

```

1 ##### Exercise 2 #####
2
3 ## Question a: Load the file zeta.csv of income data into R.
4 ```{r}
5 setwd("/Users/tenzintashi/Downloads/CSc 460 - DS/Assignment")
6 # file <- file.choose()
7 # zeta <- read.csv(file)
8 zeta <- read.csv(file='Dataset/zeta.csv', header=TRUE, sep=",")
9 View(zeta)
10 ```
11 ## Question b: Change the column names of your data frame so that zeta becomes zipCode and

```

	X	zcta	sex	meanage	meaneducation	meanemployment	meanhouseholdincome
1	1	602	F	37.40335	10.912822	0.7400294	18533.84
2	2	602	M	35.93574	10.692618	1.3438833	18533.84
3	3	604	F	31.80943	13.913371	1.0858555	40784.49
4	4	604	M	31.10425	14.264654	1.6025594	40784.49
5	5	606	F	35.99079	10.097773	0.6287526	17496.53
6	6	606	M	35.99027	9.927995	1.2169643	17496.53
7	7	610	F	37.26014	10.969157	0.8543247	19416.41
8	8	610	M	36.80649	10.850965	1.2955239	19416.41
9	9	612	F	40.42732	11.575772	0.7815393	21607.34
10	10	612	M	37.69992	11.391897	1.2723553	21607.34

```

11 ## Question b: Change the column names of your data frame so that zeta becomes zipCode and
12 ## meanhouseholdincome becomes income.
13 ```{r}
14 names(zeta)[2] <- "zipCode"
15 names(zeta)[7] <- "income"
16 View(zeta)
17 ```

```

	X	zipCode	sex	meanage	meaneducation	meanemployment	income
1	1	602	F	37.40335	10.912822	0.7400294	18533.84
2	2	602	M	35.93574	10.692618	1.3438833	18533.84
3	3	604	F	31.80943	13.913371	1.0858555	40784.49
4	4	604	M	31.10425	14.264654	1.6025594	40784.49
5	5	606	F	35.99079	10.097773	0.6287526	17496.53
6	6	606	M	35.99027	9.927995	1.2169643	17496.53

```
18 ## Question c: Analyze the summary of your data.
```

```
19 ```{r}
```

```
20 summary(zeta)
```

```
21 ```
```

X	zipCode	sex	meanage	meaneducation	meanemployment
Min. : 1	Min. : 601	Length:64076	Min. : 0.00	Min. : 0.00	Min. : 0.000
1st Qu.:16020	1st Qu.:27305	Class :character	1st Qu.: 36.65	1st Qu.:11.91	1st Qu.:1.542
Median :32038	Median :49909	Mode :character	Median : 39.30	Median :12.46	Median :1.813
Mean :32038	Mean :49801		Mean : 39.68	Mean :12.53	Mean :1.787
3rd Qu.:48057	3rd Qu.:72007		3rd Qu.: 42.28	3rd Qu.:13.11	3rd Qu.:2.077
Max. :64076	Max. :99950		Max. :137.08	Max. :19.00	Max. :3.000

income

Min. : 0

1st Qu.: 37642

Median : 44163

Mean : 48245

3rd Qu.: 54373

Max. :250000

```
22 ## What are the mean and median average incomes?
```

```
23 ```{r}
```

```
24 # Mean average Income = 48245
```

```
25 # Median average Income = 44163
```

```
26 ```
```

```
27 ## Question d: Plot a scatter plot of the data. Although this graph is not too informative,
```

```
28 ## do you see any outlier values? If so, what are they?
```

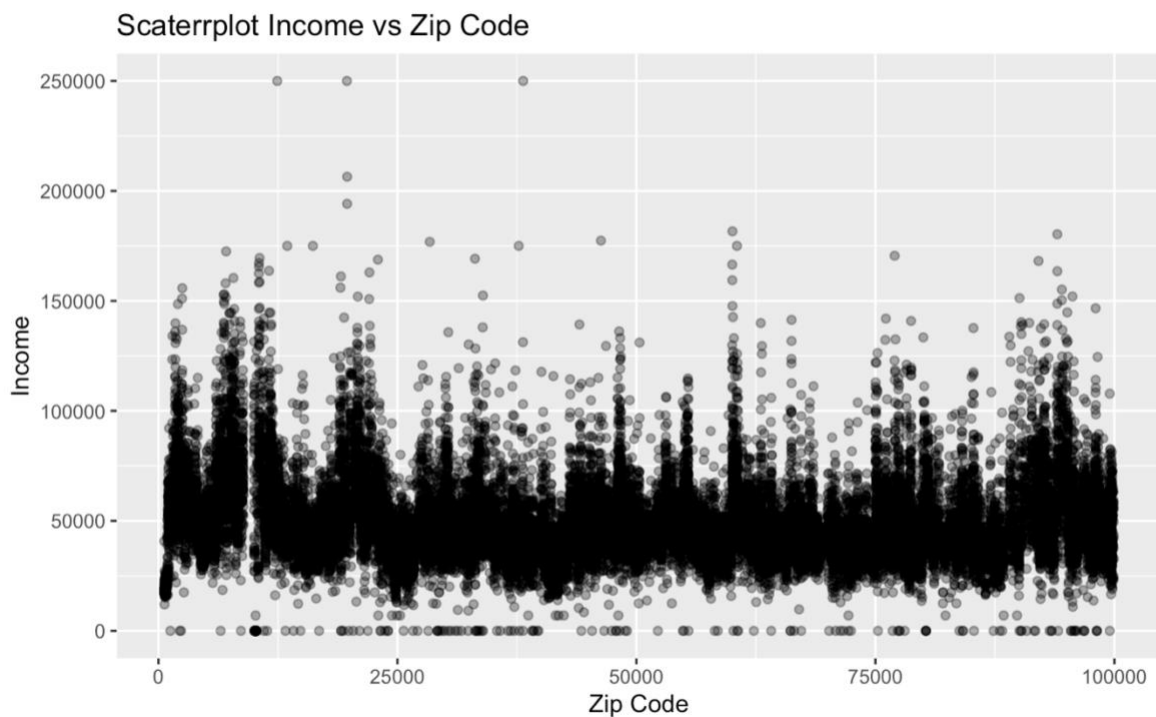
```
29
```

```
30 ```{r}
```

```
31 library(ggplot2)
```

```
32 ggplot(zeta,aes(x= zipCode, y=income)) +geom_point(alpha=0.2) +labs(x="Zip  
Code",y="Income",title="Scaterrplot Income vs Zip Code")
```

```
33 ```
```



```
34 There seem to be two outlier values are 0 and 250000
```

```
35
```

```

36 ## Question e: In order to omit outliers, create a subset of the data so that: $7,000 < income < $200,000,
37 ## What's your new mean?
38 ```{r}
39 newData <- subset(zeta, income <200000 & income >7000)
40 summary(newData)
41 ```

