

Mail Case Study

Rana Sahrir
10/01/2021

Investigate the dataset

summary(mall)				
<pre>## CustomerID Age Annual.Income.k... ## Min. : 1.89 Length:598 Min. :35.00 Min. : 15.00 ## 1st Qu.: 58.75 Class :character 1st Qu.:28.75 1st Qu.: 41.50 ## Median :108.50 Mode :character Median :36.00 Median : 61.50 ## Mean :188.50 Mean :38.85 Mean : 69.56 ## 3rd Qu.:158.25 3rd Qu.:49.00 3rd Qu.: 78.00 ## Max. :286.00 Max. :78.00 Max. :137.00 ## Spending.Score..1.100. ## Min. : 1.08 ## 1st Qu.:14.75 ## Median :59.08 ## Mean :50.28 ## 3rd Qu.:73.08 ## Max. :99.08</pre>				
str(mall)				
<pre>## 'data.frame': 208 obs. of 5 variables: ## \$ CustomerID : int 1 2 3 4 5 6 7 8 9 10 ... ## \$ Gender : chr "Male" "Male" "Female" "Female" ... ## \$ Age : int 18 21 28 23 21 22 35 23 34 38 ... ## \$ Annual.Income.k... : int 15 15 16 16 17 17 18 18 19 19 ... ## \$ Spending.Score..1.100.: int 39 81 6 77 40 76 6 94 3 72 ...</pre>				


Prepare the data for analysis

<pre># data preparation #normalize data normalize <- function(x) {return ((x - min(x)) / (max(x) - min(x)))} mall\$age=normalize(mall\$age) mall\$Annual.Income.k.=normalize(mall\$Annual.Income.k.) mall\$Spending.Score..1.100.= normalize(mall\$Spending.Score..1.100.) #remove unecessary columns mall\$CustomerID=NULL #remove duplicates mall=unique(mall)</pre>				
--	--	--	--	--

Hierarchical Clustering

<pre>#Start the Hierarchical clustering #Compute distance distances=dist(mall[,2:4], method="euclidean") clusterall=hclust(distances, method="ward.D")</pre>				
---	--	--	--	--

Plot The Dendrogram

<p>Cluster Dendrogram</p>  <p>hclust ("ward.D")</p> <pre>#design clusters clustergroups= cutree(clusterall, k=6) #investigate the clusters apply(mall\$Age, clustergroups, mean)</pre> <pre>## 1 2 3 4 5 6 ## 45. 21.739 24.88952 25.97143 54.59574 32.69231 41.68571</pre> <pre>tapply(mall\$Annual.Income.k..., clustergroups, mean)</pre> <pre>## 1 2 3 4 5 6 ## 26.30435 25.61905 55.60900 54.46899 86.53848 88.22857</pre> <pre>tapply(mall\$Spending.Score..1.100., clustergroups, mean)</pre> <pre>## 1 2 3 4 5 6 ## 29.92304 80.23810 48.77343 50.19149 82.12821 17.28571</pre>				
---	--	--	--	--

