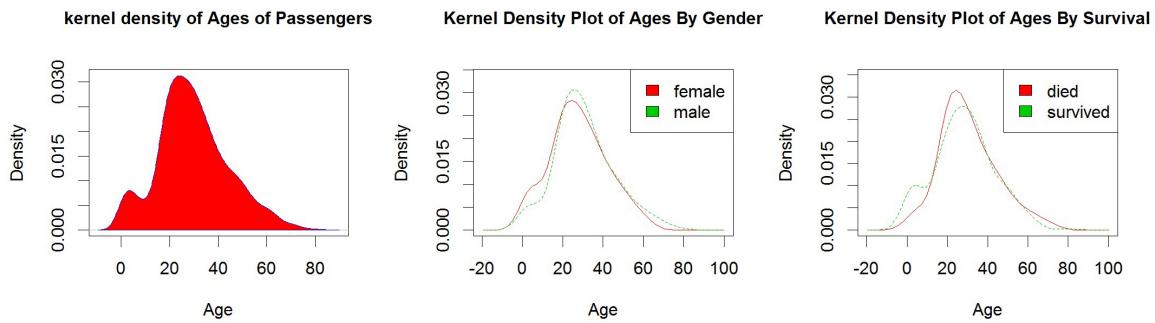


Figure 1.8

## 1.6 Feature Engineering

Let's add more features to our data. In this specific exercise we will derive from age, whether or not a person is a child or an adult. For the purpose of this example we will define a child as someone who is below 13 years of age.

```
1 ## Multi dimensional comparison
2 Child <- titanic$Age # Isolating age.
3 ## Now we need to create categories: NA = Unknown, 1 =
   Child, 2 = Adult
4 ## Every age below 13 (exclusive) is classified into age
   group 1
5 Child[Child<13] <- 1
6 ## Every child 13 or above is classified into age group 2
7 Child[Child>=13] <- 2
8 ## Use labels instead of 0's and 1's
9 Child[Child==1] <- "Child"
10 Child[Child==2] <- "Adult"
```

Merging the new age column back into the dataset.

```
1 # Appends the new column to the titanic dataset
2 titanic_with_child_column <- cbind(titanic, Child)
3 # Removes rows where age is NA
4 titanic_with_child_column <- titanic_with_child_column[!is
   .na(titanic_with_child_column$Child),]
5 ## Show part of the data frame with new feature: Child
6 head(titanic_with_child_column[,c(1:3, 13)], 4)
7 ##   PassengerId Survived Pclass Child
8 ## 1           1     died      3  Adult
9 ## 2           2 survived     1  Adult
10 ## 3          3 survived     3  Adult
11 ## 4          4 survived     1  Adult
```

Now let's look at the survival rate of our newly created column. Did children have higher survival rates over adults? (Figure: 1.9)

```

1 par(mfrow = c(1, 2), mar = c(0, 1, 2, 1), oma = c(0, 1, 0,
   1), cex=1.5)
2 pie(table(titanic_with_child_column[titanic_with_child_
  column$Child=="Child", "Survived"]),
  labels=c("Dead","Survived"), main="Survival Portion
  of Child")
4 pie(table(titanic_with_child_column[titanic_with_child_
  column$Child=="Adult", "Survived"]),
  labels=c("Dead","Survived"), main="Survival Portion
  of Adult")

```

## Survival Portion of Child   Survival Portion of Adult

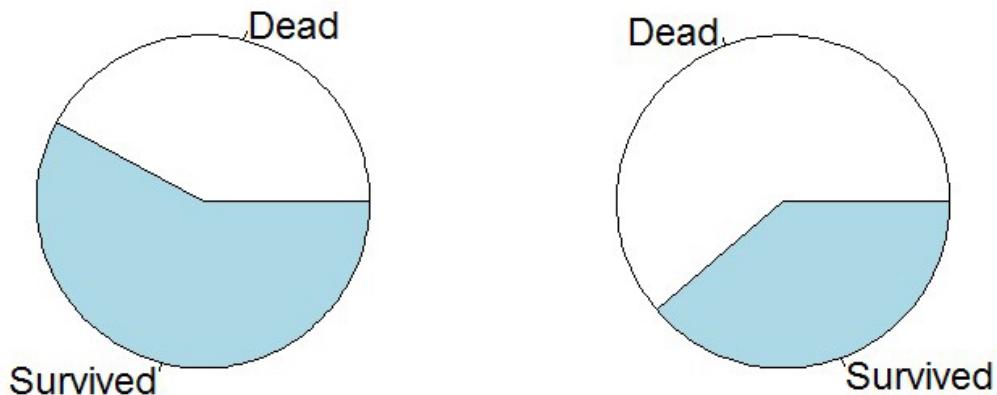


Figure 1.9

Child vs Gender (Figure: 1.10)

```

1 child.gender <- table(titanic_with_child_column$Gender,
2   titanic_with_child_column$Child)
3 gender.survived <- table(titanic_with_child_column$Survived,
4   titanic_with_child_column$Gender)
5 child.survived <- table(titanic_with_child_column$Survived
6   , titanic_with_child_column$Child)
7 barplot(child.gender, main="Gender Proportion by Child or
8   Adult",
9   col=c("red","blue"), legend = rownames(child.
10   gender))
11 barplot(gender.survived, main="Survived Proportion by
12   Gender",
13   col=c("red","blue"), legend = rownames(gender.
14   survived))
15 barplot(child.survived, main="Survived Proportion by Child
16   or Adult",
17   col=c("red","blue"), legend = rownames(child.
18   survived))

```

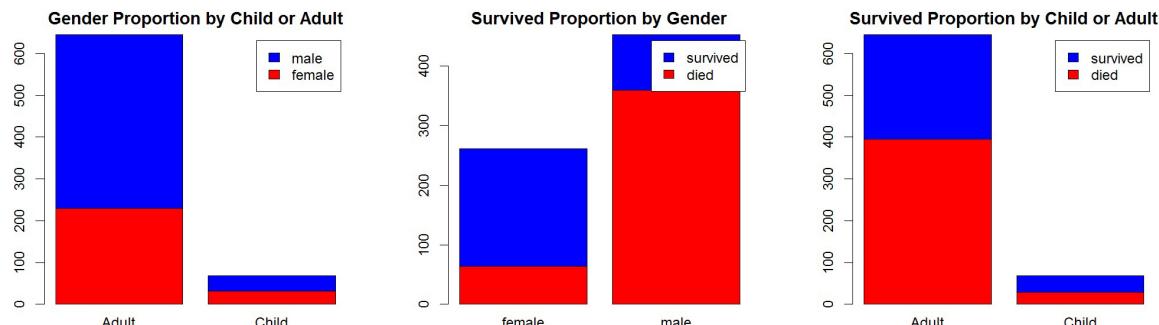


Figure 1.10

## 1.7 Advanced Visualizations in R

- **Base graphics:** constructed piecemeal. Conceptually simpler and allows plotting to mirror the thought process.
- **Lattice graphics:** entire plots created in a simple function call.
- **ggplot2 graphics:** an implementation of the Grammar of Graphics by Leland Wilkinson. Combines concepts from both base and lattice graphics. (Need to install ggplot2 library)

