

HW1 - Intro to Machine Learning

Saurabh Madaan

09/15/2012

1 Part 1 - HW01pb1data.csv

1.1 Attributes

The dataset has 800 observations of 5 variables. The first 3 columns are *integers* whereas the 4th and 5th columns are *factors*.

In R, this can be seen using the following commands:

```
1 setwd("/Users/Saurabh/Documents/ML_UCSC/week1/HW")
2 data<-read.csv("HW01pb1data.csv",header=FALSE)
3 class(data)
4 str(data)
5 #800 obs. of 5 variables. V1-3 are int, V4-5 are Factors
```

1.2 Reason for Categorical Variables

Columns 4 and 5 mostly have integers in them, but when one looks in to the *levels*, it is quickly visible that they also have strings “thirty five” and “twenty five”, respectively. Consequently, *R* treats them as factors.

Below are the R commands which provide more details:

```
9 l4<-levels(data[,4])
10 #has integers, and "thirty five"
11 which(data[,4]=="thirty_five")
12 #[1] 405
13
14 l5<-levels(data[,5])
15 #has integers, and "twenty five"
16 which(data[,5]=="twenty_five")
```

```
17 # [1] 531
```

1.3 Plots for numeric and categorical variables

1. Plot of the 1st column (numerical data)

```
41 plot (data[, 1])
```

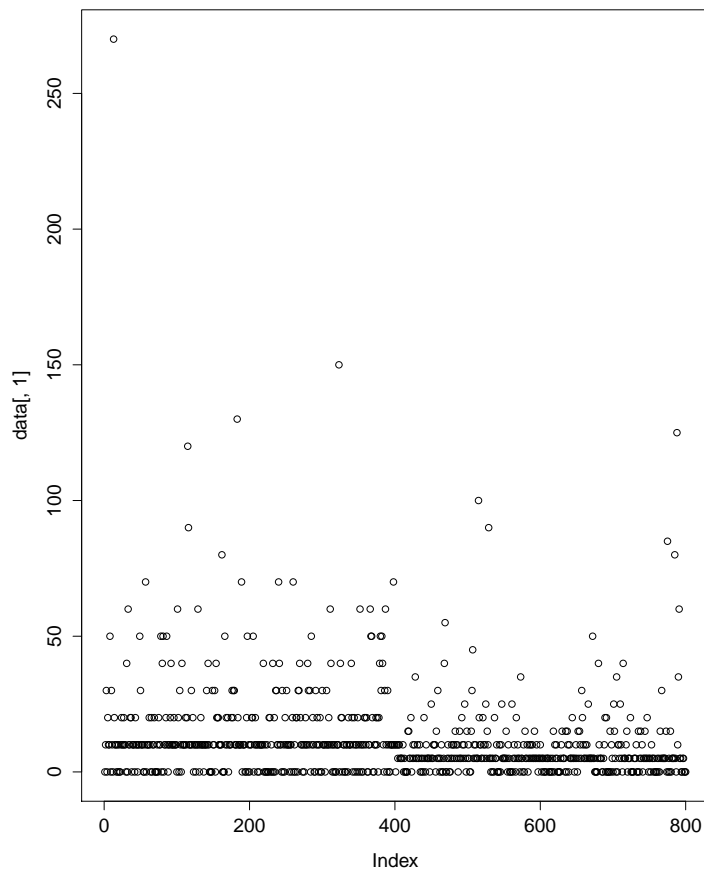


Figure 1: Plotting the numerical,1st column data

The y-axis of this plot are the values of the first column data, plotted against the index (row number) at which they occur in the dataset. This is a *scatter-plot* for a numerical data type

