

```

1  Course: Exploratory_Data_Analysis
2  Lesson: Hierarchical_Clustering
3
4  - Class: text
5  Output: "Hierarchical_Clustering. (Slides for this and other Data Science courses may
be found at github https://github.com/DataScienceSpecialization/courses/. If you care
to use them, they must be downloaded as a zip file and viewed locally. This lesson
corresponds to 04_ExploratoryAnalysis/hierarchicalClustering.)"
6
7
8  - Class: text
9  Output: In this lesson we'll learn about hierarchical clustering, a simple way of
quickly examining and displaying multi-dimensional data. This technique is usually
most useful in the early stages of analysis when you're trying to get an
understanding of the data, e.g., finding some pattern or relationship between
different factors or variables. As the name suggests hierarchical clustering creates
a hierarchy of clusters.
10
11 - Class: text
12 Output: Clustering organizes data points that are close into groups. So obvious
questions are "How do we define close?", "How do we group things?", and "How do we
interpret the grouping?" Cluster analysis is a very important topic in data analysis.
13
14 - Class: figure
15 Output: To give you an idea of what we're talking about, consider these random points
we generated. We'll use them to demonstrate hierarchical clustering in this lesson.
We'll do this in several steps, but first we have to clarify our terms and concepts.
16 Figure: ranPoints.R
17 FigureType: new
18
19 - Class: text
20 Output: Hierarchical clustering is an agglomerative, or bottom-up, approach. From
Wikipedia (http://en.wikipedia.org/wiki/Hierarchical\_clustering), we learn that in
this method, "each observation starts in its own cluster, and pairs of clusters are
merged as one moves up the hierarchy." This means that we'll find the closest two
points and put them together in one cluster, then find the next closest pair in the
updated picture, and so forth. We'll repeat this process until we reach a reasonable
stopping place.
21
22 - Class: text
23 Output: Note the word "reasonable". There's a lot of flexibility in this field and
how you perform your analysis depends on your problem. Again, Wikipedia tells us,
"one can decide to stop clustering either when the clusters are too far apart to be
merged (distance criterion) or when there is a sufficiently small number of clusters
(number criterion).".
24
25 - Class: text
26 Output: First, how do we define close? This is the most important step and there are
several possibilities depending on the questions you're trying to answer and the data
you have. Distance or similarity are usually the metrics used.
27
28 - Class: mult_question
29 Output: In the given plot which pair points would you first cluster? Use distance as
the metric.
30 AnswerChoices: 7 and 8; 1 and 4; 5 and 6; 10 and 12
31 CorrectAnswer: 5 and 6
32 AnswerTests: omnitest(correctVal='5 and 6')
33 Hint: The choices aren't very close. Look at the picture for the answer.
34
35 - Class: text
36 Output: It's pretty obvious that out of the 4 choices, the pair 5 and 6 were the
closest together. However, there are several ways to measure distance or similarity.
Euclidean distance and correlation similarity are continuous measures, while
Manhattan distance is a binary measure. In this lesson we'll just briefly discuss the
first and last of these. It's important that you use a measure of distance that fits
your problem.
37
38 - Class: figure
39 Output: Euclidean distance is what you learned about in high school algebra. Given

```

two points on a plane, (x_1, y_1) and (x_2, y_2) , the Euclidean distance is the square root of the sums of the squares of the distances between the two x-coordinates $(x_1 - x_2)$ and the two y-coordinates $(y_1 - y_2)$. You probably recognize this as an application of the Pythagorean theorem which yields the length of the hypotenuse of a right triangle.

Figure: showEuclid.R

FigureType: new

- **Class:** text

Output: It shouldn't be hard to believe that this generalizes to more than two dimensions as shown in the formula at the bottom of the picture shown here.

- **Class:** text

Output: Euclidean distance is distance "as the crow flies". Many applications, however, can't realistically use crow-flying distance. Cars, for instance, have to follow roads.

- **Class:** figure

Output: In this case, we can use Manhattan or city block distance (also known as a taxicab metric). This picture, copied from http://en.wikipedia.org/wiki/Taxicab_geometry, shows what this means.

Figure: showTaxi.R

FigureType: new

- **Class:** text

Output: You want to travel from the point at the lower left to the one on the top right. The shortest distance is the Euclidean (the green line), but you're limited to the grid, so you have to follow a path similar to those shown in red, blue, or yellow. These all have the same length (12) which is the number of small gray segments covered by their paths.