A list of interactive visualization in R can be found at:

<http://ouzor.github.io/blog/2014/11/21/interactive-visualizations.html>

Base plotting system with airquality data. Note that air quality data ships natively with R. (Figure: 1.11)

```
1 par(cex=1.5) ## increase the size of texts  
2 ## scatter plot  
3 plot(x = airquality$Temp, y = airquality$Ozone)
```

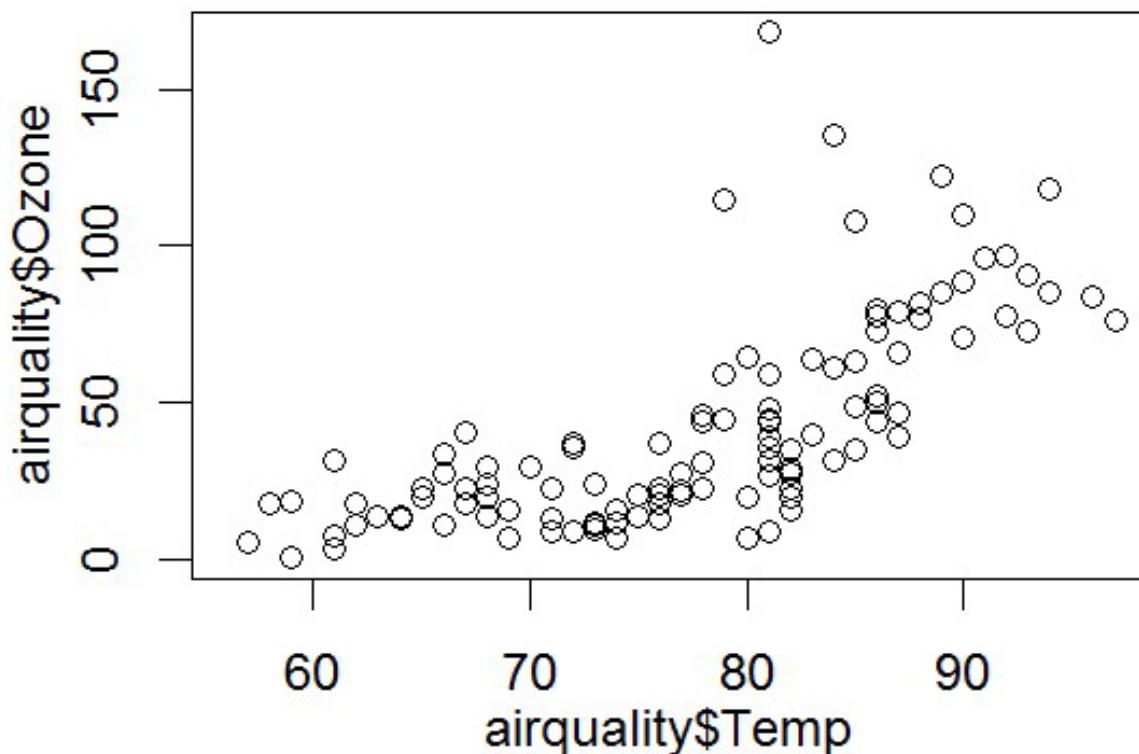


Figure 1.11

Base plotting system. (Figure: 1.12)

```

1 # par() function is used to specify global graphics
2 # parameters that affect all plots in an R session.
3 # Type ?par to see all parameters
4 par(mfrow=c(1, 2), mar=c(4, 4, 2, 1), oma=c(0, 0, 2, 0),
5   cex=1.5)
6 # mfrow=c(1, 2): the figures will be drawn in an 1x2
7 # array by row. mar/oma: A numerical vector of the
8 # form (bottom, left, top, right), which gives the
9 # number of lines of margin to be specified on the
10 # four sides of the plot for mar, and size of the outer
11 # margins in lines of text for oma.
12 with(airquality,
13   {
14     plot(Wind, Ozone, main="Ozone and Wind")
15     plot(Temp, Ozone, main="Ozone and Temperature")
16     mtext(
17       "Ozone and Weather in New York City",
18       outer=TRUE
19     )
20   }
21 )

```

For the difference and relation of mar and oma, check out:

<http://research.stowers-institute.org/mcm/efg/R/Graphics/Basics/mar-oma/index.htm>

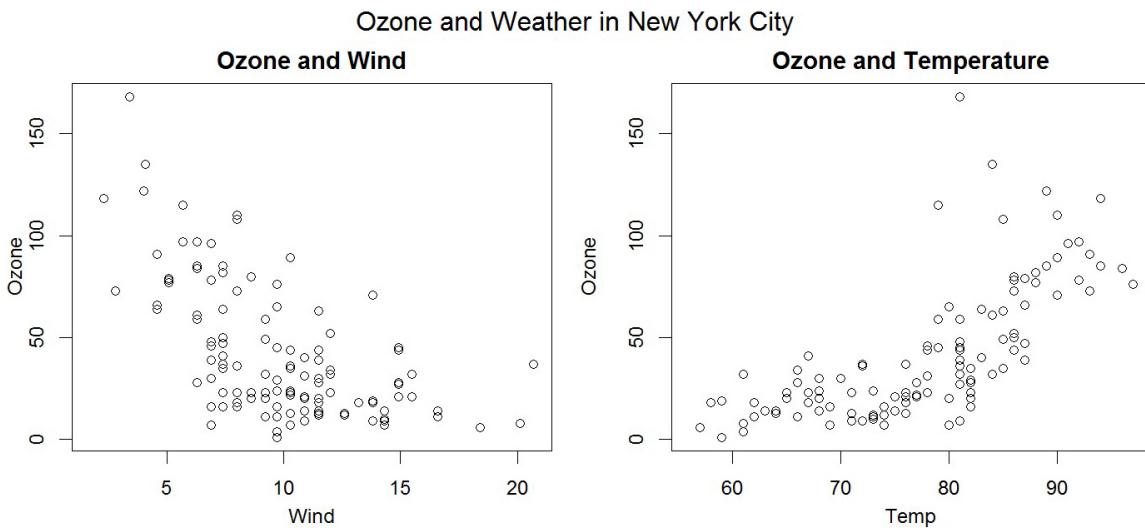


Figure 1.12

Plotting functions (high level) **PHASE ONE: Mount a canvas panel on the easel, and draw the draft (Initialize a plot).**

- **boxplot()**: a boxplot show the distribution of a vector. It is very useful to example the distribution of different variables.
- **plot()**: one of the most frequently used plotting functions in R.
- **barplot()**: create a bar plot with vertical or horizontal bars.
- **hist()**: compute a histogram of the given data values.
- **pie()**: draw a pie chart.

