

```

1   Course: Exploratory_Data_Analysis
2   Lesson: K_Means_Clustering
3
4   - Class: text
5   Output: "K_Means_Clustering. (Slides for this and other Data Science courses may be
6   found at github https://github.com/DataScienceSpecialization/courses/. If you care to
7   use them, they must be downloaded as a zip file and viewed locally. This lesson
8   corresponds to 04_ExploratoryAnalysis/kmeansClustering.)"
9
10
11  - Class: text
12  Output: In this lesson we'll learn about k-means clustering, another simple way of
13  examining and organizing multi-dimensional data. As with hierarchical clustering,
14  this technique is most useful in the early stages of analysis when you're trying to
15  get an understanding of the data, e.g., finding some pattern or relationship between
16  different factors or variables.
17
18  - Class: text
19  Output: R documentation tells us that the k-means method "aims to partition the
20  points into k groups such that the sum of squares from points to the assigned cluster
21  centres is minimized."
22
23  - Class: text
24  Output: Since clustering organizes data points that are close into groups we'll
25  assume we've decided on a measure of distance, e.g., Euclidean.
26
27  - Class: figure
28  Output: To illustrate the method, we'll use these random points we generated,
29  familiar to you if you've already gone through the hierarchical clustering lesson.
30  We'll demonstrate k-means clustering in several steps, but first we'll explain the
31  general idea.
32  Figure: ranPoints.R
33  FigureType: new
34
35  - Class: text
36  Output: As we said, k-means is a partitioning approach which requires that you first
37  guess how many clusters you have (or want). Once you fix this number, you randomly
38  create a "centroid" (a phantom point) for each cluster and assign each point or
39  observation in your dataset to the centroid to which it is closest. Once each point
40  is assigned a centroid, you readjust the centroid's position by making it the average
41  of the points assigned to it.
42
43  - Class: text
44  Output: Once you have repositioned the centroids, you must recalculate the distance
45  of the observations to the centroids and reassign any, if necessary, to the centroid
46  closest to them. Again, once the reassignments are done, readjust the positions of
47  the centroids based on the new cluster membership. The process stops once you reach
48  an iteration in which no adjustments are made or when you've reached some
49  predetermined maximum number of iterations.
50
51  - Class: mult_question
52  Output: As described, what does this process require?
53  AnswerChoices: A defined distance metric; A number of clusters; An initial guess as
54  to cluster centroids; All of the others
55  CorrectAnswer: All of the others
56  AnswerTests: omnitest(correctVal='All of the others')
57  Hint: Which choice includes all the others.
58
59  - Class: mult_question
60  Output: So k-means clustering requires some distance metric (say Euclidean), a
61  hypothesized fixed number of clusters, and an initial guess as to cluster centroids.
62  As described, what does this process produce?
63  AnswerChoices: A final estimate of cluster centroids; An assignment of each point to
64  a cluster; All of the others
65  CorrectAnswer: All of the others
66  AnswerTests: omnitest(correctVal='All of the others')
67  Hint: Which choice includes all the others.
68
69  - Class: text

```

```

43 Output: When it's finished k-means clustering returns a final position of each
cluster's centroid as well as the assignment of each data point or observation to a
cluster.
44
45 - Class: text
46 Output: Now we'll step through this process using our random points as our data. The
coordinates of these are stored in 2 vectors, x and y. We eyeball the display and
guess that there are 3 clusters. We'll pick 3 positions of centroids, one for each
cluster.
47
48 - Class: cmd_question
49 Output: We've created two 3-long vectors for you, cx and cy. These respectively hold
the x- and y- coordinates for 3 proposed centroids. For convenience, we've also
stored them in a 2 by 3 matrix cmat. The x coordinates are in the first row and the y
coordinates in the second. Look at cmat now.
50 CorrectAnswer: cmat
51 AnswerTests: omnitest(correctExpr='cmat')
52 Hint: Type cmat at the command prompt.
53
54 - Class: cmd_question
55 Output: The coordinates of these points are (1,2), (1.8,1) and (2.5,1.5). We'll add
these centroids to the plot of our points. Do this by calling the R command points
with 6 arguments. The first 2 are cx and cy, and the third is col set equal to the
concatenation of 3 colors, "red", "orange", and "purple". The fourth argument is pch
set equal to 3 (a plus sign), the fifth is cex set equal to 2 (expansion of
character), and the final is lwd (line width) also set equal to 2.