- **Class:** text

Output: More formally, Manhattan distance is the sum of the absolute values of the distances between each coordinate, so the distance between the points (x_1, y_1) and (x_2, y_2) is $|x_1 - x_2| + |y_1 - y_2|$. As with Euclidean distance, this too generalizes to more than 2 dimensions.

- **Class:** figure

Output: Now we'll go back to our random points. You might have noticed that these points don't really look randomly positioned, and in fact, they're not. They were actually generated as 3 distinct clusters. We've put the coordinates of these points in a data frame for you, called `dataFrame`.

Figure: ranPoints.R

FigureType: new

- **Class:** text

Output: We'll use this `dataFrame` to demonstrate an agglomerative (bottom-up) technique of hierarchical clustering and create a dendrogram. This is an abstract picture (or graph) which shows how the 12 points in our dataset cluster together. Two clusters (initially, these are points) that are close are connected with a line, We'll use Euclidean distance as our metric of closeness.

- **Class:** cmd_question

Output: Run the R command `dist` with the argument `dataFrame` to compute the distances between all pairs of these points. By default `dist` uses Euclidean distance as its metric, but other metrics such as Manhattan, are available. Just use the default.

CorrectAnswer: `dist(dataFrame)`

AnswerTests: `omnittest(correctExpr='dist(dataFrame)')`

Hint: Type `dist(dataFrame)` at the command prompt.

- **Class:** text

Output: You see that the output is a lower triangular matrix with rows numbered from 2 to 12 and columns numbered from 1 to 11. Entry (i, j) indicates the distance between points i and j . Clearly you need only a lower triangular matrix since the distance between points i and j equals that between j and i .

- **Class:** mult_question

Output: From the output of `dist`, what is the minimum distance between two points?

AnswerChoices: 0.0815; 0.08317; -0.0700; 0.1085

CorrectAnswer: 0.0815

```

81  AnswerTests: omnitest(correctVal='0.0815')
82  Hint: Recall a previous question where points 5 and 6 looked close.
83
84  - Class: figure
85  Output: So 0.0815 (units are unspecified) between points 5 and 6 is the shortest
      distance. We can put these points in a single cluster and look for another close pair
      of points.
86  Figure: cluster56.R
87  FigureType: new
88
89  - Class: mult_question
90  Output: Looking at the picture, what would be another good pair of points to put in
      another cluster given that 5 and 6 are already clustered?
91  AnswerChoices: 7 and the cluster containing 5 and 6; 10 and 11; 1 and 4; 7 and 8
92  CorrectAnswer: 10 and 11
93  AnswerTests: omnitest(correctVal='10 and 11')
94  Hint: Which of the choices looks closest on the picture?
95
96  - Class: figure
97  Output: So 10 and 11 are another pair of points that would be in a second cluster.
      We'll start creating our dendrogram now. Here're the original plot and two beginning
      pieces of the dendrogram.
98  Figure: startDendro.R
99  FigureType: new
100
101  - Class: cmd_question
102  Output: We can keep going like this in the obvious way and pair up individual
      points, but as luck would have it, R provides a simple function which you can call
      which creates a dendrogram for you. It's called hclust() and takes as an argument the
      pairwise distance matrix which we looked at before. We've stored this matrix for you
      in a variable called distxy. Run hclust now with distxy as its argument and put the
      result in the variable hc.
103  CorrectAnswer: hc <- hclust(distxy)
104  AnswerTests: expr_creates_var("hc"); omnitest(correctExpr='hc <- hclust(distxy)')
105  Hint: Type hc <- hclust(distxy) at the command prompt.
106
107  - Class: figure
108  Output: You're probably curious and want to see hc.
109  Figure: clearPlot.R
110  FigureType: new
111
112  - Class: cmd_question
113  Output: Call the R function plot with one argument, hc.
114  CorrectAnswer: plot(hc)
115  AnswerTests: omnitest(correctExpr='plot(hc)')
116  Hint: Type plot(hc) at the command prompt.
117
118  - Class: cmd_question
119  Output: Nice plot, right? R's plot conveniently labeled everything for you. The
      points we saw are the leaves at the bottom of the graph, 5 and 6 are connected, as
      are 10 and 11. Moreover, we see that the original 3 groupings of points are closest
      together as leaves on the picture. That's reassuring. Now call plot again, this time
      with the argument as.dendrogram(hc).
120  CorrectAnswer: plot(as.dendrogram(hc))
121  AnswerTests: omnitest(correctExpr='plot(as.dendrogram(hc))')
122  Hint: Type plot(as.dendrogram(hc)) at the command prompt.
123
124  - Class: cmd_question
125  Output: The essentials are the same, but the labels are missing and the leaves
      (original points) are all printed at the same level. Notice that the vertical heights
      of the lines and labeling of the scale on the left edge give some indication of
      distance. Use the R command abline to draw a horizontal blue line at 1.5 on this
      plot. Recall that this requires 2 arguments, h=1.5 and col="blue".
126  CorrectAnswer: abline(h=1.5,col="blue")
127  AnswerTests: omnitest(correctExpr='abline(h=1.5,col="blue")')
128  Hint: Type abline(h=1.5,col="blue") at the command prompt.
129
130  - Class: cmd_question
131  Output: We see that this blue line intersects 3 vertical lines and this tells us that