2. Plot of the 4th column (categorical data)

```
44 plot(data[,4])
```

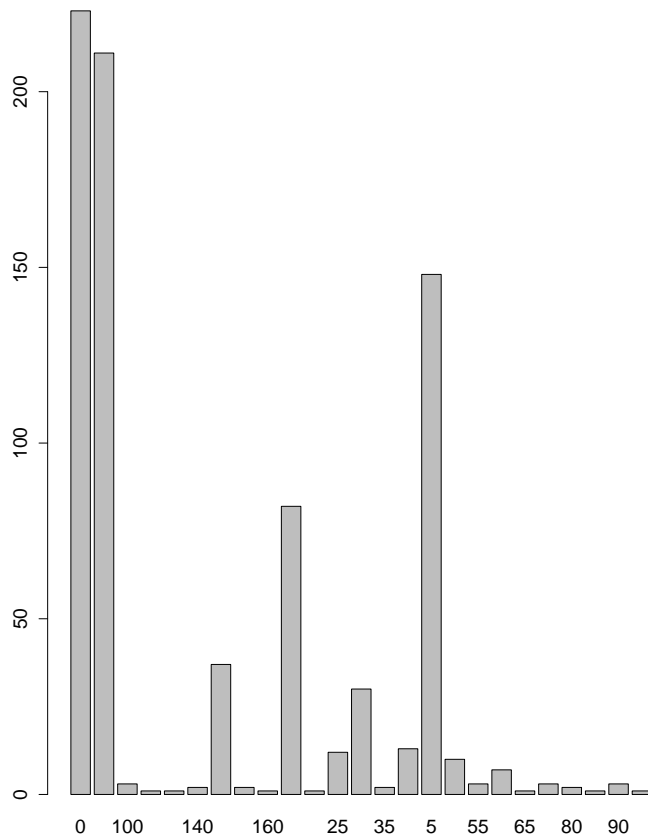


Figure 2: Plotting the categorical, 4th column data

This graph is a distribution: it is a *histogram* with the x-axis showing the values in the 4th column of our dataset, and the y-axis depicts the frequency with which values occur in a specified range. This graph is default since the 4th column data happens to be a categorical variable in our dataset.

2 Part 2 - HW01pb2data.csv

2.1 Extract a random sample of 10k observations

I used the *sample* command in *R* to obtain a smaller sample composed of 10,000 random records from our dataset.

Below are the R commands used to read and structure the data:

```
1 setwd("/Users/smadaan/Documents/ML_UCSC/week1/HW")
2 data<-read.csv("HW01pb2data.csv",header=FALSE)
3 str(data)
4 # 'data.frame': 2000000 obs. of 1 variable:
5
6 nrow(data)
7 #[1] 2000000
8
9 # selecting a sample of 10,000 random records
10 ss<-seq(1,nrow(data))
11 rand.ind<-sample(ss,10000,replace=F) #set of random indices for subset
12 small_data<-data[,1][rand.ind]
13 length(small_data)
14 #[1] 10000
```

2.2 Descriptive Stats on Sample

Below are the commands used to compute the mean, max, variance and 1st quartile from the data.

The resulting values are included as comments in the below code (mean =9.41002, max=16.93748, var=4.004991).

```
19 #mean, max and other descriptive stats
20 mean(small_data)
21 #[1] 9.41002
22 max(small_data)
23 #[1] 16.93748
24 var(small_data)
25 #[1] 4.004991
26 quantile(small_data,0.25)
27 #8.079612
```

2.3 Descriptive Stats for the Entire data

Below is the *R* code to compute descriptive stats for the entire dataset.

```
32 mean(data)
33 #9.451468
34 max(data)
35 # [1] 18.96657
36 var(data)
37 # 4.001822
38 quantile(data[,1],0.25)
39 # 8.10388
```

The population values for these parameters are very close to those that we found for the smaller sample. The table below shows the comparison.