Remember to use **?plot** or **str(plot)**, etc. to check the arguments when you want to make more personalized plots. A **tutorial** of base plotting system with more details: <http://bcb.dfci.harvard.edu/~aedin/courses/BiocDec2011/2.Plotting.pdf>

Plotting functions (low level) **PHASE TWO: Add more details on your canvas, and make an artwork (Add more on an existing plot)**

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

Save your artwork R can generate graphics (of varying levels of quality) on almost any type of display or printing device. Like:

- **postscript()**: for printing on PostScript printers, or creating PostScript graphics files.
- **pdf()**: produces a PDF file, which can also be included into PDF files.
- **jpeg()**: produces a bitmap JPEG file, best used for image plots. **helpDevices** for a list of them all. Simple example:

```
1 ## png(filename = 'plot1.png', width = 480, height = 480,
2 # units = 'px')
3 ## plot(x, y)
4 ## dev.off()
```

Example: boxplot and histogram

```
1 ## the layout
2 par(mfrow = c(2, 1), mar = c(2, 0, 2, 0), oma = c(0, 0, 0,
3 # 0))
4 ## histogram at the top
5 hist(airquality$Ozone, breaks=12, main = "Histogram of
6 # Ozone")
7 ## box plot below for comparison
8 boxplot(airquality$Ozone, horizontal=TRUE, main = "Box
9 # plot of Ozone")
```

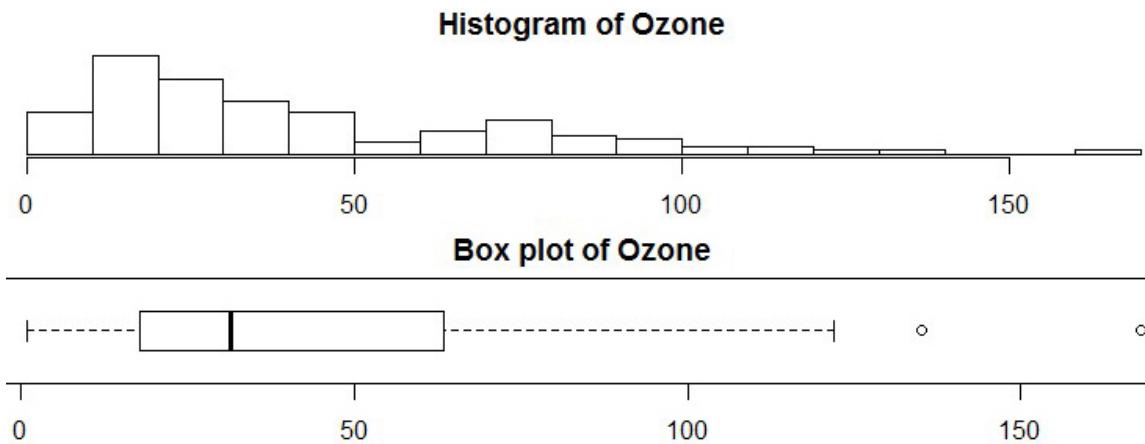


Figure 1.13

### Lattice plotting system

```

1 ## install.packages("lattice") # install the lattice
  library if you haven't
2 library(lattice) # load the library
3 ## Warning: package 'lattice' was built under R version
  3.1.2
4 ## set the seed so our plots are the same
5 set.seed(10)
6 ## use rnorm(): generate 100 random numbers
7 x <- rnorm(100)
8 ## use rep() function to generate a vector, with first 25
  elements are 1,
9 ## second 25 elements are 2, ...
10 f <- rep(1:4, each = 25)
11 y <- x + f - f * x + rnorm(100, sd = 0.5)
12 # name the levels of the factor: first 25 elements are in
  Group 1,
13 ## second 25 elements are in Group 2, ...
14 f <- factor(f, labels = c("Group 1", "Group 2", "Group 3",
  "Group 4"))
15 xyplot(y ~ x | f)

```

(Figure: 1.14)