```

using the distance 1.5 (unspecified units) gives us 3 clusters (1 through 4), (9 through 12), and (5 through 8). We call this a "cut" of our dendrogram. Now cut the dendrogram by drawing a red horizontal line at .4.

CorrectAnswer: `abline(h=.4,col="red")`

AnswerTests: `omnittest(correctExpr='abline(h=.4,col="red")')`

Hint: Type `abline(h=.4,col="red")` at the command prompt.

- **Class:** `cmd_question`

Output: How many clusters are there with a cut at this distance?

CorrectAnswer: 5

AnswerTests: `equiv_val(5)`

Hint: How many vertical lines does this red line cross?

- **Class:** `cmd_question`

Output: We see that by cutting at .4 we have 5 clusters, indicating that this distance is small enough to break up our original grouping of points. If we drew a horizontal line at .05, how many clusters would we get

CorrectAnswer: 12

AnswerTests: `equiv_val(12)`

Hint: Recall that our shortest distance was around .08, so a distance smaller than that would make all the points their own private clusters.

- **Class:** `cmd_question`

Output: Try it now (draw a horizontal line at .05) and make the line green.

CorrectAnswer: `abline(h=.05,col="green")`

AnswerTests: `abline(h=.05,col="green")`

Hint: Type `abline(h=.05,col="green")` at the command prompt.

- **Class:** `text`

Output: So the number of clusters in your data depends on where you draw the line! (We said there's a lot of flexibility here.) Now that we've seen the practice, let's go back to some "theory". Notice that the two original groupings, 5 through 8, and 9 through 12, are connected with a horizontal line near the top of the display. You're probably wondering how distances between clusters of points are measured.

- **Class:** `text`

Output: There are several ways to do this. We'll just mention two. The first is called complete linkage and it says that if you're trying to measure a distance between two clusters, take the greatest distance between the pairs of points in those two clusters. Obviously such pairs contain one point from each cluster.

- **Class:** `figure`

Output: So if we were measuring the distance between the two clusters of points (1 through 4) and (5 through 8), using complete linkage as the metric we would use the distance between points 4 and 8 as the measure since this is the largest distance between the pairs of those groups.

Figure: `complete.R`

FigureType: `new`

- **Class:** `figure`

Output: The distance between the two clusters of points (9 through 12) and (5 through 8), using complete linkage as the metric, is the distance between points 11 and 8 since this is the largest distance between the pairs of those groups.

Figure: `complete2.R`

FigureType: `new`

- **Class:** `figure`

Output: As luck would have it, the distance between the two clusters of points (9 through 12) and (1 through 4), using complete linkage as the metric, is the distance between points 11 and 4.

Figure: `complete3.R`

FigureType: `new`

- **Class:** `cmd_question`

Output: We've created the dataframe `dFsm` for you containing these 3 points, 4, 8, and 11. Run `dist` on `dFsm` to see what the smallest distance between these 3 points is.