Statistic	Small Sample	Entire-Dataset
Mean	9.41002	9.451468
Max	16.93748	18.96657
Variance	4.004991	4.001822
1st Quartile	8.079612	8.10388

Table 1: Comparison of statistics between the sample and entire-population datasets

3 Analysis of Ocean-view and Desert Home Prices

3.1 Box Plots for Pricing Data

R code for box plots:

```
5 boxplot(data.desert[,1], data.ocean[,1], col=c("red", "blue"),
6         main="House Box Plots: Comparison of House Prices",
7         names=c("Desert View", "Ocean View"), ylab="Prices (in thousand
           dollars)")
```

The BoxPlot shows the distribution of prices between the two kinds of houses. Desert View houses have a much lower price than Ocean-view houses. The edges of the box indicate the 1st and last quartiles— the median value (bold line) is much closer to the 25th percentile for desert view houses, indicating that the median is lower than the mean for desert houses. Finally, although the desert view houses are less expensive on average, the most expensive desert homes have a higher price (over 2.5 million dollars) than the most expensive ocean-view houses (less than 2.5 million dollars).

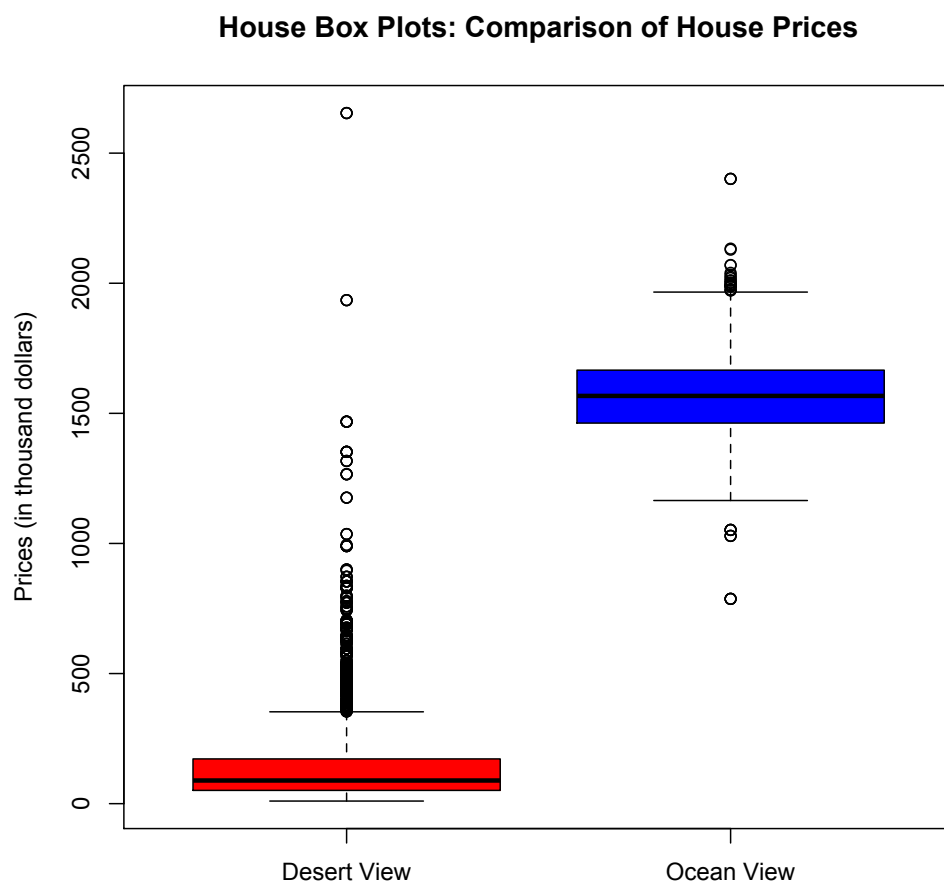


Figure 3: Comparison of prices between desert and ocean-view and desert houses

3.2 Frequency Histogram for Ocean-view Houses

R code for histogram:

```
11 hist(data.ocean[,1], breaks=seq(from=0,to=3000,by=500),  
12      xlab="Price_(in_thousand_dollars)", main="Price_Distribution_  
      of_Ocean-view_Houses")
```