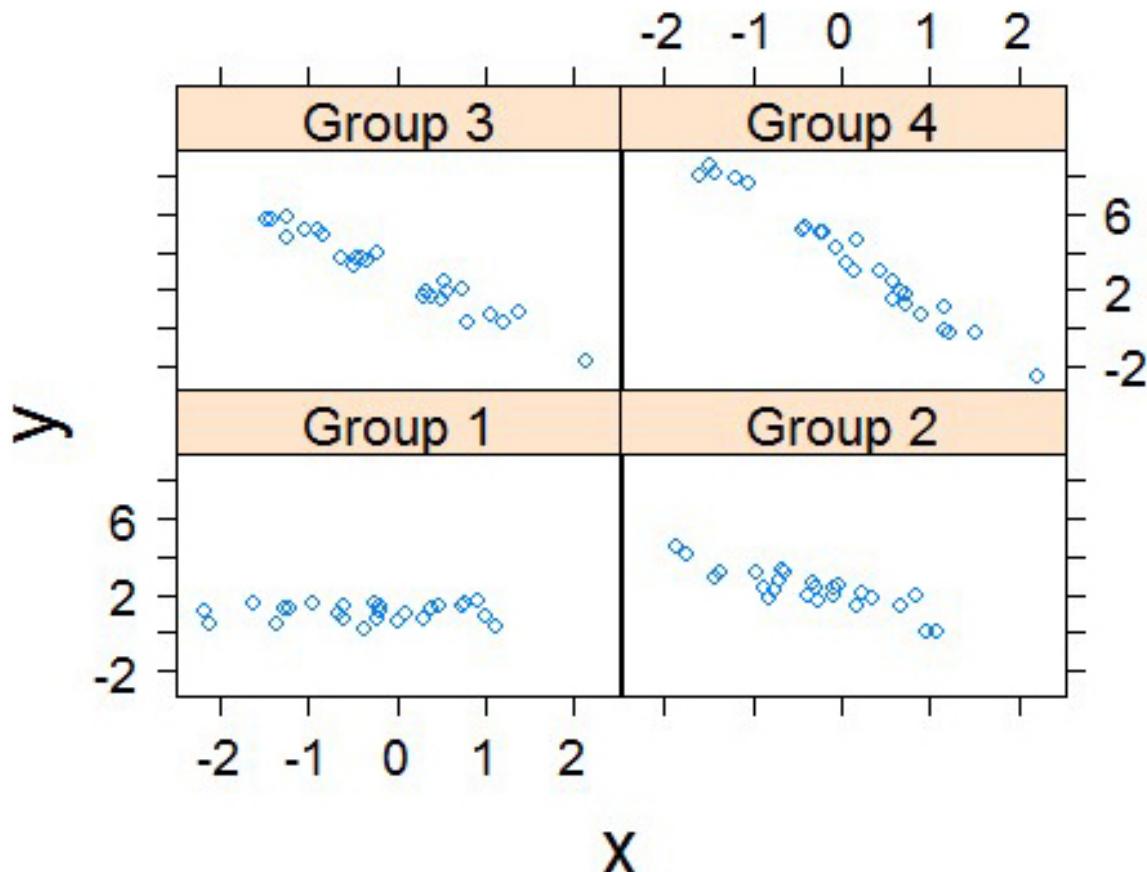


Figure 1.14

Want more on the plot? Customize the panel function:

```

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: 1.15)

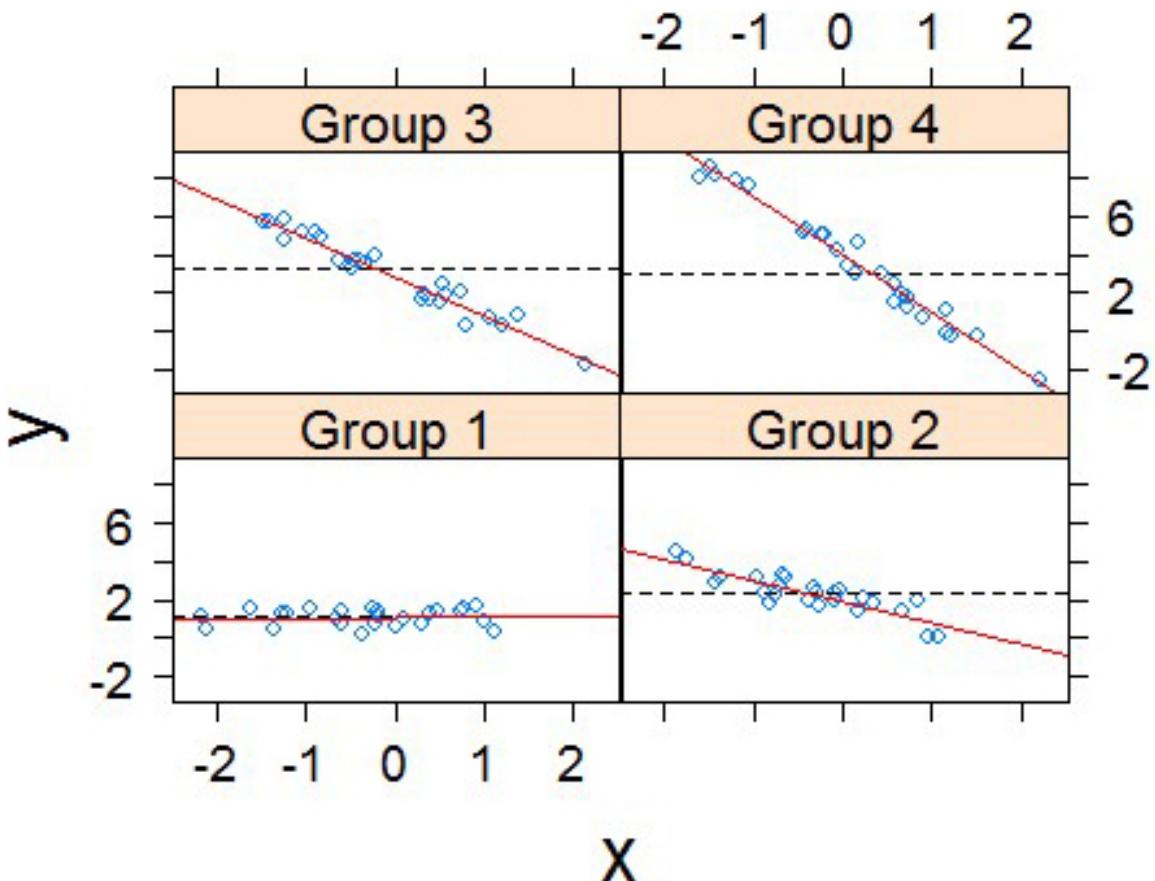


Figure 1.15

### Plotting functions

- **xypplot()**: main function for creating scatterplots
- **bwplot()**: box and whiskers plots (box plots)
- **histogram()**: histograms
- **stripplot()**: box plot with actual points
- **dotplot()**: plot dots on “violin strings”
- **splom()**: scatterplot matrix (like pairs() in base plotting system)
- **levelplot()/contourplot()**: plotting image data

Very useful when we want a lot...

```
1 pairs(iris) ## iris is a data set in R
```

(Figure: 1.16)

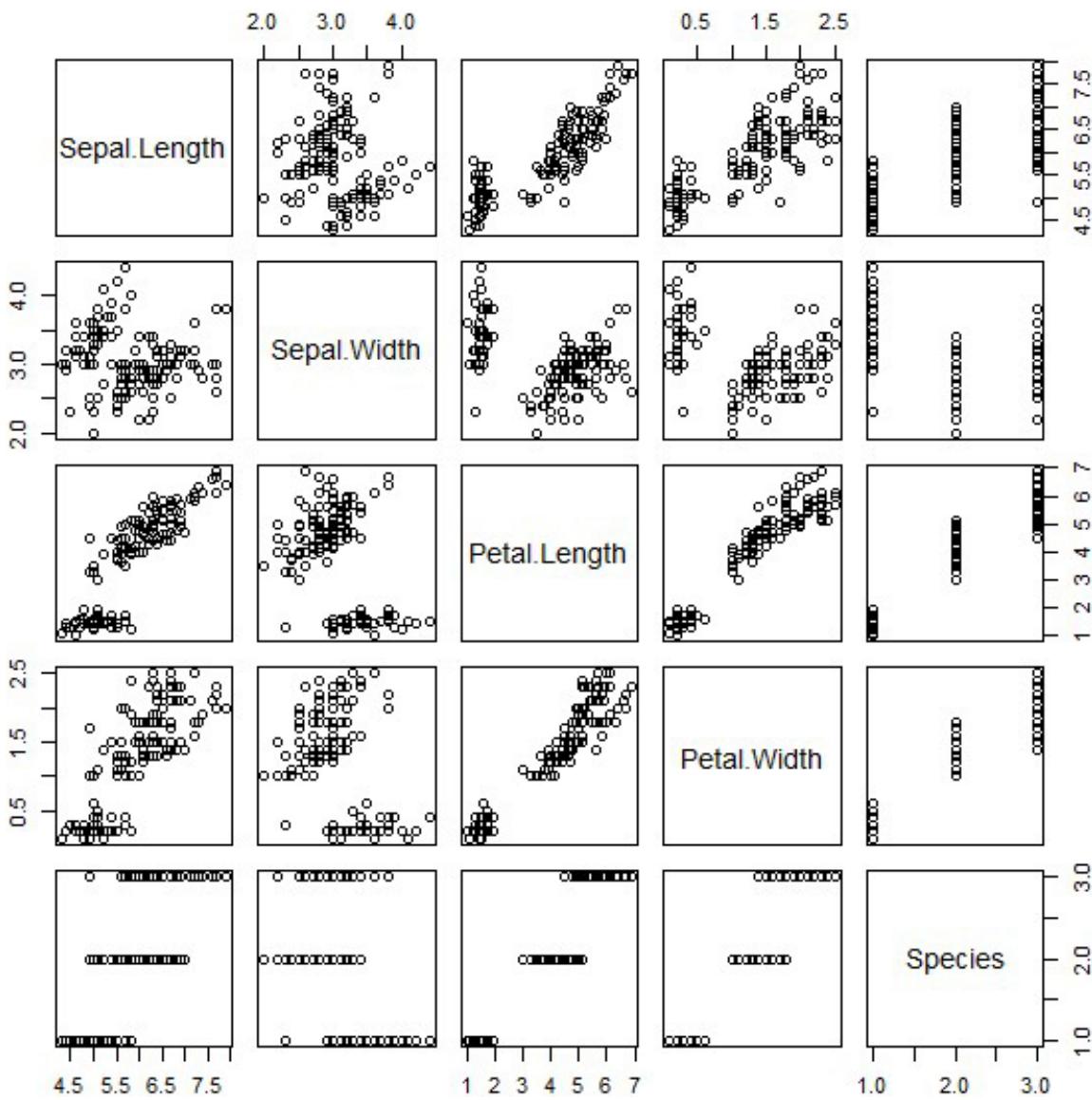


Figure 1.16

## ggplot2

- An implementation of the **Grammar of Graphics** by Leland Wilkinson
- Written by Hadley Wickham (while he was a graduate student at Iowa State)
- A “third” graphics system for R (along with base and lattice) Available from CRAN via `install.packages()` web site: <http://ggplot2.org>(better documentation)
- Grammar of graphics represents the abstraction of graphics ideas/objects Think “verb”, “noun”, “adjective” for graphics “Shorten” the distance from mind to page
- Two main functions: `qplot()` hides what goes on underneath, which is okay for most operations `ggplot()` is the core function and very flexible for doing this `qplot()` cannot do

qplot function The `qplot()` function is the analog to `plot()` but with many built-in features Syntax somewhere in between base/lattice Difficult to be customized (don’t bother, use full ggplot2 power)

in that case)