```

```

      X      zipCode      sex      meanage      meaneducation      meanemployment
Min.   : 1      Min.   : 601      Length:63742      Min.   : 0.00      Min.   : 0.00      Min.   :0.000
1st Qu.:16046      1st Qu.:27333      Class :character      1st Qu.: 36.69      1st Qu.:11.91      1st Qu.:1.546
Median :32076      Median :49935      Mode  :character      Median : 39.31      Median :12.46      Median :1.816
Mean   :32051      Mean   :49817                      Mean   : 39.78      Mean   :12.57      Mean   :1.795
3rd Qu.:48047      3rd Qu.:72003                      3rd Qu.: 42.28      3rd Qu.:13.11      3rd Qu.:2.078
Max.   :64076      Max.   :99950                      Max.   :133.11      Max.   :19.00      Max.   :3.000

      income
Min.   : 8465
1st Qu.: 37755
Median : 44234
Mean   : 48465
3rd Qu.: 54444
Max.   :194135

```

```

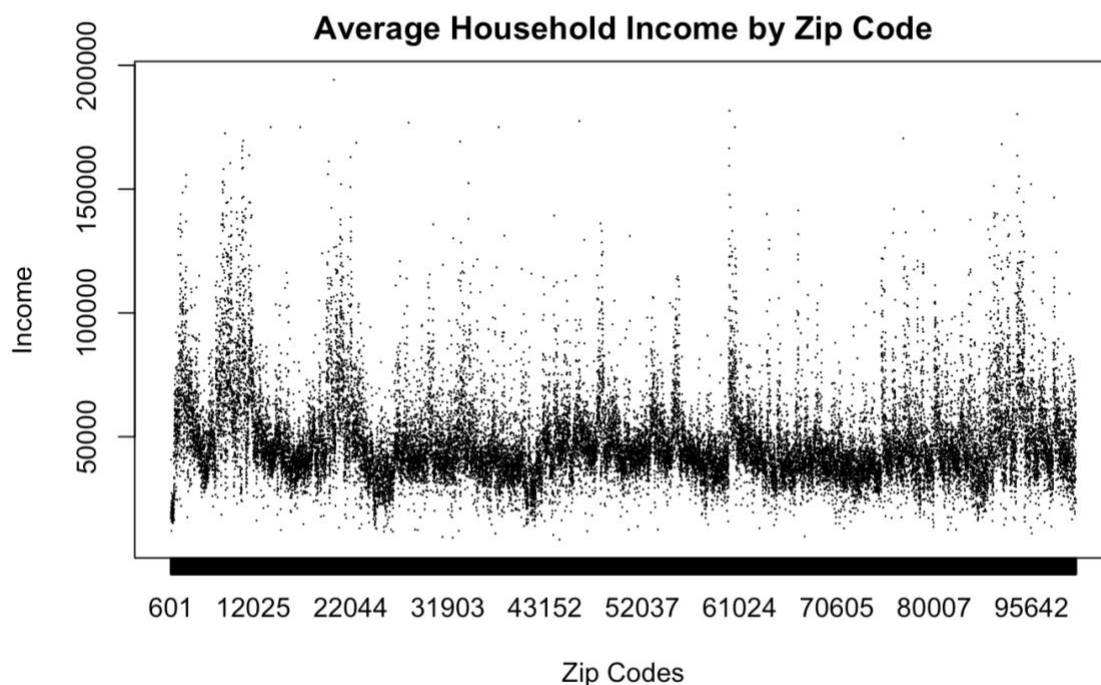
42 ```{r}
43 # New Mean average Income = 48465
44 ```

```

```

45 ## Question f: Create a simple box plot of your data. Be sure to add a title and label the axes.
46 ```{r}
47 boxplot(col="white", data = newData, income ~ zipCode, main = "Average Household Income by Zip Code", xlab =
  "Zip Codes", ylab = "Income")
48
49 ```

```



50

51 ▾ ## Question g: In the box plot you created, notice that all of the income data is pushed towards

52 ▾ ## the bottom of the graph because most average incomes tend to be low. Create a new box plot