CorrectAnswer: `dist(dFsm)`

AnswerTests: `omnittest(correctExpr='dist(dFsm)')`

Hint: Type `dist(dFsm)` at the command prompt.

```

180
181 - Class: figure
182 Output: We see that the smallest distance is between points 2 and 3 in this reduced
set, (these are actually points 8 and 11 in the original set), indicating that the
two clusters these points represent ((5 through 8) and (9 through 12) respectively)
would be joined (at a distance of 1.869) before being connected with the third
cluster (1 through 4). This is consistent with the dendrogram we plotted.
183 Figure: dendro.R
184 FigureType: new
185
186 - Class: figure
187 Output: The second way to measure a distance between two clusters that we'll just
mention is called average linkage. First you compute an "average" point in each
cluster (think of it as the cluster's center of gravity). You do this by computing
the mean (average) x and y coordinates of the points in the cluster.
188 Figure: average.R
189 FigureType: new
190
191 - Class: figure
192 Output: Then you compute the distances between each cluster average to compute the
intercluster distance.
193 Figure: average2.R
194 FigureType: new
195
196 - Class: cmd_question
197 Output: Now look at the hierarchical cluster we created before, hc.
198 CorrectAnswer: hc
199 AnswerTests: omnitest(correctExpr='hc')
200 Hint: Type hc at the command prompt.
201
202 - Class: mult_question
203 Output: Which type of linkage did hclust() use to agglomerate clusters?
204 AnswerChoices: average; complete
205 CorrectAnswer: complete
206 AnswerTests: omnitest(correctVal='complete')
207 Hint: Look at the output when you looked at hc. What was the cluster method?
208
209 - Class: text
210 Output: In our simple set of data, the average and complete linkages aren't that
different, but in more complicated datasets the type of linkage you use could affect
how your data clusters. It is a good idea to experiment with different methods of
linkage to see the varying ways your data groups. This will help you determine the
best way to continue with your analysis.
211
212 - Class: text
213 Output: The last method of visualizing data we'll mention in this lesson concerns
heat maps. Wikipedia (http://en.wikipedia.org/wiki/Heat\_map) tells us a heat map is
"a graphical representation of data where the individual values contained in a matrix
are represented as colors. ... Heat maps originated in 2D displays of the values in a
data matrix. Larger values were represented by small dark gray or black squares
(pixels) and smaller values by lighter squares."
214
215 - Class: text
216 Output: You've probably seen many examples of heat maps, for instance weather radar
and displays of ocean salinity. From Wikipedia
(http://en.wikipedia.org/wiki/Heat\_map) we learn that heat maps are often used in
molecular biology "to represent the level of expression of many genes across a number
of comparable samples (e.g. cells in different states, samples from different
patients) as they are obtained from DNA microarrays."
217
218 - Class: figure
219 Output: We won't say too much on this topic, but a very nice concise tutorial on
creating heatmaps in R exists at
http://sebastianraschka.com/Articles/heatmaps\_in\_r.html#clustering. Here's an image
from the tutorial to start you thinking about the topic. It shows a sample heat map
with a dendrogram on the left edge mapping the relationship between the rows. The
legend at the top shows how colors relate to values.
220 Figure: showheat.R
221 FigureType: new