56 CorrectAnswer: points(cx,cy,col=c("red","orange","purple"),pch=3,cex=2,lwd=2)
57 AnswerTests:
omnitest(correctExpr='points(cx,cy,col=c("red","orange","purple"),pch=3,cex=2,lwd=2)')
58 Hint: Type points(cx,cy,col=c("red","orange","purple"),pch=3,cex=2,lwd=2) at the
command prompt.
59
60 - Class: text
61 Output: We see the first centroid (1,2) is in red. The second (1.8,1), to the right
and below the first, is orange, and the final centroid (2.5,1.5), the furthest to the
right, is purple.
62
63 - Class: mult_question
64 Output: Now we have to calculate distances between each point and every centroid.
There are 12 data points and 3 centroids. How many distances do we have to calculate?
65 AnswerChoices: 15; 36; 9; 108
66 CorrectAnswer: 36
67 AnswerTests: omnitest(correctVal='36')
68 Hint: The distance between each point and one centroid means 12 distances have to be
calculated for each centroid. This has to be done for all 3 centroids.
69
70 - Class: cmd_question
71 Output: We've written a function for you called mdist which takes 4 arguments. The
vectors of data points (x and y) are the first two and the two vectors of centroid
coordinates (cx and cy) are the last two. Call mdist now with these arguments.
72 CorrectAnswer: mdist(x,y,cx,cy)
73 AnswerTests: omnitest(correctExpr='mdist(x,y,cx,cy)')
74 Hint: Type mdist(x,y,cx,cy) at the command prompt.
75
76
77 - Class: mult_question
78 Output: We've stored these distances in the matrix distTmp for you. Now we have to
assign a cluster to each point. To do that we'll look at each column and ?
79 AnswerChoices: pick the minimum entry; pick the maximum entry; add up the 3 entries.
80 CorrectAnswer: pick the minimum entry
81 AnswerTests: omnitest(correctVal='pick the minimum entry')
82 Hint: We assign each point to the centroid closest to it. Recall that the matrix
holds distances.
83
84 - Class: mult_question
85 Output: From the distTmp entries, which cluster would point 6 be assigned to?
86 AnswerChoices: 1; 2; 3; none of the above
87 CorrectAnswer: 3
88 AnswerTests: omnitest(correctVal='3')

```

```

89     Hint: Which row in column 6 has the lowest value?
90
91
92 - Class: cmd_question
93 Output: R has a handy function which.min which you can apply to ALL the columns of
distTmp with one call. Simply call the R function apply with 3 arguments. The first
is distTmp, the second is 2 meaning the columns of distTmp, and the third is
which.min, the function you want to apply to the columns of distTmp. Try this now.
94 CorrectAnswer: apply(distTmp,2,which.min)
95 AnswerTests: omnitest(correctExpr='apply(distTmp,2,which.min)')
96 Hint: Type apply(distTmp,2,which.min) at the command prompt.
97
98 - Class: text
99 Output: You can see that you were right and the 6th entry is indeed 3 as you answered
before. We see the first 3 entries were assigned to the second (orange) cluster and
only 2 points (4 and 8) were assigned to the first (red) cluster.
100
101 - Class: cmd_question
102 Output: We've stored the vector of cluster colors ("red","orange","purple") in the
array cols1 for you and we've also stored the cluster assignments in the array
newClust. Let's color the 12 data points according to their assignments. Again, use
the command points with 5 arguments. The first 2 are x and y. The third is pch set to
19, the fourth is cex set to 2, and the last, col is set to cols1[newClust].