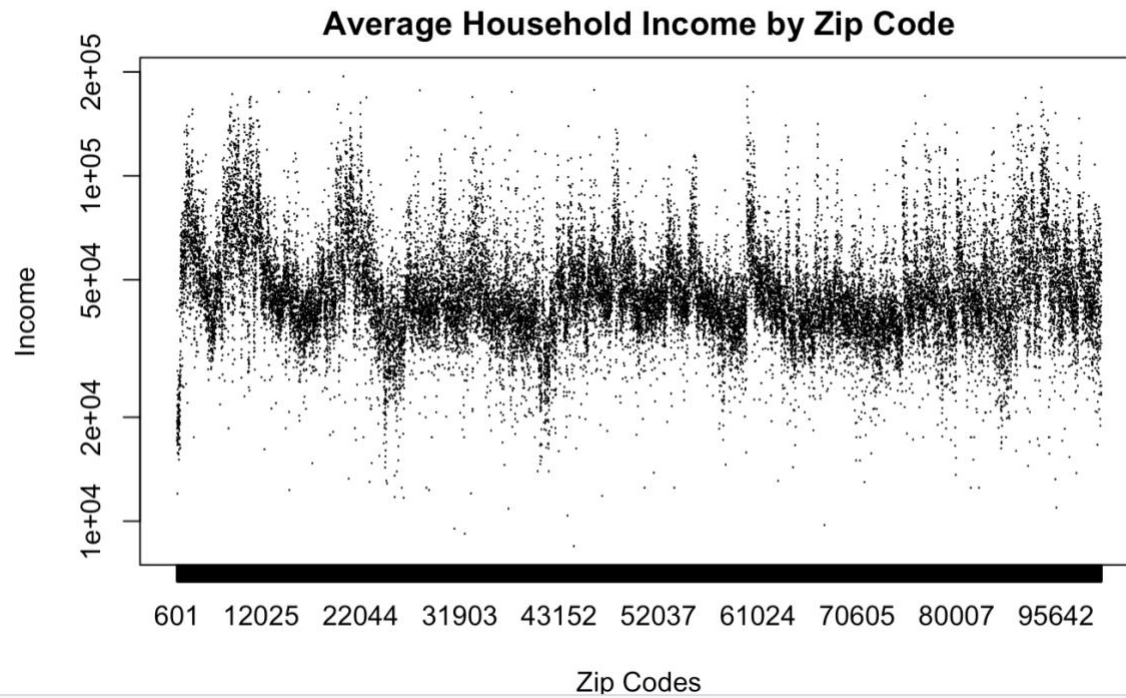
53 ▾ ## where the y- axis uses a log scale. Be sure to add a title and label the axes.

54 ▾ ```{r}

55 #Create a new box plot where the y-axis uses a log scale

56 boxplot(col="white", data = newData, income ~ zipCode, main = "Average Household Income by Zip Code", xlab =
"Zip Codes", ylab = "Income", log='y')

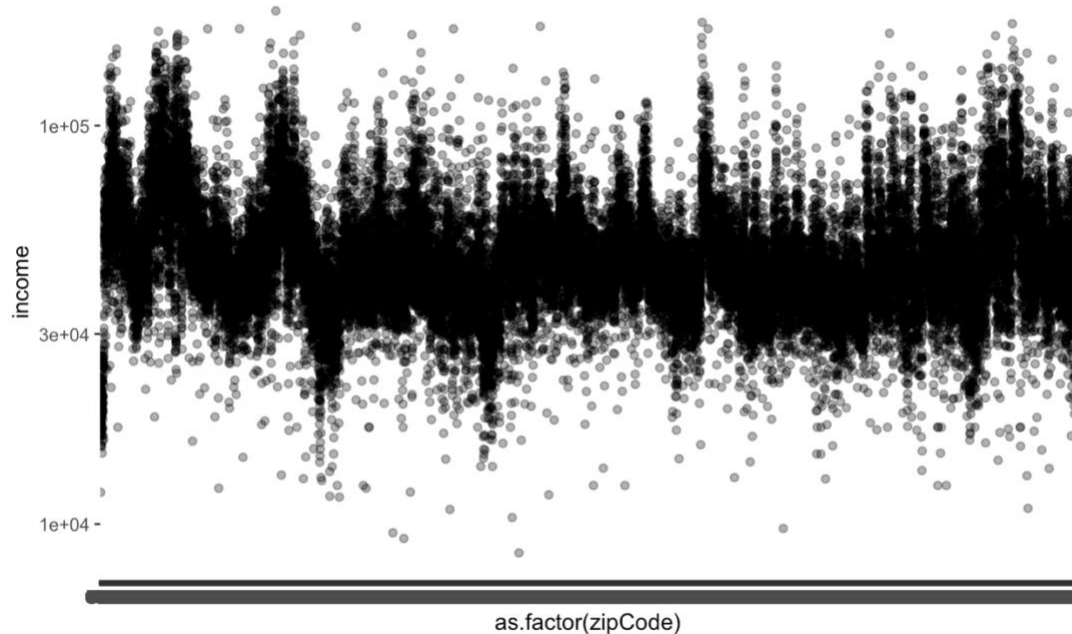
57 ▴ ```




```

57 # Question h: Use the ggplot library in R, which enables you to create graphs with several different
58 # types of plots layered over each other. Be sure to read the documentation for ggplot
59 # and load the library ggplot2 (you may have to install this package into R).
60 {r}
61 library(ggplot2)
62 # Make a ggplot that consists of just a scatter plot using the function geom_point() with position = "jitter"
63 ggplot(newData, aes(x = as.factor(zipCode), y = income)) + geom_point(position = "jitter", alpha = 0.2) + scale_y_log10()
64 {r}