Cluster #1 characteristics

G1 = subset(mall, clustergroups==1)																																																				
<table> <tr> <th>Gender</th> <th>Age</th> <th>Annual.Income.k..</th> <th>Spending.Score..1.100.</th> </tr> <tr> <th><chr></th> <th><dbl></th> <th><dbl></th> <th><dbl></th> </tr> <tr> <td>1 Male</td> <td>19</td> <td>15</td> <td>39</td> </tr> <tr> <td>3 Female</td> <td>20</td> <td>16</td> <td>6</td> </tr> <tr> <td>5 Female</td> <td>31</td> <td>17</td> <td>40</td> </tr> <tr> <td>7 Female</td> <td>35</td> <td>18</td> <td>6</td> </tr> <tr> <td>9 Male</td> <td>64</td> <td>19</td> <td>3</td> </tr> <tr> <td>11 Male</td> <td>67</td> <td>19</td> <td>14</td> </tr> <tr> <td>13 Female</td> <td>58</td> <td>20</td> <td>15</td> </tr> <tr> <td>15 Male</td> <td>37</td> <td>20</td> <td>13</td> </tr> <tr> <td>17 Female</td> <td>35</td> <td>21</td> <td>35</td> </tr> <tr> <td>19 Male</td> <td>52</td> <td>23</td> <td>28</td> </tr> </table>					Gender	Age	Annual.Income.k..	Spending.Score..1.100.	<chr>	<dbl>	<dbl>	<dbl>	1 Male	19	15	39	3 Female	20	16	6	5 Female	31	17	40	7 Female	35	18	6	9 Male	64	19	3	11 Male	67	19	14	13 Female	58	20	15	15 Male	37	20	13	17 Female	35	21	35	19 Male	52	23	28
Gender	Age	Annual.Income.k..	Spending.Score..1.100.																																																	
<chr>	<dbl>	<dbl>	<dbl>																																																	
1 Male	19	15	39																																																	
3 Female	20	16	6																																																	
5 Female	31	17	40																																																	
7 Female	35	18	6																																																	
9 Male	64	19	3																																																	
11 Male	67	19	14																																																	
13 Female	58	20	15																																																	
15 Male	37	20	13																																																	
17 Female	35	21	35																																																	
19 Male	52	23	28																																																	
1-10 of 23 rows																																																				
<div> <div>Previous</div> <div>123</div> <div>Next</div> </div>																																																				
Show G1 summary																																																				
summary(G1)																																																				
<pre>## Gender Age Annual.Income.k... Spending.Score..1.100. ## Length:23 Min. :19.00 Min. :35.00 Min. : 3.00 ## Class :character 1st Qu.:35.50 1st Qu.:19.50 1st Qu.: 9.50 ## Mode :character Median :46.00 Median :25.00 Median :17.00 ## Mean :45.22 Mean :26.3 Mean :20.91 ## 3rd Qu.:53.50 3rd Qu.:33.0 3rd Qu.:33.50 ## Max. :67.00 Max. :39.0 Max. :40.00</pre>																																																				
table(G1\$Gender)																																																				
<pre>## ## Female Male ## 14 9</pre>																																																				

Spending score is low

Annual income is average

However, both spending score and annual income are small. The mean value of age in this cluster is 45, however this cluster have a various ages in it resulting in a high range of ages within the cluster.

Females is two times more than males

Cluster #2 characteristics

Spending score is low

Annual income is average

However, both spending score and annual income are small. The mean value of age in this cluster is 45, however this cluster have a various ages in it resulting in a high range of ages within the cluster.

Females is two times more than males

Cluster #2 characteristics

0	Female	23	18	94
1	Female	30	19	72
2	Female	35	19	99
3	Female	24	20	77
4	Male	22	20	79
5	Male	20	21	66
6	Female	35	23	98

1-10 of 21 rows

Previous123Next

summary(G2)

```
##      Gender      Age      Annual.Income.k... Spending.Score..1.100.  
## Length:21      Min.    :18.00   Min.    :15.00   Min.    :65.00  
## Class :character  1st Qu.:21.00   1st Qu.:19.00   1st Qu.:73.00  
## Mode :character  Median :23.00   Median :24.00   Median :77.00  
##      Mean    :24.81   Mean   :25.62   Mean   :80.24  
##      3rd Qu.:29.00   3rd Qu.:33.00   3rd Qu.:87.00  
##      Max.    :35.00   Max.    :39.00   Max.    :99.00
```

table(G2\$Gender)

```
##  
## Female Male  
##      13      8
```

Spending score is High

Annual Income is average

However, both spending score and annual income are small. The mean value of age in this cluster is between 25-35 age group. And females is more than 60% of the cluster