```

1 ## install.packages("ggplot2") # need to install ggplot2
  first
2 library(ggplot2) # load the library
3 levels(mpg$drv) <- c("4wd", "front-wheel drv", "rear-wheel
  drv") # rename the levels
4 theme_set(theme_gray(base_size = 20)) # change text size
5 qplot(displ, hwy, data = mpg, facets = .~drv)

```

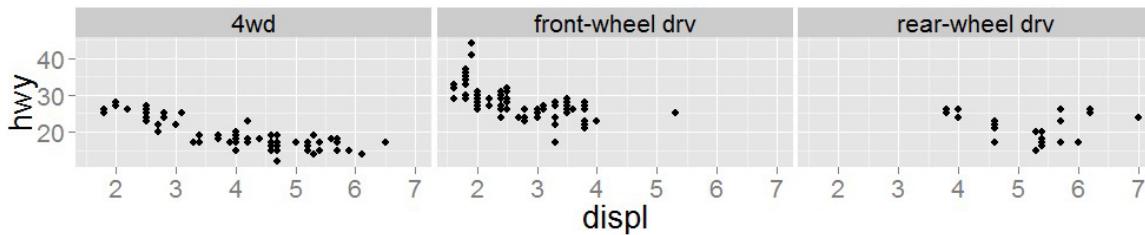


Figure 1.17

ggplot function When building plots in ggplot2 (ggplot, rather than using qplot) The “**artist’s palette**” model may be the closest analogy Plots are built up in layers \* **Step I: Input the data noun:** the data

```

1 library(ggplot2) ## need to install and load this library
2 g <- ggplot(iris, aes(Sepal.Length, Sepal.Width)) ## this
  would not show you add plot

```

**Step II: Add layers adjective:** describe the type of plot you will produce.

```

1 g <- g + geom_point() + geom_smooth(method = "lm") + facet
  _grid(. ~ Species) +
  theme(text = element_text(size = 20))
2 g
3
4 #geom_smooth(method='lm'): add linear regression line;
5 #geom\_\_grid(. \~ Species): lay out panels in a grid by
  Species;
6 #theme(): revision of the appearance (It is actually the
  step III)

```

(Figure: 1.18)

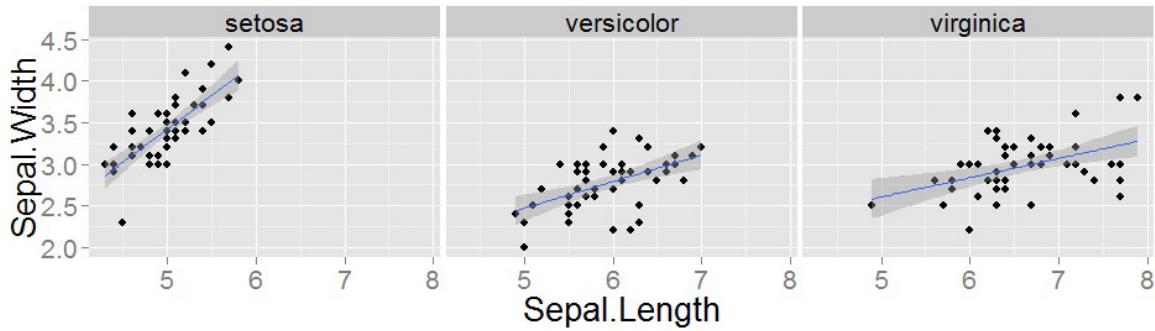


Figure 1.18

**Step III: Add metadata and annotation adjective:** control the mapping between data and aesthetics.

```

1 g + ggtitle("Sepal length vs. width for different species"
  ) + theme_bw() +
  theme(text = element_text(size = 18)) # add "verb"
2
3 #ggtitle(): add the title
4 #theme_bw: choose the classic dark-on-light ggplot2 theme.
5

```

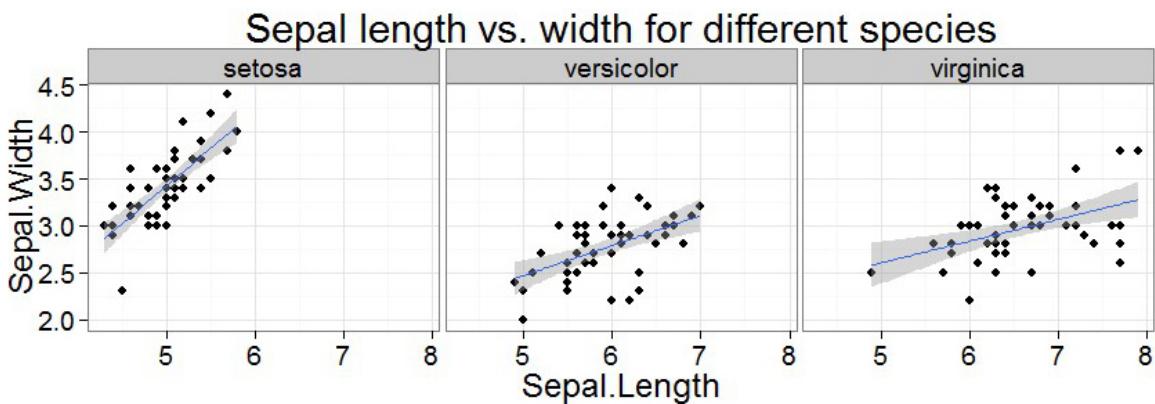


Figure 1.19

Great documentation of ggplot with all functions in step II and III and demos: <http://docs.ggplot2.org/current/>

## Overview of the diamond data