```



```

74 # Question i: Make a ggplot that consists of just a scatter plot using the function geom_point()
75 # with position = "jitter" so that the data points are grouped by zip code. Be sure to use ggplot's
76 # function for taking the log10 of the y-axis. (Hint: for geom_point, have alpha=0.2).
77 {r}
78 library(ggplot2)
79 ggplot(newData, aes(x = as.factor(newData$zipCode), y = newData$income)) + geom_point(aes(colour = factor(zipCode)), position =
80 'jitter', alpha = 0.2) + geom_boxplot(alpha = 0.1, outlier.size = -Inf) + scale_y_log10() + labs(color = "Region", x = "Zip
81 Code", y = "Income", title = "Average Income by Zip Code") + theme(plot.title = element_text(size = 11, face = "plain", hjust = 0.5))
82 {r}
83 # Question j: Create a new ggplot by adding a box plot layer to your previous graph.
84 # To do this, add the ggplot function geom_boxplot(). Also, add color to the scatter plot so
85 # that data points between different zip codes are different colors. Be sure to label the axes and
86 # add a title to the graph. (Hint: for geom_boxplot, have alpha=0.1 and outlier.size=0).
87 {r}
88 library(ggplot2)
89 ggplot(newData, aes(x = as.factor(zipCode), y = income)) + geom_point(aes(colour = factor(zipCode)), position = 'jitter', alpha = 0.2) +
90 geom_boxplot(alpha = 0.1, outlier.size = 0) + scale_y_log10() + ylab("Income") + xlab("Zip Code") + ggtitle("Average Income by Zip
91 Code") + labs(color = "Region") + theme(plot.title = element_text(size = 11, face = "plain", hjust = 0.5))
92 {r}
93 # Question k: What can you conclude from this data analysis/visualization?
94 {r}
95 # - It is important to visualize your data in different ways.
96 # - Visualization enables you to better understand what your data is telling you.
97 # - Visualization enables you to better communicate your results to stakeholders.
98 # - Zip codes starting in 0 (New England) and 9 (West Coast) have higher average household incomes.
99 {r}

```



```

96 ▾ ## Question k: What can you conclude from this data analysis/visualization?
97 ▾ ```{r}
98 # - It is important to visualize your data in different ways.
99
100 # - Visualization enables you to better understand what your data is telling you.
101
102 # - Visualization enables you to better communicate your results to stakeholders.
103
104 # - Zip codes starting in 0 (New England) and 9 (West Coast) have higher average household incomes.
105 ▴ ```
106
107

```

```

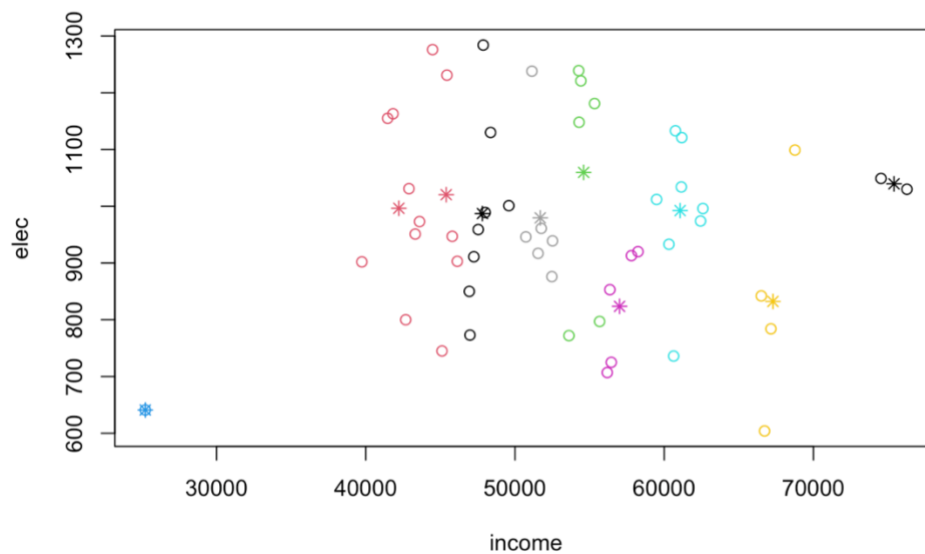
100 ▾
101 ▾ ##### Exercise 2.2 #####
102