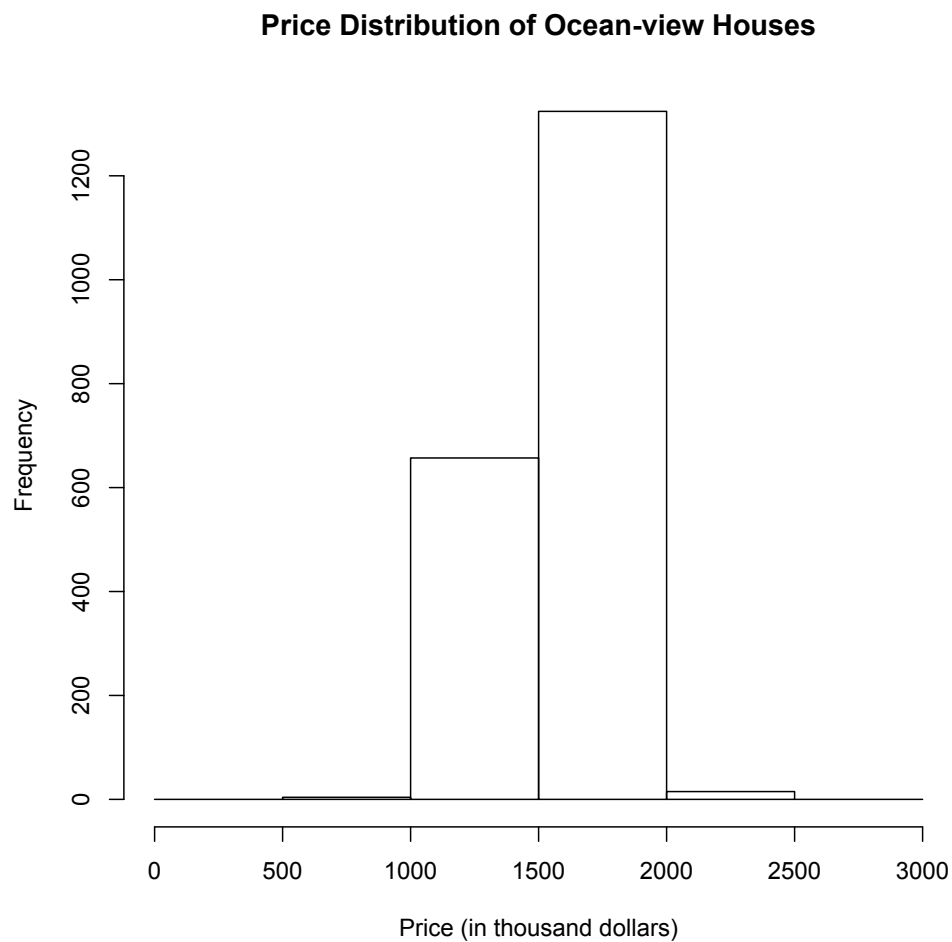


Figure 4: Histogram for prices of the ocean-view houses (0 to 3 million dollars)

3.3 ECDF for House Prices

R code for ecdf plots:

```
11 hist (data.ocean[,1], breaks=seq(from=0,to=3000,by=500),  
12       xlab="Price_(in_thousand_dollars)", main="Price_Distribution_  
      of_Ocean-view_Houses")
```