```
1 data(diamonds) # loading diamonds data set
2 head(diamonds, 9) # first few rows of diamond data set
3 #   carat      cut color clarity depth table price     x
4 #       y      z
5 # 1 0.23    Ideal     E     SI2  61.5    55    326 3.95
6 # 2 0.21  Premium     E     SI1  59.8    61    326 3.89
7 # 3 0.23      Good     E     VS1  56.9    65    327 4.05
8 # 4 0.29  Premium     I     VS2  62.4    58    334 4.20
9 # 5 0.31      Good     J     SI2  63.3    58    335 4.34
10 # 6 0.24 Very Good     J    VVS2  62.8    57    336 3.94
11 # 7 0.24 Very Good     I    VVS1  62.3    57    336 3.95
12 # 8 0.26 Very Good     H     SI1  61.9    55    337 4.07
13 # 9 0.22      Fair     E     VS2  65.1    61    337 3.87
```

Histogram of carat

```
1 library(ggplot2)
2 ggplot(data=diamonds) + geom_histogram(aes(x=carat),
   binwidth=1/5) +
3   theme(text = element_text(size = 25))
```

(Figure: 1.20)

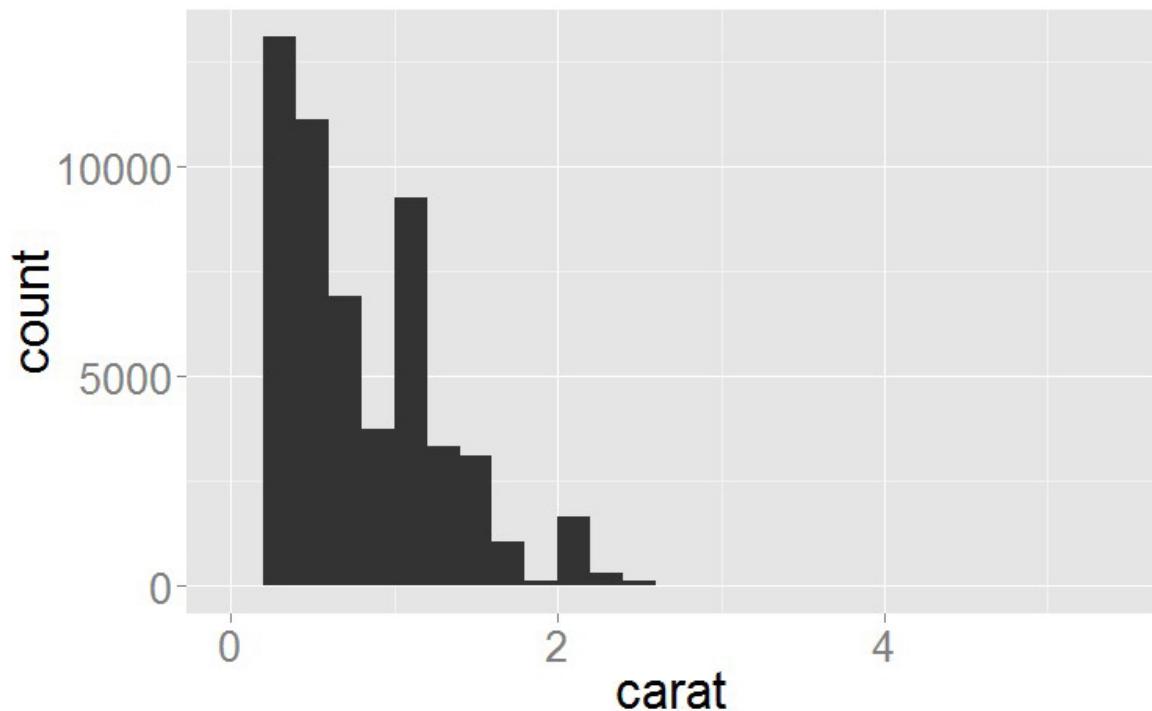


Figure 1.20

Density plot of carat

```
1 ggplot(data=diamonds) +
2   geom_density(aes(x=carat), fill="gray50") +
3   theme(text = element_text(size = 25))
```

(Figure: 1.21)

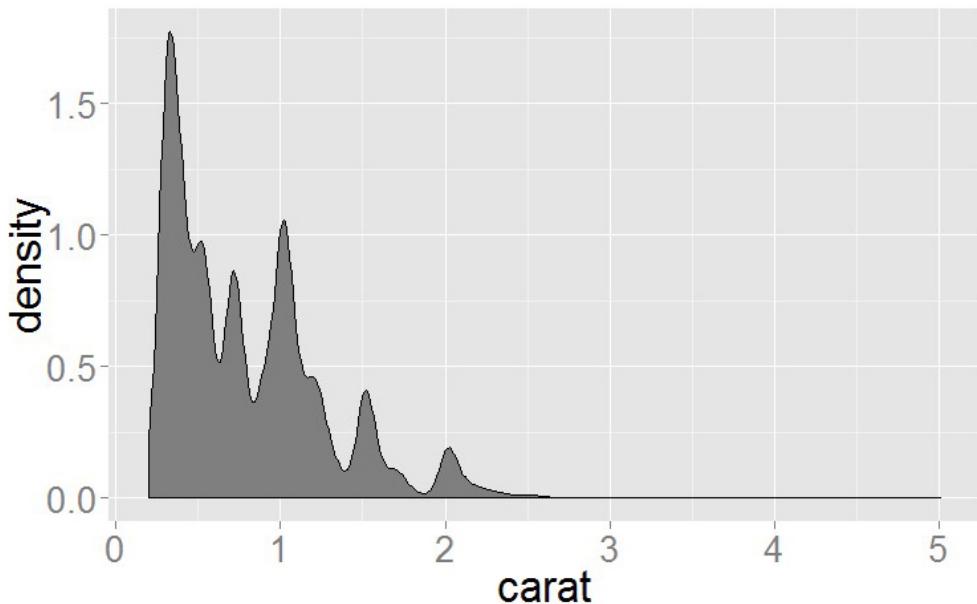


Figure 1.21

Scatter plots (carat vs. price)