```

```

103 ▾ ## Question a: Cluster the data and plot all 52 data points, along with the centroids.
104 ▾ ## Mark all data points and centroids belonging to a given cluster with their own color. Here, let k=10.
105 ▾ ```{r}
106 load('income_elec_state.Rdata')
107 k = kmeans(income_elec_state, 10)
108 plot(income_elec_state, col = k$cluster)
109 points(k$centers, col=1:10, pch=8)
110 ▴ ```

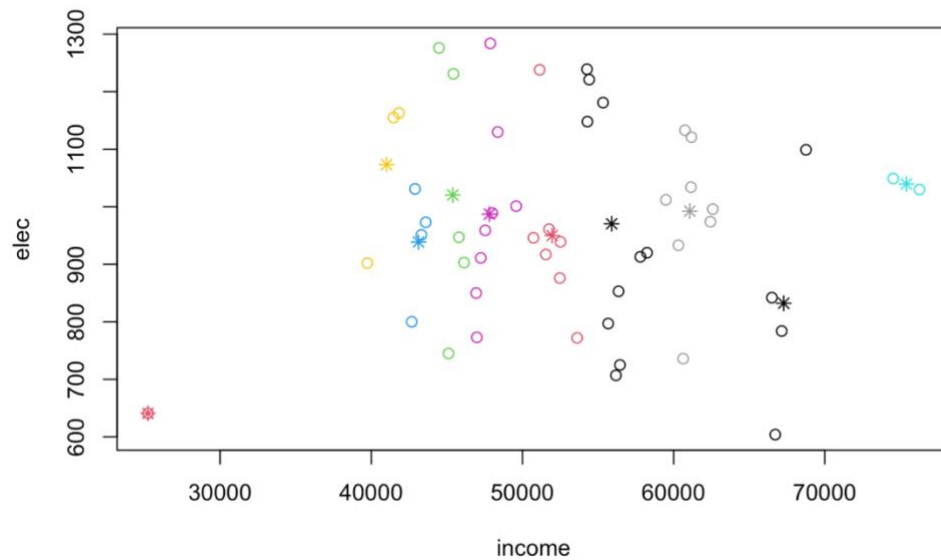
```



```

111 ## Question b: Repeat step (a) several times. What can change each time you cluster the data?
112 ## Why? How do you prevent these changes from occurring?
113 ```{r}
114 k = kmeans(income_elec_state, 10)
115 plot(income_elec_state, col = k$cluster)
116 points(k$centers, col=1:10, pch=8)
117 ^

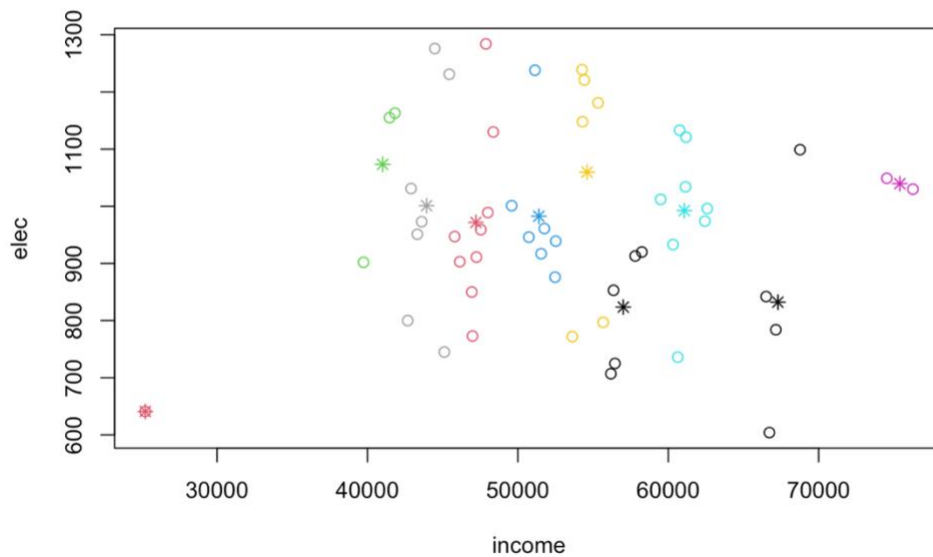
```



```

118 
119 ```{r}
120 k = kmeans(income_elec_state, 10)
121 plot(income_elec_state, col = k$cluster)
122 points(k$centers, col=1:10, pch=8)
123 ^

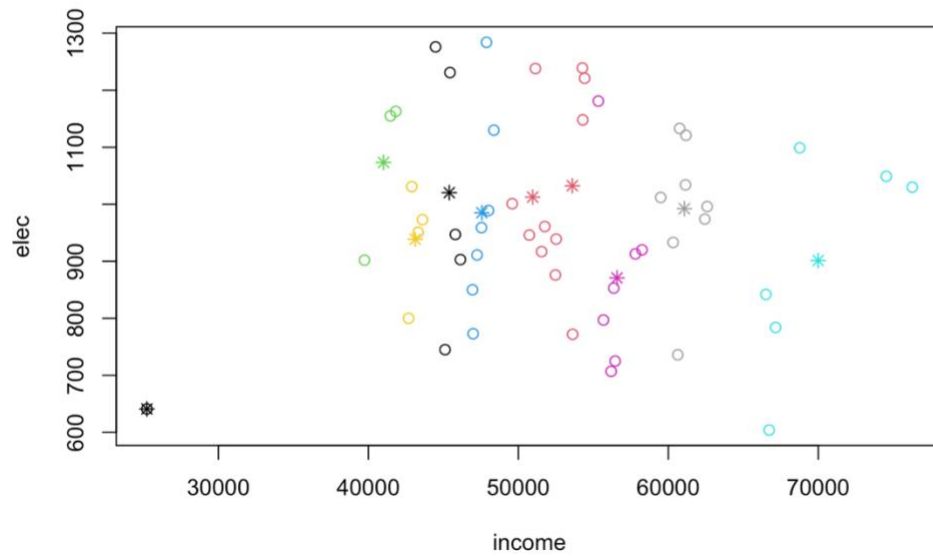
```




```

124
125 ````{r}
126 k = kmeans(income_elec_state, 10)
127 plot(income_elec_state, col = k$cluster)
128 points(k$centers, col=1:10, pch=8)
129 ````

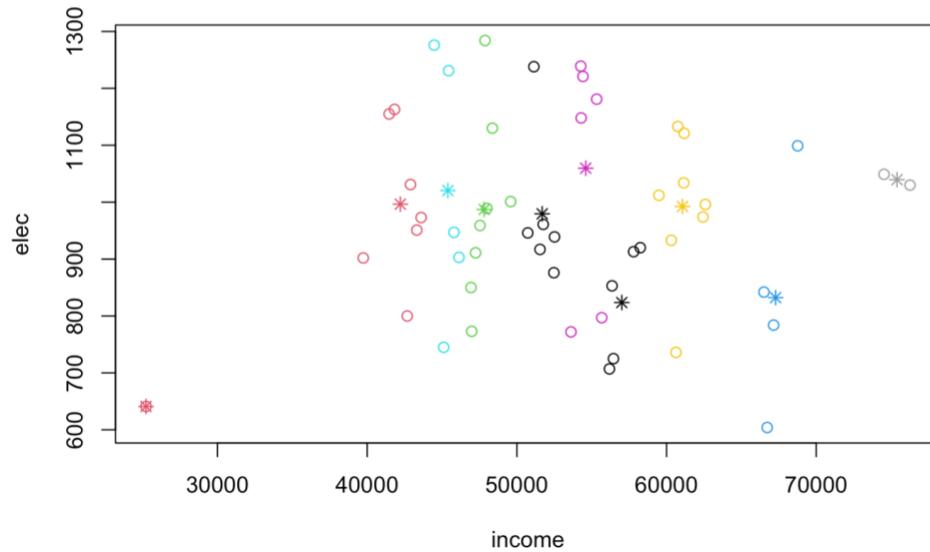
```



```

130
131 ```{r}
132 # Sizes, centers' position, sum of squares of clusters can change after each time repeat above step.
133 # Because by default nstart = 1: having only one random starting set can result in different
134 # clusterings over multiple runs.
135
136 # To prevent these changes from occurring, we can:
137 # - Increase "nstart" to improve the likelihood of obtaining the globally optimal clustering.