103 CorrectAnswer: points(x,y,pch=19,cex=2,col=cols1[newClust])
104 AnswerTests: omnitest(correctExpr='points(x,y,pch=19,cex=2,col=cols1[newClust])')
105 Hint: Type points(x,y,pch=19,cex=2,col=cols1[newClust]) at the command prompt.
106
107 - Class: text
108 Output: Now we have to recalculate our centroids so they are the average (center of
gravity) of the cluster of points assigned to them. We have to do the x and y
coordinates separately. We'll do the x coordinate first. Recall that the vectors x
and y hold the respective coordinates of our 12 data points.
109
110 - Class: cmd_question
111 Output: We can use the R function tapply which applies "a function over a ragged
array". This means that every element of the array is assigned a factor and the
function is applied to subsets of the array (identified by the factor vector). This
allows us to take advantage of the factor vector newClust we calculated. Call tapply
now with 3 arguments, x (the data), newClust (the factor array), and mean (the
function to apply).
112 CorrectAnswer: tapply(x,newClust,mean)
113 AnswerTests: omnitest(correctExpr='tapply(x,newClust,mean)')
114 Hint: Type tapply(x,newClust,mean) at the command prompt.
115
116 - Class: cmd_question
117 Output: Repeat the call, except now apply it to the vector y instead of x.
118 CorrectAnswer: tapply(y,newClust,mean)
119 AnswerTests: omnitest(correctExpr='tapply(y,newClust,mean)')
120 Hint: Type tapply(y,newClust,mean) at the command prompt.
121
122
123 - Class: cmd_question
124 Output: Now that we have new x and new y coordinates for the 3 centroids we can plot
them. We've stored off the coordinates for you in variables newCx and newCy. Use the
R command points with these as the first 2 arguments. In addition, use the arguments
col set equal to cols1, pch equal to 8, cex equal to 2 and lwd also equal to 2.
125 CorrectAnswer: points(newCx,newCy,col=cols1,pch=8,cex=2,lwd=2)
126 AnswerTests: omnitest(correctExpr='points(newCx,newCy,col=cols1,pch=8,cex=2,lwd=2)')
127 Hint: Type points(newCx,newCy,col=cols1,pch=8,cex=2,lwd=2) at the command prompt.
128
129 - Class: cmd_question
130 Output: We see how the centroids have moved closer to their respective clusters. This
is especially true of the second (orange) cluster. Now call the distance function
mdist with the 4 arguments x, y, newCx, and newCy. This will allow us to reassign the
data points to new clusters if necessary.
131 CorrectAnswer: mdist(x,y,newCx,newCy)
132 AnswerTests: omnitest(correctExpr='mdist(x,y,newCx,newCy)')
133 Hint: Type mdist(x,y,newCx,newCy) at the command prompt.
134

```

```

135 - Class: mult_question
136 Output: We've stored off this new matrix of distances in the matrix distTmp2 for you. Recall that the first cluster is red, the second orange and the third purple. Look closely at columns 4 and 7 of distTmp2. What will happen to points 4 and 7?
137 AnswerChoices: Nothing; They will both change to cluster 2; They will both change clusters; They're the only points that won't change clusters
138 CorrectAnswer: They will both change clusters
139 AnswerTests: omnitest(correctVal='They will both change clusters')
140 Hint: Two of the choices are obviously wrong. That leaves two possibilities which are similar. Look carefully at the numbers in columns 4 and 7 to see where the minimum values are.

141
142 - Class: cmd_question
143 Output: Now call apply with 3 arguments, distTmp2, 2, and which.min to find the new cluster assignments for the points.
144 CorrectAnswer: apply(distTmp2,2,which.min)
145 AnswerTests: omnitest(correctExpr='apply(distTmp2,2,which.min)')
146 Hint: Type apply(distTmp2,2,which.min) at the command prompt.
147
148 - Class: cmd_question
149 Output: We've stored off the new cluster assignments in a vector of factors called newClust2. Use the R function points to recolor the points with their new assignments. Again, there are 5 arguments, x and y are first, followed by pch set to 19, cex to 2, and col to cols1[newClust2].