Cluster #3 characteristics

Spending score is High

Annual income is average

However, both spending score and annual income are small. The mean value of age in this cluster is between 25-35 age group. And females is more than 60% of the cluster

Cluster #3 characteristics

48	Female	27		40		47
49	Female	29		40		42
50	Female	31		40		42
52	Male	33		42		60
53	Female	31		43		55
59	Female	27		46		51
62	Male	19		46		55
66	Male	18		48		59
1-10 of 35 rows						Previous 1 2 3 4 Next
summary(G3)						
##	Gender	Age	Annual.Income.k..	Spending.Score..1.100.		
##	Length:35	Min. :18.00	Min. :35.00	Min. :29.00		
##	Class :character	1st Qu.:28.50	1st Qu.:47.0	1st Qu.:42.00		
##	Mode :character	Median :32.00	Median :60.0	Median :50.00		
##		Mean :25.97	Mean :54.47	Mean :48.77		
##		3rd Qu.:31.00	3rd Qu.:63.5	3rd Qu.:55.00		
##		Max. :40.00	Max. :76.0	Max. :61.00		
table(G3\$Gender)						
##						
##	Female	Male				
##	23	12				
Spending score is Average						
Annual Income is relatively high						

Spending score is Average

Annual income is relatively high

This cluster have a low spending score relative to its annual income. The mean value of age in this cluster is 25, however this cluster have a various ages in it resulting in a high range of ages within the cluster.

Females is two times more than males

Gender		Age	Annual.Income.k..	Spending.Score..1.100.
<chr>		<dbl>	<dbl>	<dbl>
47	Female	50	40	55
51	Female	49	42	52
54	Male	59	43	60
55	Female	50	43	45
56	Male	47	43	41
57	Female	51	44	50
58	Male	69	44	46
60	Male	53	46	46
61	Male	70	46	56
63	Female	67	47	52
1-10 of 47 rows				
Previous 1 2 3 4 5 Next				
summary(G4)				
<pre>## Gender Age Annual.Income.k... Spending.Score..1.100. ## Length:47 Min. :34.0 Min. :40.00 Min. :41.00 ## Class :character 1st Qu.:48.0 1st Qu.:48.00 1st Qu.:46.00 ## Mode :character Median :51.0 Median :54.00 Median :50.00 ## Mean :54.6 Mean :54.47 Mean :50.19 ## 3rd Qu.:64.0 3rd Qu.:62.00 3rd Qu.:55.00 ## Max. :78.0 Max. :69.00 Max. :60.00</pre>				
table(G4\$Gender)				
<pre>## ## Female Male ## 26 21</pre>				

Spending score is Average

Annual income is above average

This cluster have a low spending score relative to its annual income. The mean value of age in this cluster is 25, however this cluster have a various ages in it resulting in a high range of ages within the cluster.

Females is slightly higher than males

G5 = subset(mall, clustergroups==5)
G5

	Gender	Age	Annual.Income.k..	Spending.Score..1.100.
	<chr>	<dbl>	<dbl>	<dbl>
124	Male	39	69	91
126	Female	31	70	77
128	Male	40	71	95
130	Male	38	71	75
132	Male	39	71	75
134	Female	31	72	71
136	Female	29	73	88
138	Male	32	73	73
140	Female	35	74	72
142	Male	32	75	93

1-10 of 39 rows

Previous1234Next

summary(G5)

```
##      Gender      Age      Annual.Income.k... Spending.Score..1.100.  
## Length:39      Min.    :27.00   Min.    :69.00   Min.    :63.00  
## Class :character  1st Qu.:39.00   1st Qu.:75.50   1st Qu.:74.50  
## Mode :character  Median :32.00   Median :79.00   Median :83.00  
##      Mean    :32.69   Mean   :86.54   Mean   :82.13  
##      3rd Qu.:35.50   3rd Qu.:95.00   3rd Qu.:90.00  
##      Max.    :48.00   Max.    :137.00   Max.    :97.00
```

table(G5\$Gender)

```
##  
## Female Male  
##      11      18
```

Spending score is High

Annual income is high

This cluster have a high spending score as well as high annual income. The mean value of age in this cluster is 32, with no high variability in the age range