138 # - Increasing the "iter.max" parameter reduces the likelihood that the kmeans algorithm terminates
139 # prematurely.
140 ```
141
142 ```{r}
143 k = kmeans(income_elec_state, 10, nstart=100, iter.max = 50)
144 plot(income_elec_state, col = k$cluster)
145 points(k$centers, col=1:10, pch=8)
146 ```

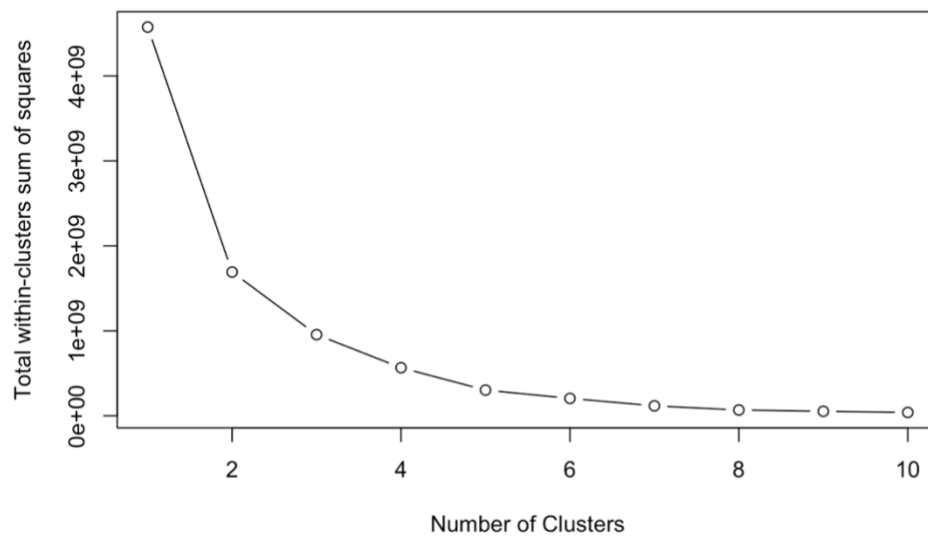
```



```

147 ## Question c: Once you've accounted for the issues in the previous step, determine a reasonable
148 ## value of k. Why would you suggest this value of k?
149 ```{r}
150 wss = numeric(10)
151 for (i in 1:10) wss[i] = sum(kmeans(income_elec_state, centers=i, nstart = 100, iter.max = 50)$tot.withinss)
152 plot(1:10, wss, type="b", xlab="Number of Clusters", ylab="Total within-clusters sum of squares")
153 ```

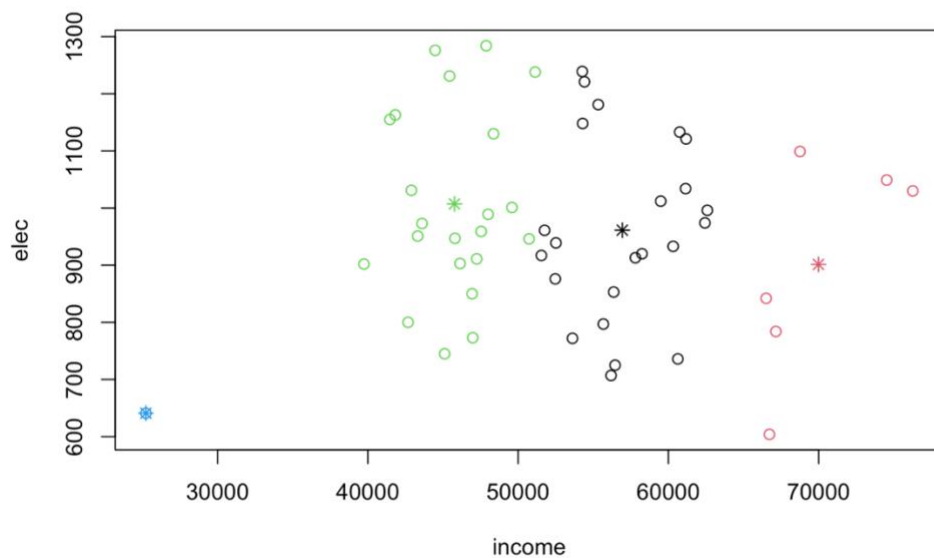
```



```

154 We can see here "elbow" = 4
155
157 With k=4
158 ```{r}
159 k = kmeans(income_elec_state, 4, nstart=100, iter.max = 50)
160 plot(income_elec_state, col = k$cluster)
161 points(k$centers, col=1:4, pch=8)
162
163 ```

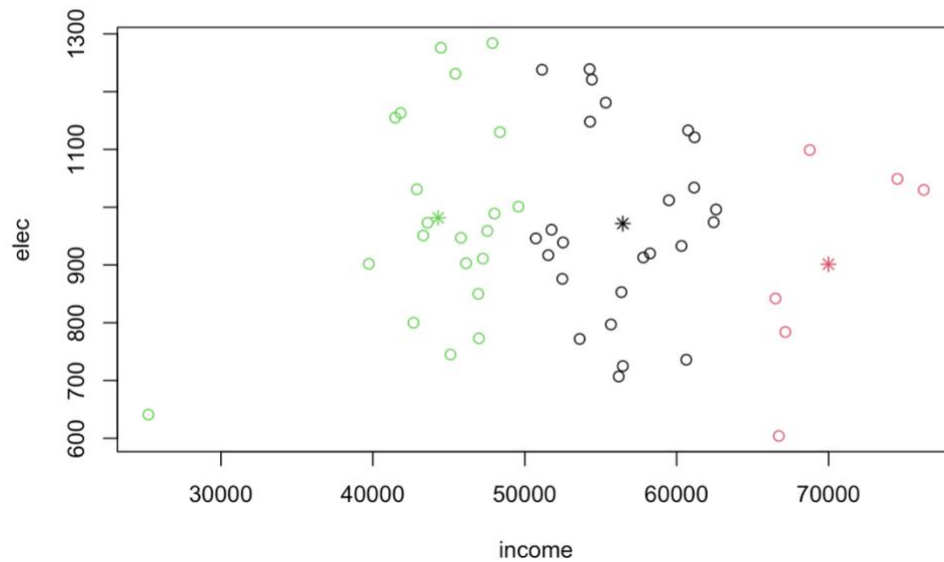
```



```

164 Repeat the modeling with k=3
165
166 ```{r}
167 k = kmeans(income_elec_state, 3, nstart=100, iter.max = 50)
168 plot(income_elec_state, col = k$cluster)
169 points(k$centers, col=1:4, pch=8)
170 ^

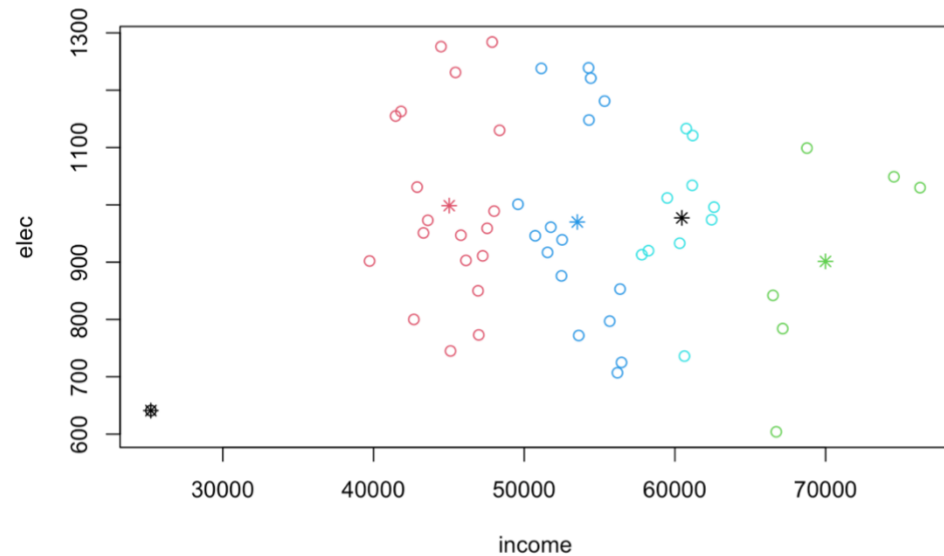
```



```

171 Repeat the modeling with k=5
172
173 ```{r}
174 k = kmeans(income_elec_state, 5, nstart=100, iter.max = 50)
175 plot(income_elec_state, col = k$cluster)
176 points(k$centers, col=1:4, pch=8)
177 ^