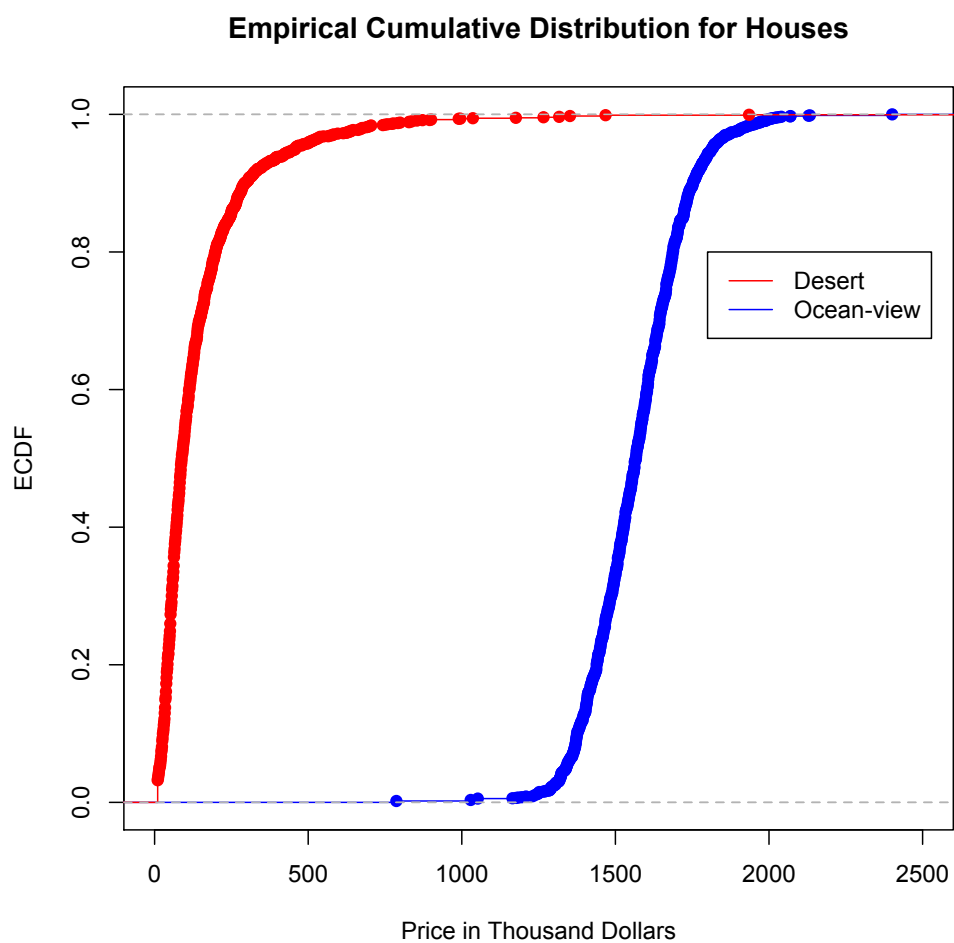


Figure 5: ECDF for Desert and Ocean-view Houses

4 Analyzing Data for Orange Trees

4.1 Age vs. Circumference relation

We first explore the data to identify the 3 variables: Tree, age, circumference. Using the *summary* command on the dataset outputs useful information, including the (min,max) parameters for age and circumference, and the different levels for 'Tree' (1,2,3,4,5).

x-axis (age) ranges from 118-152, and y-axis (circumference) ranges from 30-214.

Below is the *R* code for the same:

```
4 #— explore the data
5 head(orange)
6 levels(orange[,1])
7 str(orange)
8 summary(orange)
9
10 # 3 variables: Tree(factor), age (numeric), circumference (numeric)
11 # range for age: 118–1582
12 # range for circumference: 30–214
13
14 #— plot age vs circumference
15 plot(orange$age, orange$circumference, xlab="Age_of_Tree_(days_since_
      1968/12/31)",
16      ylab="Trunk_Circumference_(mm)", pch=20, main="Circumference_vs.
      Age_for_Trees",
17      col=orange$Tree)
18 legend('bottomright', legend = levels(factor(orange$Tree)),
19      text.col=seq_along(levels(orange$Tree)), title="Tree_
      Type")
```

The resulting plot shows the variation of tree circumference with age, for the different groups/levels of trees.