```
1 ggplot(diamonds, aes(x=carat, y=price)) + geom_point() +
2   theme(text = element_text(size = 25))
```

(Figure: 1.22)

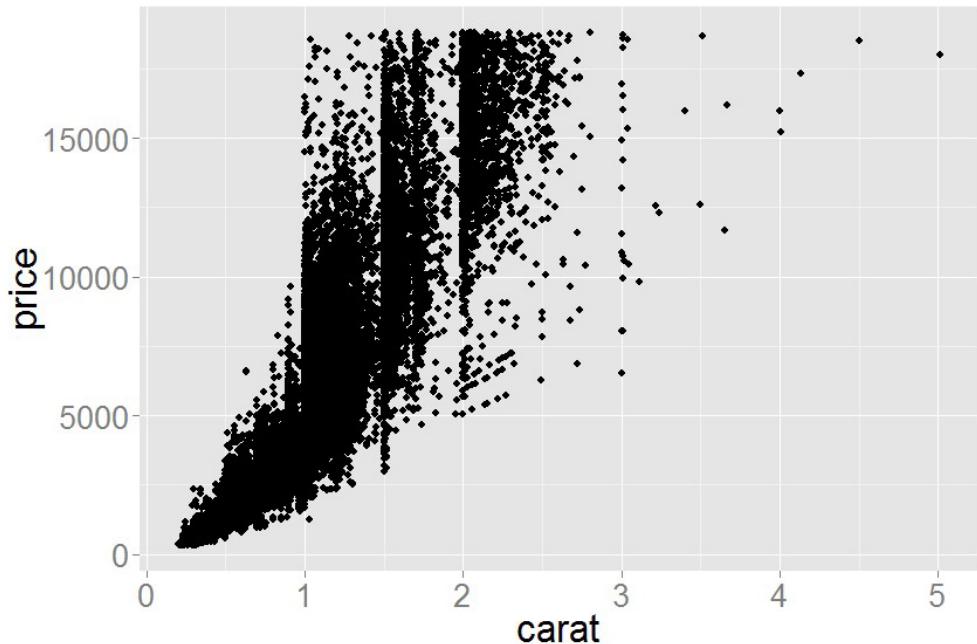


Figure 1.22

Carat with colors

```
1 g = ggplot(diamonds, aes(x=carat, y=price)) # saving first
   layer as variable
2 # rendering first layer and adding another layer
3 g + geom_point(aes(color=color)) + theme(text = element_
  text(size = 25))
```

(Figure: 1.23)

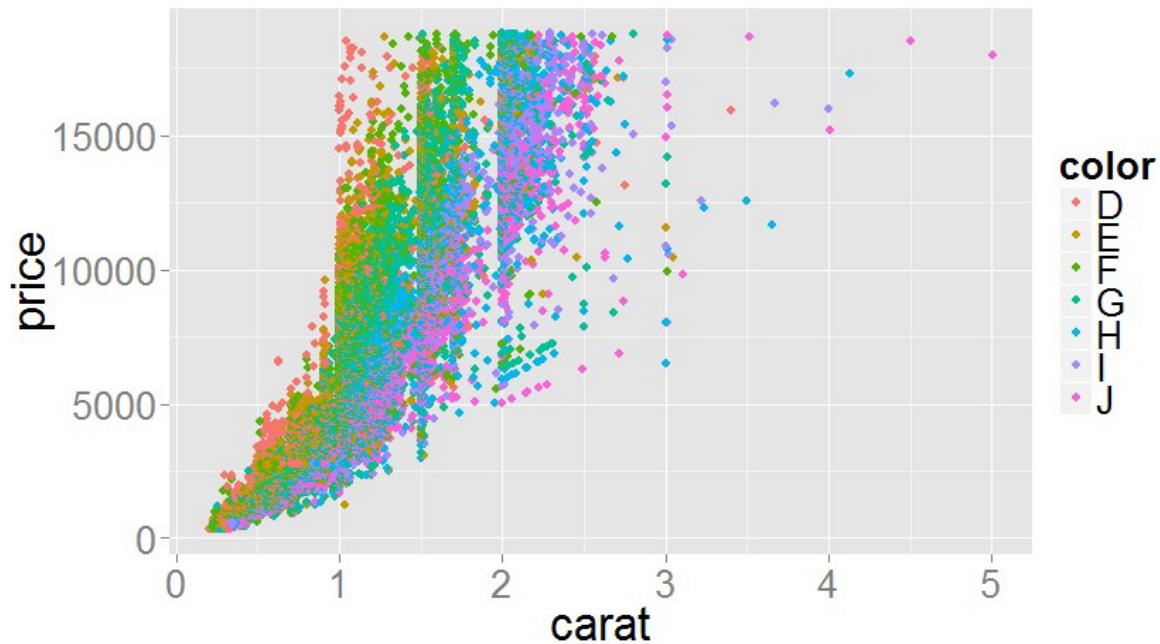


Figure 1.23

Carat with colors (more details)

```
1 g + geom_point(aes(color=color)) + facet_wrap(~color) +
2   theme(text = element_text(size = 20))
```

(Figure: 1.24)

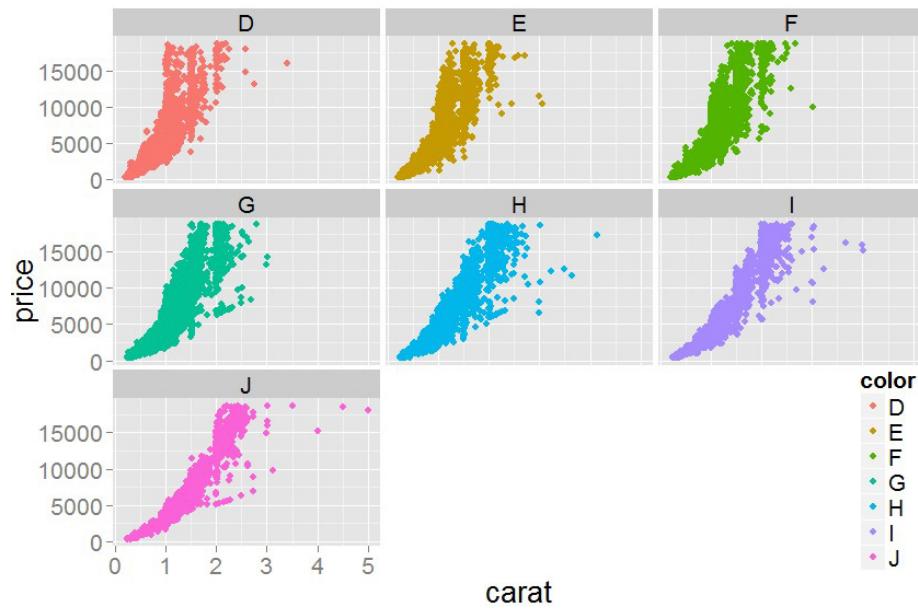


Figure 1.24

Let's consider cut and clarity (Figure: 1.25)

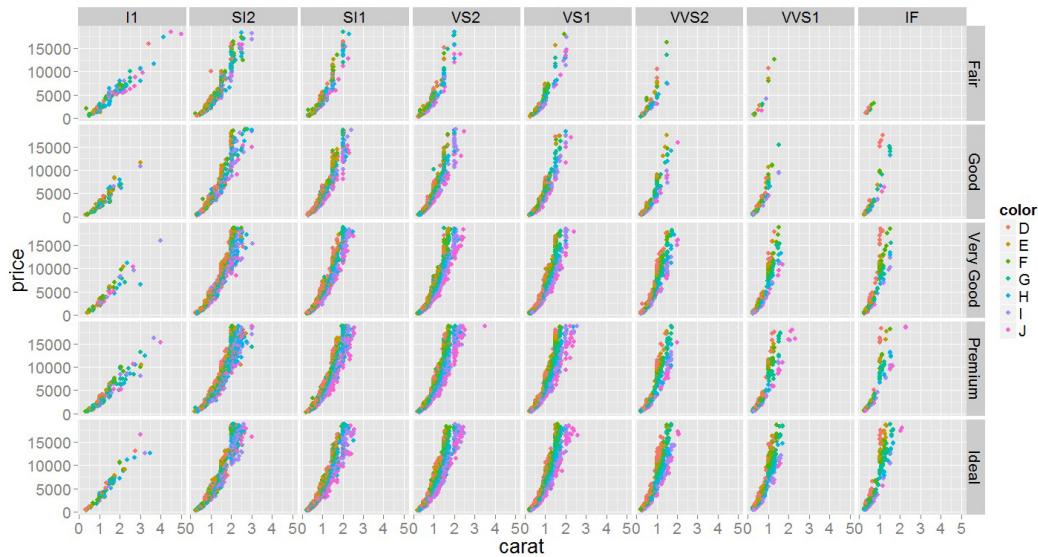


Figure 1.25

Your turn! What is your knowledge of diamond's price after exploring this data?

Fare matters