Females to males ration is negligible

Females to males ratio is negligible

G6 = subset(mall, clustergroups==6)
G6

	Gender	Age	Annual.Income.k..	Spending.Score..1.100.
	<chr>	<dbl>	<dbl>	<dbl>
127	Male	43	71	35
129	Male	59	71	11
131	Male	47	71	9
135	Male	20	73	5
137	Female	44	73	7
139	Male	19	74	10
141	Female	57	75	5
145	Male	25	77	12
147	Male	48	77	36
149	Female	34	78	22

1-10 of 35 rows

Previous1234Next

summary(G6)

```
##      Gender      Age      Annual.Income.k... Spending.Score..1.100.  
## Length:35      Min.    :19.00   Min.    :71.00   Min.    : 3.00  
## Class :character  1st Qu.:35.00   1st Qu.:77.50   1st Qu.:18.00  
## Mode :character  Median :43.00   Median :85.00   Median :16.00  
##      Mean    :41.69   Mean   :88.23   Mean   :19.52  
##      3rd Qu.:47.50   3rd Qu.:97.50   3rd Qu.:23.50  
##      Max.    :59.00   Max.    :137.00   Max.    :37.00
```

table(G6\$Gender)

```
##  
## Female Male  
##      15      20
```

Spending score is low

Annual income is low

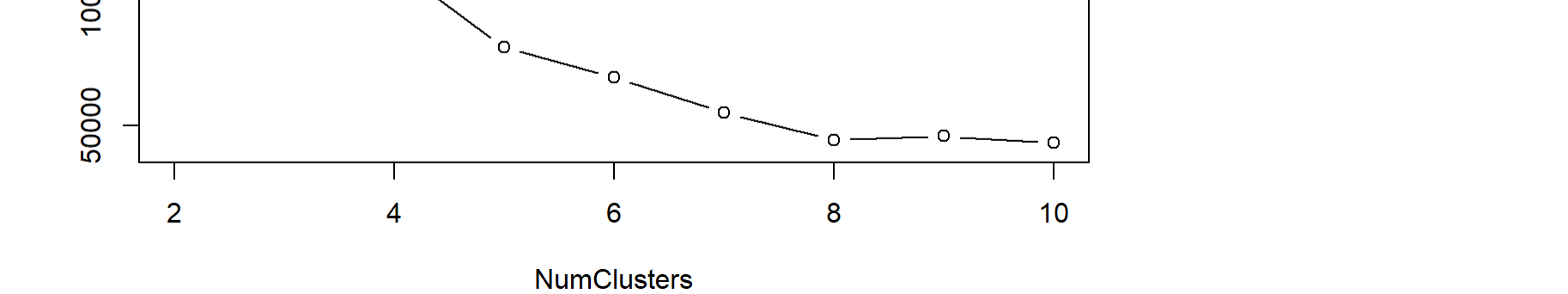
This cluster have a low spending score vs a annual income. The mean value of age in this cluster is 41, however this cluster have a various ages in it resulting in a high range of ages within the cluster.

This is the only cluster that has males higher than femals.