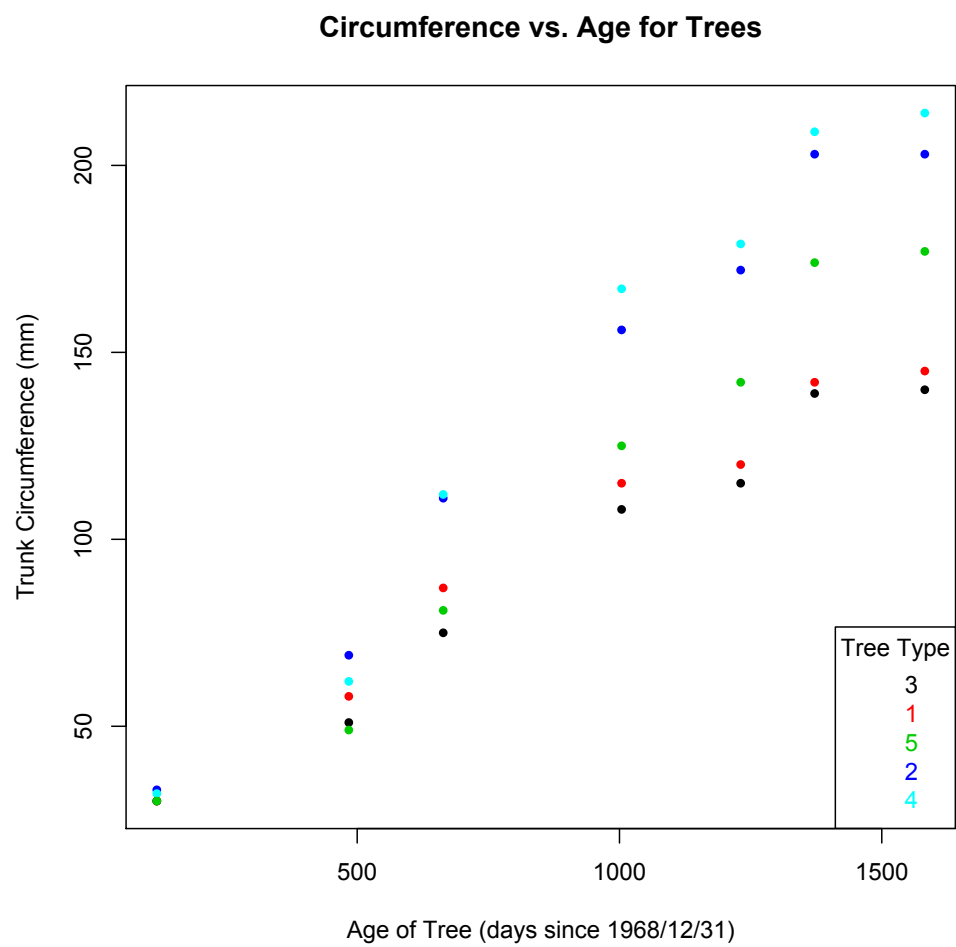


Figure 6: Age vs. Circumference, for the 5 groups of Orange Trees

4.2 Correlation between age and circumference for 1st Tree type

The two variables (age, circumference) are strongly correlated, with a correlation value of 0.9854675.

R code:

```
23 orange.1<-orange[which(orange$Tree==1),]  
24 cor(orange.1$age,orange.1$circumference)  
25 # [1] 0.9854675
```

4.3 Covariance and Correlation for Each Tree-type

Approach: Create an empty data frame with the different tree groups as the first column. Then, leverage the *by* command to calculate the parameters for each group. This saves having to individually calculate and *merge* results.

Below is the *R* code with results:

```
29 names(orange)  
30 t.levels<-sort(levels(orange$Tree))  
31  
32 stats<-data.frame(matrix(nrow=length(t.levels),ncol=0))  
33 stats$Tree<-t.levels  
34 stats$COVARIANCE<-as.matrix(by(orange, orange$Tree,  
35                               function(x){cov(x$age,x$circumference)}))  
36 stats$CORRELATION<-as.matrix(by(orange, orange$Tree,  
37                               function(x){cor(x$age,x$circumference)}))  
38 stats  
39 '  
40 >_stats  
41 _Tree_COVARIANCE_CORRELATION  
42 1____1____22239.83____0.9881766  
43 2____2____22340.07____0.9854675  
44 3____3____30442.81____0.9877376  
45 4____4____34290.45____0.9873624  
46 5____5____37062.62____0.9844610  
47 '
```

4.4 Effect of Adding 10 to each Circumference value

Correlation and covariance remain *unchanged*.