```
1 library(ggplot2)
2 ggplot(titanic_with_child_column, aes(y=Fare, x=Survived))
+   +
3 geom_boxplot() + facet_grid(Gender ~ Child) +
4 theme(text = element_text(size = 25))
```

(Figure: 1.26)

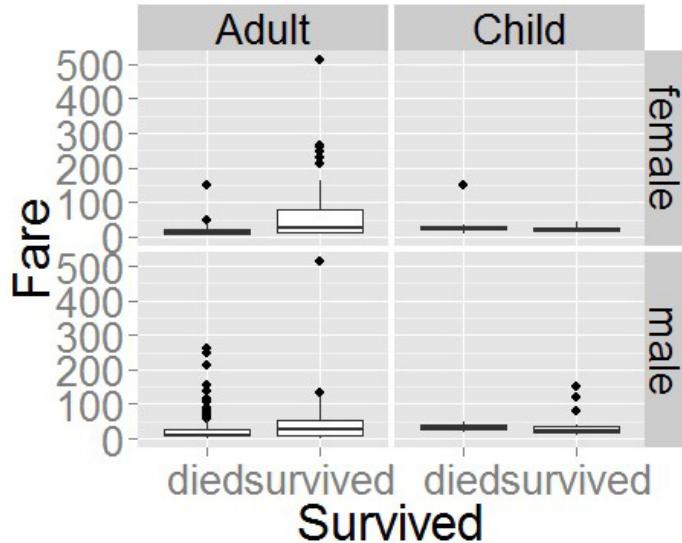


Figure 1.26

How about fare, ship class, port embarkation? (Figure: 1.27)

```
1 library(ggplot2)
2 titanic$Pclass = as.factor(titanic$Pclass)
3 ggplot(titanic, aes(y=Fare, x=Pclass)) + geom_boxplot() +
4   facet_grid(~ Embarked) + theme(text = element_text(size = 25))
```

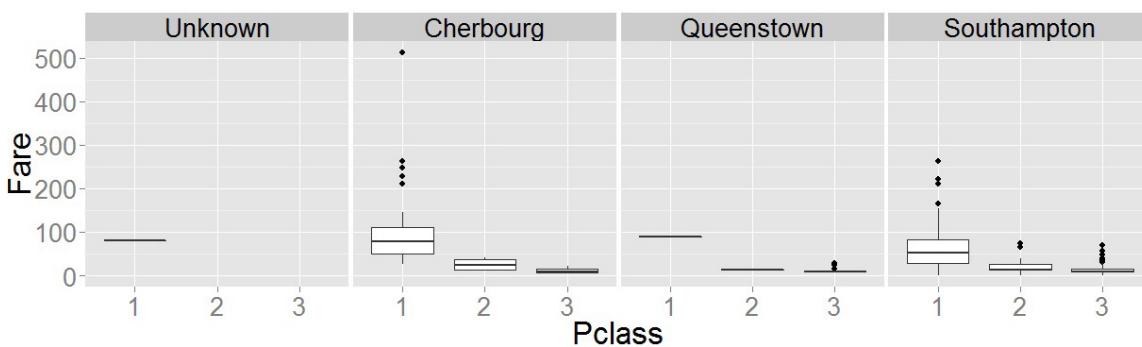


Figure 1.27

# COVERAGE INCLUDES

**Data Science Dojo's Handbook** is a collection of exercises that are designed to teach you the basics of data science and data engineering through hands-on experience. The book is accessibly written to guide anyone interested in learning more about data science.

Getting started should be easy. This Handbook will serve as an introduction to the fundamental concepts of data science and engineering, simplifying even the most complex topics to ground your knowledge in a strong foundation. Rather than focusing on the inaccessible math and programming of data science, that often deters most beginners, this book will give users the crash course into data mining theory.

Learn the best practices of the industry with the most up-to-date data science tools. Readers will be exposed to a variety of technologies such as Azure Machine Learning Studio, Cloud-based SQL databases, Hive, Hadoop, stream processors, and data ingestion hubs.

The hands-on approach in this book will cover both data science and data engineering; intertwining the two concepts to encourage future data scientists to cultivate a set of complementary skills.

## ABOUT US

**Data Science Dojo** is a Seattle company that offers bootcamps, corporate trainings, consulting services, and community events in all things data science. The company started in 2014, stemming from the realization that data science and data engineering were being treated as separate entities. Data Science Dojo instead treats the concepts holistically by encouraging a hands-on, accessible approach to learning. Data Science Dojo's carefully crafted curriculum and community of data scientists have enabled individuals from all around the world to unleash the data scientists within.

- Azure Machine Learning Studio
- Importing & writing datasets
- Data exploration & visualization
- Preprocessing data & missing values cleansing
- Data partitioning & cross validation
- Classification algorithms
- Regression algorithms
- Train, score, & model evaluation
- Deploying a predictive model as a web service
- Text entity extraction
- SQL databasing in the cloud (Azure SQL Database)
- Hadoop & HiveQL (HDInsight)
- Setting up big data pipelines in the cloud
- Analytics on streaming data
- Data ingestion (Azure Event Hubs)
- Stream processing (Azure Stream Analytics & StreamQL)
- Internet of Things