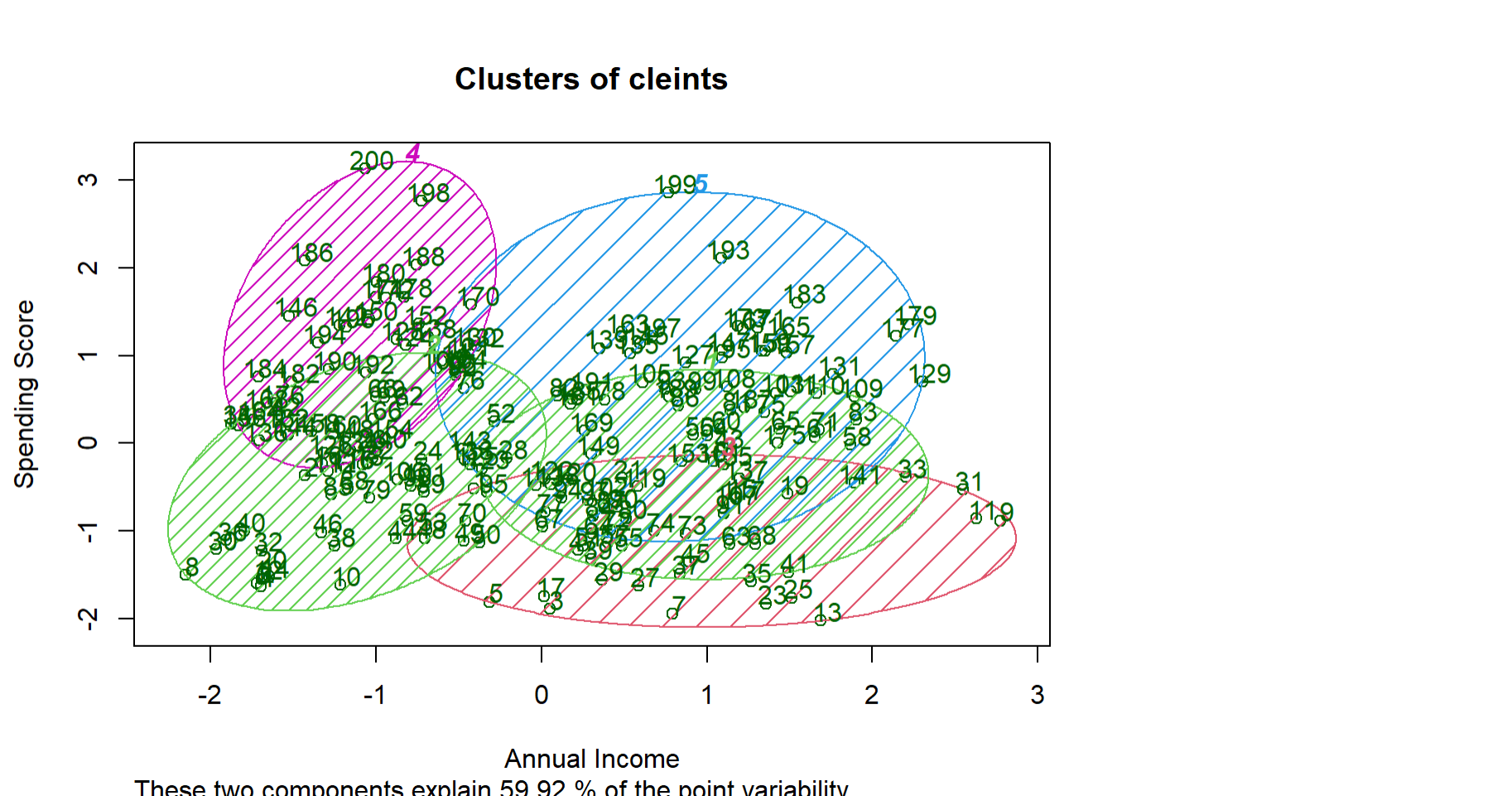
K-Means Clustering

<pre># Run k-means k<5 set.seed(1) KMC = kmeans(mall[,2:4], centers = k, iter.max = 1000) str(KMC)</pre>				
<pre>## List of 9 ## \$ cluster : Named int [1:208] 3 5 3 5 3 5 3 5 3 5 ... ## \$ attr(,"names")= chr [1:208] "1" "2" "3" "4" ... ## \$ centers : num [1:5, 1:3] 52.7 49.4 44.1 32.7 24.8 ... ## \$ attr(,"dimnames")=list of 2 ## .. \$: chr [1:5] "1" "2" "3" "4" ... ## .. \$: chr [1:3] "Age" "Annual.Income.k.." "Spending.Score..1.100." ## \$ totss : num 308813 ## \$ withinss : num [1:5] 11866 18994 7732 13972 27587 ## \$ betweenss : num 70961 ## \$ size : int [1:5] 52 38 21 39 50 ## \$ iter : int 3 ## \$ ifault : int 0 ## \$ attr(,"class")= chr "kmeans"</pre>				
<pre>mallClusters = KMC\$cluster #etermine best number of clusters KMC1 = kmeans(mall[,2:4], centers = 2, iter.max = 1000) KMC2 = kmeans(mall[,2:4], centers = 3, iter.max = 1000) KMC3 = kmeans(mall[,2:4], centers = 4, iter.max = 1000) KMC4 = kmeans(mall[,2:4], centers = 5, iter.max = 1000) KMC5 = kmeans(mall[,2:4], centers = 6, iter.max = 1000) KMC6 = kmeans(mall[,2:4], centers = 7, iter.max = 1000) KMC7 = kmeans(mall[,2:4], centers = 8, iter.max = 1000) KMC8 = kmeans(mall[,2:4], centers = 9, iter.max = 1000) KMC9 = kmeans(mall[,2:4], centers = 10, iter.max = 1000)</pre>				

elbow method



Clusters

<p>Clusters of clients</p>  <p>These two components explain 59.92 % of the point variability.</p>																																							
aggregate(mall, by=list(cluster=KMC\$cluster), mean)																																							