Below is the *R* code with results:

```

59 new.stats<-data.frame(matrix(nrow=length(t.levels),ncol=0))
60 new.stats$Tree<-t.levels
61 new.stats$COVARIANCE<-as.matrix(by(orange, orange$Tree,
62                                     function(x){cov(x$age,x$circumference+10)}))
63 new.stats$CORRELATION<-as.matrix(by(orange, orange$Tree,
64                                     function(x){cor(x$age,x$circumference+10)}))
65 new.stats
66 # no change in covariance or correlation
67 ,
68 >_new.stats
69 _Tree_COVARIANCE_CORRELATION
70 1____1____22239.83____0.9881766
71 2____2____22340.07____0.9854675
72 3____3____30442.81____0.9877376
73 4____4____34290.45____0.9873624
74 5____5____37062.62____0.9844610
75 ,

```

4.5 Effect of Doubling Circumference

Covariance *doubles* and correlation remains *unchanged*.

Below is the *R* code with results:

```

79 new2.stats<-data.frame(matrix(nrow=length(t.levels),ncol=0))
80 new2.stats$Tree<-t.levels
81 new2.stats$COVARIANCE<-as.matrix(by(orange, orange$Tree,
82                                     function(x){cov(x$age,x$circumference*2)}))
83 new2.stats$CORRELATION<-as.matrix(by(orange, orange$Tree,
84                                     function(x){cor(x$age,x$circumference*2)}))
85 new2.stats
86 ,
87 >_new2.stats
88 _Tree_COVARIANCE_CORRELATION
89 1____1____44479.67____0.9881766
90 2____2____44680.14____0.9854675
91 3____3____60885.62____0.9877376
92 4____4____68580.90____0.9873624
93 5____5____74125.24____0.9844610
94 ,
95 # covariance doubles, correlation remains same

```

4.6 Effect of Multiplying Circumference by -2

Covariance becomes negative (and double in magnitude compared to the original dataset), and correlation becomes *negative* of its original value. Results:

New covariance = original covariance * (-2)

New correlation = -original correlation

Below is the *R* code with results:

```
99 new3.stats<-data.frame(matrix(nrow=length(t.levels),ncol=0))
100 new3.stats$Tree<-t.levels
101 new3.stats$COVARIANCE<-as.matrix(by(orange, orange$Tree,
102                                     function(x){cov(x$age,x$circumference*(-2))}))
103 new3.stats$CORRELATION<-as.matrix(by(orange, orange$Tree,
104                                     function(x){cor(x$age,x$circumference*(-2))}))
105 new3.stats
106 '
107 >_new3.stats
108 _Tree_COVARIANCE_CORRELATION
109 1____1____-44479.67____-0.9881766
110 2____2____-44680.14____-0.9854675
111 3____3____-60885.62____-0.9877376
112 4____4____-68580.90____-0.9873624
113 5____5____-74125.24____-0.9844610
114 '
115 # new.corr = orig.corr*-2, correlation becomes -ve
```

5 Desert Homes: Revisited

5.1 Calculating the Median Price

Approach: *R* does not have a native *median* function, so we define a custom *median* function to calculate the median home price. This results in:

Below is the *R* code with results:

```
11
12 median = function (data.vector)
13 {
14     median=NA;
15     l = length(data.vector);
16     sv<-sort(data.vector);
```

```

17         median=ifelse ( is.integer ( 1/2 ) , sv [ 1/2 ] , sv [ 1/2+1 ] ) ;
18         return ( median ) ;
19     }
20
21
22     median ( data . desert [ , 1 ] )
23     #89
24     mean ( data . desert [ , 1 ] )
25     #144.0348

```

Median value: 89

Mean value: 144.0348

(note: the values are in units of thousand dollars).

Median is *smaller* than the mean.

5.2 Characteristics of the Distribution

We found earlier that the median is lower than the mean. This generally indicates that the distribution is positively, or *right-skewed*.