150 CorrectAnswer: points(x,y,pch=19,cex=2,col=cols1[newClust2])
151 AnswerTests: omnitest(correctExpr='points(x,y,pch=19,cex=2,col=cols1[newClust2])')
152 Hint: Type points(x,y,pch=19,cex=2,col=cols1[newClust2]) at the command prompt.
153
154 - Class: text
155 Output: Notice that points 4 and 7 both changed clusters, 4 moved from 1 to 2 (red to orange), and point 7 switched from 3 to 2 (purple to red).

156
157 - Class: cmd_question
158 Output: Now use tapply to find the x coordinate of the new centroid. Recall there are 3 arguments, x, newClust2, and mean.
159 CorrectAnswer: tapply(x,newClust2,mean)
160 AnswerTests: omnitest(correctExpr='tapply(x,newClust2,mean)')
161 Hint: Type tapply(x,newClust2,mean) at the command prompt.
162
163 - Class: cmd_question
164 Output: Do the same to find the new y coordinate.
165 CorrectAnswer: tapply(y,newClust2,mean)
166 AnswerTests: omnitest(correctExpr='tapply(y,newClust2,mean)')
167 Hint: Type tapply(y,newClust2,mean) at the command prompt.
168
169 - Class: cmd_question
170 Output: We've stored off these coordinates for you in the variables finalCx and finalCy. Plot these new centroids using the points function with 6 arguments. The first 2 are finalCx and finalCy. The argument col should equal cols1, pch should equal 9, cex 2 and lwd 2.
171 CorrectAnswer: points(finalCx,finalCy,col=cols1,pch=9,cex=2,lwd=2)
172 AnswerTests: omnitest(correctExpr='points(finalCx,finalCy,col=cols1,pch=9,cex=2,lwd=2)')
173 Hint: Type points(finalCx,finalCy,col=cols1,pch=9,cex=2,lwd=2) at the command prompt.
174
175 - Class: text
176 Output: It should be obvious that if we continued this process points 5 through 8 would all turn red, while points 1 through 4 stay orange, and points 9 through 12 purple.

177
178 - Class: text
179 Output: Now that you've gone through an example step by step, you'll be relieved to hear that R provides a command to do all this work for you. Unsurprisingly it's called kmeans and, although it has several parameters, we'll just mention four. These are x, (the numeric matrix of data), centers, iter.max, and nstart. The second of these (centers) can be either a number of clusters or a set of initial centroids. The third, iter.max, specifies the maximum number of iterations to go through, and nstart is the number of random starts you want to try if you specify centers as a number.

```

```

181 - Class: cmd_question
182 Output: Call kmeans now with 2 arguments, dataFrame (which holds the x and y
coordinates of our 12 points) and centers set equal to 3.
183 CorrectAnswer: kmeans(dataFrame,centers=3)
184 AnswerTests: omnitest(correctExpr='kmeans(dataFrame,centers=3)')
185 Hint: Type kmeans(dataFrame,centers=3) at the command prompt.
186
187 - Class: cmd_question
188 Output: The program returns the information that the data clustered into 3 clusters
each of size 4. It also returns the coordinates of the 3 cluster means, a vector
named cluster indicating how the 12 points were partitioned into the clusters, and
the sum of squares within each cluster. It also shows all the available components
returned by the function. We've stored off this data for you in a kmeans object
called kmObj. Look at kmObj$iter to see how many iterations the algorithm went through.
189 CorrectAnswer: kmObj$iter
190 AnswerTests: omnitest(correctExpr='kmObj$iter')
191 Hint: Type kmObj$iter at the command prompt.
192
193 - Class: cmd_question
194 Output: Two iterations as we did before. We just want to emphasize how you can access
the information available to you. Let's plot the data points color coded according to
their cluster. This was stored in kmObj$cluster. Run plot with 5 arguments. The data,
x and y, are the first two; the third, col is set equal to kmObj$cluster, and the
last two are pch and cex. The first of these should be set to 19 and the last to 2.
195 CorrectAnswer: plot(x,y,col=kmObj$cluster,pch=19,cex=2)
196 AnswerTests: omnitest(correctExpr='plot(x,y,col=kmObj$cluster,pch=19,cex=2)')
197 Hint: Type plot(x,y,col=kmObj$cluster,pch=19,cex=2) at the command prompt.
198
199 - Class: cmd_question
200 Output: Now add the centroids which are stored in kmObj$centers. Use the points
function with 5 arguments. The first two are kmObj$centers and
col=c("black","red","green"). The last three, pch, cex, and lwd, should all equal 3.