```

```

222
223
224 - Class: cmd_question
225 Output: R provides a handy function to produce heat maps. It's called heatmap. We've
put the point data we've been using throughout this lesson in a matrix. Call heatmap
now with 2 arguments. The first is dataMatrix and the second is col set equal to
cm.colors(25). This last is optional, but we like the colors better than the default
ones.
226 CorrectAnswer: heatmap(dataMatrix,col=cm.colors(25))
227 AnswerTests: omnitest(correctExpr='heatmap(dataMatrix,col=cm.colors(25))')
228 Hint: Type heatmap(dataMatrix,col=cm.colors(25)) at the command prompt.
229
230 - Class: text
231 Output: We see an interesting display of sorts. This is a very simple heat map -
simple because the data isn't very complex. The rows and columns are grouped together
as shown by colors. The top rows (labeled 5, 6, and 7) seem to be in the same group
(same colors) while 8 is next to them but colored differently. This matches the
dendrogram shown on the left edge. Similarly, 9, 12, 11, and 10 are grouped together
(row-wise) along with 3 and 2. These are followed by 1 and 4 which are in a separate
group. Column data is treated independently of rows but is also grouped.
232
233 - Class: cmd_question
234 Output: We've subsetting some vehicle data from mtcars, the Motor Trend Car Road
Tests which is part of the package datasets. The data is in the matrix mt and
contains 6 factors of 11 cars. Run heatmap now with mt as its only argument.
235 CorrectAnswer: heatmap(mt)
236 AnswerTests: omnitest(correctExpr='heatmap(mt)')
237 Hint: Type heatmap(mt) at the command prompt.
238
239 - Class: cmd_question
240 Output: This looks slightly more interesting than the heatmap for the point data. It
shows a little better how the rows and columns are treated (clustered and colored)
independently of one another. To understand the disparity in color (between the left
4 columns and the right 2) look at mt now.
241 CorrectAnswer: mt
242 AnswerTests: omnitest(correctExpr='mt')
243 Hint: Type mt at the command prompt.
244
245 - Class: cmd_question
246 Output: See how four of the columns are all relatively small numbers and only two
(displacement and horsepower) are large? That explains the big difference in color columns. Now to
understand the grouping of the rows call plot with one argument, the dendrogram
object denmt we've created for you.
247 CorrectAnswer: plot(denmt)
248 AnswerTests: omnitest(correctExpr='plot(denmt)')
249 Hint: Type plot(denmt) at the command prompt.
250
251 - Class: cmd_question
252 Output: We see that this dendrogram is the one displayed at the side of the heat
map. How was this created? Recall that we generalized the distance formula for more
than 2 dimensions. We've created a distance matrix for you, distmt. Look at it now.
253 CorrectAnswer: distmt
254 AnswerTests: omnitest(correctExpr='distmt')
255 Hint: Type distmt at the command prompt.
256
257 - Class: text
258 Output: See how these distances match those in the dendrogram? So hclust really
works! Let's review now.
259
260 - Class: mult_question
261 Output: What is the purpose of hierarchical clustering?
262 AnswerChoices: Give an idea of the relationships between variables or observations;
Present a finished picture; Inspire other researchers; None of the others
263 CorrectAnswer: Give an idea of the relationships between variables or observations
264 AnswerTests: omnitest(correctVal='Give an idea of the relationships between variables
or observations')
265 Hint: Recall that this is a technique for EXPLORING data when you're first starting
to research a problem.
266

```



```

267 - Class: mult_question
268 Output: True or False? When you're doing hierarchical clustering there are strict
rules that you MUST follow.
269 AnswerChoices: True; False
270 CorrectAnswer: False
271 AnswerTests: omnitest(correctVal='False')
272 Hint: There are always choices to be made in this area.
273
274 - Class: mult_question
275 Output: True or False? There's only one way to measure distance.
276 AnswerChoices: True; False
277 CorrectAnswer: False
278 AnswerTests: omnitest(correctVal='False')
279 Hint: There are always choices to be made in this area.
280
281 - Class: mult_question
282 Output: True or False? Complete linkage is a method of computing distances between
clusters.
283 AnswerChoices: True; False
284 CorrectAnswer: True
285 AnswerTests: omnitest(correctVal='True')
286 Hint: Once a cluster contains more than one point you need a method of defining a
distance between it and other clusters.
287
288 - Class: mult_question
289 Output: True or False? Average linkage uses the maximum distance between points of
two clusters as the distance between those clusters.
290 AnswerChoices: True; False
291 CorrectAnswer: False
292 AnswerTests: omnitest(correctVal='False')
293 Hint: Average linkage uses the average or mean point as the representative of a
cluster. The distance between these average points is the distance between the
clusters.
294
295 - Class: mult_question
296 Output: True or False? The number of clusters you derive from your data depends on
the distance at which you choose to cut it.
297 AnswerChoices: True; False
298 CorrectAnswer: True
299 AnswerTests: omnitest(correctVal='True')
300 Hint: Recall our example where we drew horizontal cut lines through the dendrogram.
301
302 - Class: mult_question
303 Output: True or False? Once you decide basics, such as defining a distance metric
and linkage method, hierarchical clustering is deterministic.
304 AnswerChoices: True; False
305 CorrectAnswer: True
306 AnswerTests: omnitest(correctVal='True')
307 Hint: Once you pick your algorithm, all you have to do is apply it.
308
309 - Class: text
310 Output: Congratulations! We hope this lesson didn't fluster you or get you too heated!
311
312 - Class: mult_question
313 Output: "Would you like to receive credit for completing this course on
Coursera.org?"
314 CorrectAnswer: NULL
315 AnswerChoices: Yes;No
316 AnswerTests: coursera_on_demand()
317 Hint: ""
318
319

```