```



```

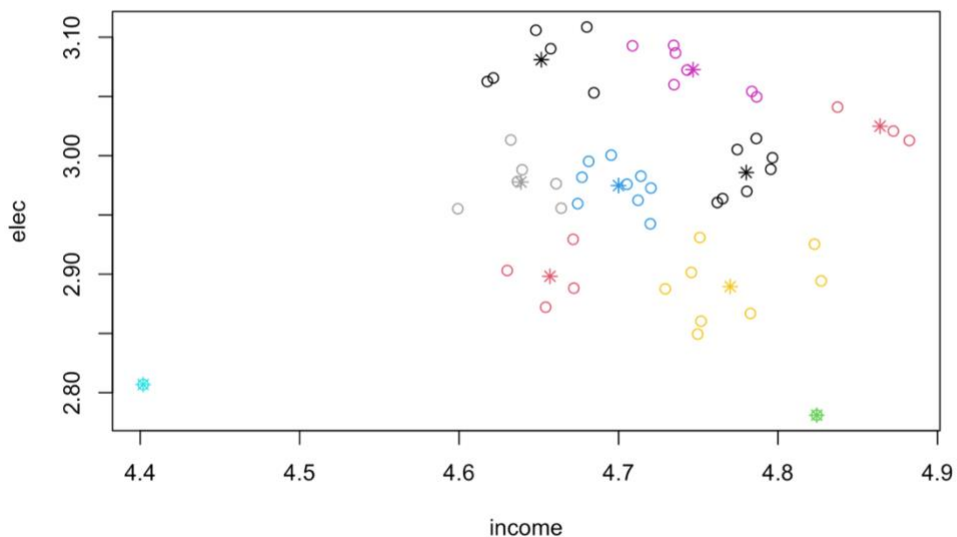
178
179 Chosen k=4. Because we see that Puerto Rico is an outlier, and should perhaps belong to its own cluster. It is the smallest k such
180 that Puerto Rico belongs to its own cluster, so this k would be a good value to suggest.