<pre>## Warning in mean.default(X[[i]], ...): argument is not numeric or logical: ## returning NA ## Warning in mean.default(X[[i]], ...): argument is not numeric or logical: ## returning NA ## Warning in mean.default(X[[i]], ...): argument is not numeric or logical: ## returning NA ## Warning in mean.default(X[[i]], ...): argument is not numeric or logical: ## returning NA ## Warning in mean.default(X[[i]], ...): argument is not numeric or logical: ## returning NA</pre>																																							
<table><thead><tr><th>cluster</th><th>Gender</th><th>Age</th><th>Annual.Income.k..</th><th>Spending.Score..1.100.</th></tr><tr><th><int></th><th><dbl></th><th><dbl></th><th><dbl></th><th><dbl></th></tr></thead><tbody><tr><td>1</td><td>NA</td><td>53.71154</td><td>54.42308</td><td>48.73077</td></tr><tr><td>2</td><td>NA</td><td>24.88000</td><td>41.46000</td><td>63.70000</td></tr><tr><td>3</td><td>NA</td><td>44.14286</td><td>25.14286</td><td>19.52861</td></tr><tr><td>4</td><td>NA</td><td>32.69231</td><td>86.53846</td><td>82.12821</td></tr><tr><td>5</td><td>NA</td><td>40.39474</td><td>87.00000</td><td>18.63158</td></tr></tbody></table>					cluster	Gender	Age	Annual.Income.k..	Spending.Score..1.100.	<int>	<dbl>	<dbl>	<dbl>	<dbl>	1	NA	53.71154	54.42308	48.73077	2	NA	24.88000	41.46000	63.70000	3	NA	44.14286	25.14286	19.52861	4	NA	32.69231	86.53846	82.12821	5	NA	40.39474	87.00000	18.63158
cluster	Gender	Age	Annual.Income.k..	Spending.Score..1.100.																																			
<int>	<dbl>	<dbl>	<dbl>	<dbl>																																			
1	NA	53.71154	54.42308	48.73077																																			
2	NA	24.88000	41.46000	63.70000																																			
3	NA	44.14286	25.14286	19.52861																																			
4	NA	32.69231	86.53846	82.12821																																			
5	NA	40.39474	87.00000	18.63158																																			
5 rows																																							
KMC\$size																																							
<pre>## [1] 52 38 21 39 38</pre>																																							