This can be further confirmed by plotting a histogram to get a sense of the distribution, as shown below.

```

29 hist ( data . desert [ , 1 ] , main="Distribution of Desert Home Prices" ,
30        xlab="Home Price" , ylab="Number of Homes" )
31 # also , median < mean , positive / right skewed

```

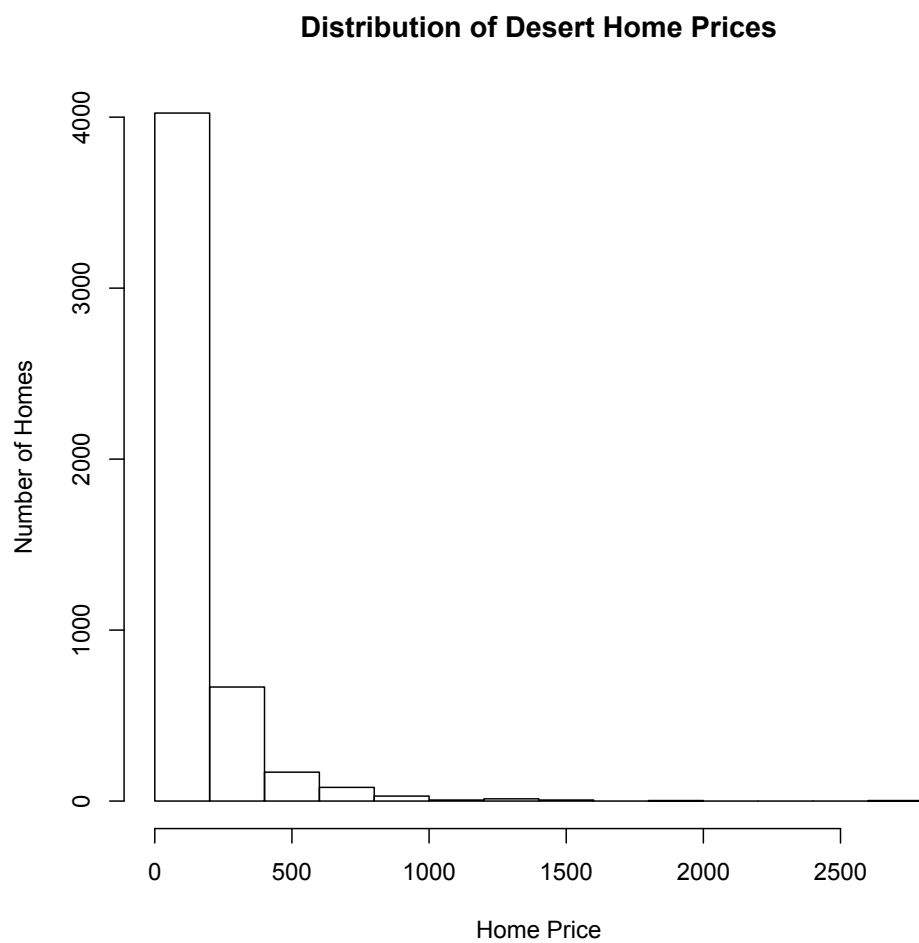


Figure 7: Histogram showing desert home prices

5.3 Effect of Higher Home Prices on Median Value

If we increase house price by 10 (thousand dollars) each, the median value *increases* by the same amount.

Specifically, the new median (99k) is higher than the old median (89k) by 10 (thousand dollars). Below is the *R* code with results:

```
35 median(data.desert[,1]+10)
36 #99
37 #median also increases by 10k
```

5.4 Effect of Doubling Home Prices on Median Value

If home prices double, the median value also *doubles*.

Specifically, the new median (178k) is twice the old median (89k). Below is the *R* code with results:

```
41 median(data.desert[,1]*2)
42 #178
43 #median also doubles
```

6 Closing Notes and References

I read material from *The Art of R Programming* by Norman Matloff for an introduction to *R*, and googled online references for the use of the *by* command.

Tools used include the *R* package for Mac, *TeXShop* for writing Latex files. I did *not* collaborate with another class member on this assignment.

Included in my submission is this *pdf* file, and all the *R* files where my code is programmed.