

```

1  setwd("C:/Users/Min Seong Kim/Dropbox/R_programming/lecture")
2
3  data.set <- read.csv("county_data.csv", stringsAsFactors = FALSE)
4
5  head(data.set)
6  tail(data.set)
7  str(data.set)
8
9  data.1 <- data.set[1:111,]
10 tail(data.1)
11
12 data.state <- factor(data.1$State, levels=c("Alabama", "Alaska", "Arizona"))
13 summary(data.state) # or table(data.state)
14
15 plot(data.state) # If plot() function is used for a single factor, it
will count frequencies (bar chart).
16 title(main="Number of Counties")
17
18 plot(data.state, data.1$IncomePerCap, range=1, las=1) # If plot() is used for
x=factor, y=numeric, it will produce boxplots.
19 boxplot(data.1$IncomePerCap~data.state, range=1, las = 2) # How the distribution of y
changes conditional on x
20
# whiskers extend to the
most extreme data point
which is no more
# than range times the
interquartile range
21
22
23 big <- NA
24 big[data.1$TotalPop>70000] <- "big"
25 big[data.1$TotalPop<=70000] <- "small"
26
27 par(mar=c(6,4,2,1)) # bottom, left, up, right
28 boxplot(data.1$IncomePerCap~big*data.state, col=c("blue", "red"), main="Income per
Capital", range=1, las = 2)
29
30 plot(data.1[,3:ncol(data.1)])
31 pairs(data.1[,3:ncol(data.1)]) # scatterplot matrix
32
33 # highlight outliers
34 data.2 <- data.set[data.set$State=="California",]
35 plot(data.2$TotalPop)
36 abline(h=10000000, lty=2)
37 index.1 <- which(data.2$TotalPop > 10000000)
38 points(index.1,data.2$TotalPop[index.1], pch=16, col="red" )
39
40 abline(h=2000000, lty=2)
41 index.2 <- which(data.2$TotalPop > 2000000 & data.2$TotalPop < 10000000 )
42 points(index.2,data.2$TotalPop[index.2], pch=16, col="blue" )
43
44 text(c(index.1,index.2),data.2$TotalPop[c(index.1,index.2)], labels =
data.2$County[c(index.1,index.2)], cex=0.8)
45 # labels: character vector specifying the text to be written
46
47 # Adding a straight line on a plot
48 proportion <- data.1$Women/data.1$TotalPop
49 plot(proportion, data.1$Unemployment, pch=20, xlab="Women/Total Population",
ylab="Unemployment rate")
50 linear.fit <- lm(data.1$Unemployment~proportion)
51 abline(linear.fit, lwd=3, col="red")
52
53 linear.fit
54 # You can also specify coefficient directly to use abline()
55 abline(a=-3, b=30, lwd=3, lty=4, col="blue")
56
57 abline(v=0.5, col="red", lwd=3, lty=2)
58 abline(h=15, col="blue", lwd=3, lty=1)
59
60 # hist(), image(): histogram and heatmap

```

```

61 # heat.colors(), topo.colors(), etc: create a color vector
62 # density(): estimate density, which can be plotted
63
64 # Plotting a histogram
65 # To plot a histogram of a numeric vector, use hist()
66
67 ave.income <- data.set$IncomePerCap
68 hist(ave.income)
69
70 # Histogram options
71 # Several options are available as arguments to hist(), such as col, freq,
72 # breaks, xlab, ylab, main
73
74 hist(ave.income, col="pink", freq=TRUE) # Frequency scale, default
75
76 hist(ave.income, col="pink", freq=FALSE, # Probability scale, and more options
77       breaks=seq(0,70000,by=10000), xlab="Income per capita", main="County Level Average
78       Income per capital")
79
80 hist(ave.income, col="pink", freq=FALSE, # Probability scale, and more options
81       breaks=seq(0,70000,by=50), xlab="Income per capita", main="County Level Average
82       Income per capital")
83
84 hist(ave.income, col="pink", freq=FALSE, # Probability scale, and more options
85       breaks=seq(0,70000,by=2000), xlab="Income per capita", main="County Level Average
86       Income per capital")
87
88 # Adding a histogram to an existing plot
89
90 # To add a histogram to an existing plot (say, another histogram), use hist()
91 # with add=TRUE
92
93 hist(ave.income + 20000, col=rgb(0.1,0.1,0.5,0.2), # Note: Using a transparent color:
94       red, green, blue, alpha(degree of transparency)
95       freq=FALSE, breaks=seq(0,90000,by=2000), add=TRUE)
96
97 # Adding a density curve to a histogram
98
99 # To estimate a density from a numeric vector, use density().
100 # This returns a list; it has components x and y, so we can actually
101 # call lines() directly on the returned object
102
103 hist(ave.income, col="pink", freq=FALSE, # Probability scale, and more options
104       breaks=seq(0,70000,by=2000), xlab="Income per capita", main="County Level Average
105       Income per capital")
106
107 density.est = density(ave.income, adjust=1.5) # 1.5 times the default bandwidth, try
108 different values
109
110 lines(density.est, lwd=3)
111
112 # Exercise
113
114 data.set <- read.csv("histogram.csv", stringsAsFactors = FALSE )
115 head(data.set)
116 data.set <- data.set[data.set$YRDATA==2013,]
117
118 # Using plot() and points() functions, highlight the 4 best schools in CT in terms of
119 expenditure per student (TOTALEXP/ENROLL).
120
121 data1 <- data.set[data.set$STATE == "Connecticut",]
122 data1 <- data1[data1$ENROLL != 0,]
123 exp_per_stu <- data1$TOTALEXP/data1$ENROLL
124
125 plot(exp_per_stu)
126
127 index1 <- which(exp_per_stu > 30)
128 points(index1,exp_per_stu[index1], pch=16, col="red" )
129

```

```
123 text(index1,exp_per_stu[index1]-1, labels = data1[index1, 3], cex=0.8)
124
125 # Make two histograms together for expenditure per student (TOTALEXP/ENROLL) of
    Connecticut and Alabama
126
127 head(data.set)
128 data.set <- data.set[data.set$ENROLL != 0,]
129 data.hist = data.set$TOTALEXP/data.set$ENROLL
130 # data.hist <- as.numeric(exp_per_stu)
131 hist(data.hist[data.set$STATE == "Connecticut"], col="pink", freq=FALSE,
132       breaks=seq(0,40,by=2), xlab="Total Expenditure", main="Education expenditure in CT
    and MI")
133 hist(data.hist[data.set$STATE == "Alabama"], col=rgb(0.3,0.5,0.5,0.5), # Note: Using a
    transparent color
134       freq=FALSE, breaks=seq(0,40,by=2), add=TRUE)
135
136
```