```

```

183- ## Question d: Convert the mean household income and mean electricity usage to a log10 scale and
184- ## cluster this transformed dataset. How has the clustering changed? Why?
185- ```{r}
186- new = log10(income_elec_state)
187- k = kmeans(new, 10, nstart=100, iter.max = 50)
188- plot(new, col = k$cluster)
189- points(k$centers, col=1:10, pch=8)
190- ```

```

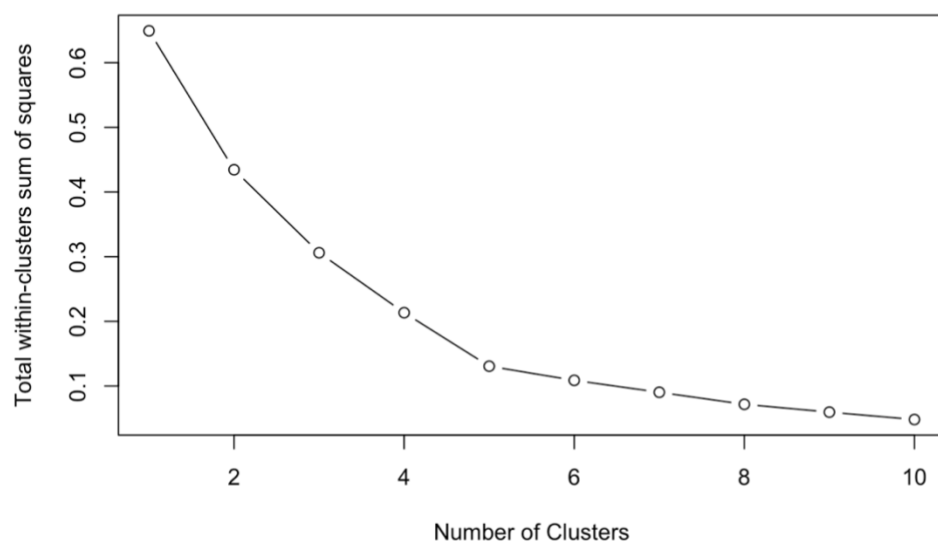


191 K-means clustering is not scale-invariant, so any adjustments made to the units of the data may impact the clustering.

```

192
193- ## Question e: Reevaluate your choice of k. Would you now choose k differently? Why or why not?
194- ```{r}
195- wss = numeric(10)
196- for (i in 1:10) wss[i] = sum(kmeans(new, centers=i, nstart = 100, iter.max = 50)$tot.withinss)
197- plot(1:10, wss, type="b", xlab="Number of Clusters", ylab="Total within-clusters sum of squares")
198- ```

```



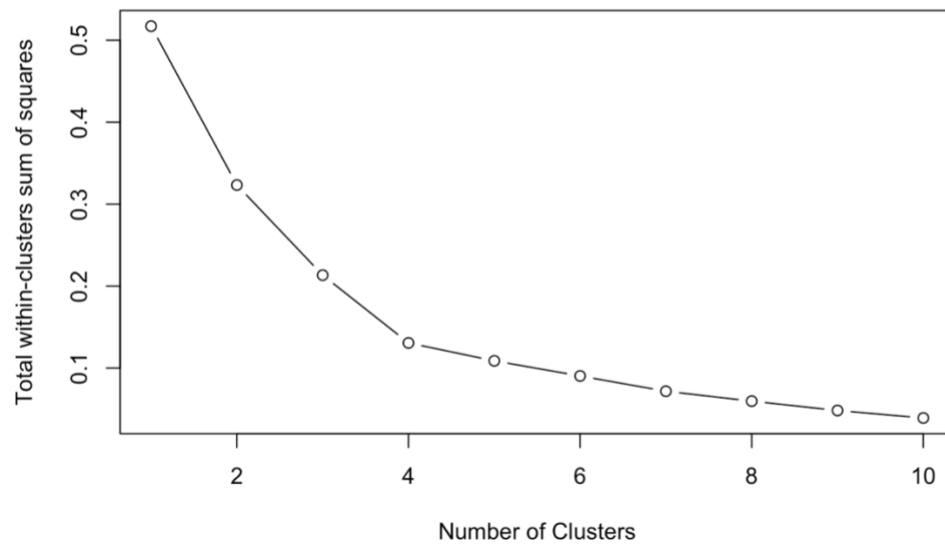
199 We see more clear elbow in the different position: k=5

200

```

201 # Question f: Have you observed an outlier in the data? Remove the outlier and, once again, reevaluate your
202 # choice of k
203 ```{r}
204 new <- subset(new, rownames(new) != "PR")
205 |
206 wss = numeric(10)
207 for (i in 1:10) wss[i] = sum(kmeans(new, centers=i, nstart = 100, iter.max = 50)$tot.withinss)
208 plot(1:10, wss, type="b", xlab="Number of Clusters", ylab="Total within-clusters sum of squares")
209 ```

```



210 After removing the outliers, it is clear that elbow on the plot change its position to smaller value. $k=4$
 211


```

212 ## Question g: Color a map of the U.S. according to the clustering you obtained. To simplify this task, use
213 ## the "maps" package and color only the 48 contiguous states and Washington D.C.
214 ```{r}
215 library(maps)
216
217 km <- kmeans(new,4,nstart = 100, iter.max = 50)
218 #Prepare vector with state order
219 map_order <- c('AL', 'AZ', 'AR', 'CA', 'CO', 'CT', 'DE', 'DC', 'FL',
220 'GA', 'ID', 'IL', 'IN', 'IA', 'KS', 'KY', 'LA', 'ME',
221 'MD', 'MA', 'MA', 'MA', 'MI', 'MI', 'MN', 'MS', 'MO',
222 'MT', 'NE', 'NV', 'NH', 'NJ', 'NM', 'NY', 'NY', 'NY',
223 'NY', 'NC', 'NC', 'NC', 'ND', 'OH', 'OK', 'OR', 'PA',
224 'RI', 'SC', 'SD', 'TN', 'TX', 'UT', 'VT', 'VA', 'VA',
225 'VA', 'WA', 'WA', 'WA', 'WA', 'WA', 'WA', 'WV', 'WI', 'WY')
226 #Prepare color vector
227 map_color <- km$cluster[map_order]
228 map('state', col = map_color,fill=TRUE)
229 ```

```