201 CorrectAnswer: points(kmObj$centers,col=c("black","red","green"),pch=3,cex=3,lwd=3)
202 AnswerTests:
omnitest(correctExpr='points(kmObj$centers,col=c("black","red","green"),pch=3,cex=3,lwd
=3)')
203 Hint: Type points(kmObj$centers,col=c("black","red","green"),pch=3,cex=3,lwd=3) at
the command prompt.
204
205
206 - Class: text
207 Output: Now for some fun! We want to show you how the output of the kmeans function
is affected by its random start (when you just ask for a number of clusters). With
random starts you might want to run the function several times to get an idea of the
relationships between your observations. We'll call kmeans with the same data points
(stored in dataFrame), but ask for 6 clusters instead of 3.
208
209
210 - Class: cmd_question
211 Output: We'll plot our data points several times and each time we'll just change the
argument col which will show us how the R function kmeans is clustering them. So,
call plot now with 5 arguments. The first 2 are x and y. The third is col set equal
to the call kmeans(dataFrame,6)$cluster. The last two (pch and cex) are set to 19 and
2 respectively.
212 CorrectAnswer: plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2)
213 AnswerTests:
omnitest(correctExpr='plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2)')
214 Hint: Type plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2) at the command
prompt.
215
216 - Class: cmd_question
217 Output: See how the points cluster? Now recall your last command and rerun it.
218 CorrectAnswer: plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2)
219 AnswerTests:
omnitest(correctExpr='plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2)')
220 Hint: Type plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2) at the command
prompt.
221
222 - Class: cmd_question

```

```

223 Output: See how the clustering has changed? As the Teletubbies would say, "Again!
Again!"
224 CorrectAnswer: plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2)
225 AnswerTests:
omnitest(correctExpr='plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2)')
226 Hint: Type plot(x,y,col=kmeans(dataFrame,6)$cluster,pch=19,cex=2) at the command
prompt.
227
228 - Class: text
229 Output: So the clustering changes with different starts. Perhaps 6 is too many
clusters? Let's review!
230
231 - Class: mult_question
232 Output: True or False? K-means clustering requires you to specify a number of
clusters before you begin.
233 AnswerChoices: True; False
234 CorrectAnswer: True
235 AnswerTests: omnitest(correctVal='True')
236 Hint: What did you provide when you called the R function?
237
238 - Class: mult_question
239 Output: True or False? K-means clustering requires you to specify a number of
iterations before you begin.
240 AnswerChoices: True; False
241 CorrectAnswer: False
242 AnswerTests: omnitest(correctVal='False')
243 Hint: What did you provide when you called the R function?
244
245 - Class: mult_question
246 Output: True or False? Every data set has a single fixed number of clusters.
247 AnswerChoices: True; False
248 CorrectAnswer: False
249 AnswerTests: omnitest(correctVal='False')
250 Hint: The number of clusters depends on your eye.
251
252 - Class: mult_question
253 Output: True or False? K-means clustering will always stop in 3 iterations
254 AnswerChoices: True; False
255 CorrectAnswer: False
256 AnswerTests: omnitest(correctVal='False')
257 Hint: The number of iterations depends on your data.
258
259
260 - Class: mult_question
261 Output: True or False? When starting kmeans with random centroids, you'll always end
up with the same final clustering.
262 AnswerChoices: True; False
263 CorrectAnswer: False
264 AnswerTests: omnitest(correctVal='False')
265 Hint: Recall the last experiment we did in the lesson, rerunning the same routine.
266
267 - Class: text
268 Output: Congratulations! We hope this means you found this lesson oK.
269
270 - Class: mult_question
271 Output: "Would you like to receive credit for completing this course on
Coursera.org?"
272 CorrectAnswer: NULL
273 AnswerChoices: Yes;No
274 AnswerTests: coursera_on_demand()
275 Hint: ""
276
277

```