

# 1 Part 9: Clustering

## 2 Summary

3 Here, we will focus on the general problem of **clustering**, the challenge of  
4 dividing objects into  $K$  groups such that similar objects are assigned to the  
5 same group.

## 1 Why Cluster?

2 Making predictions is a very natural problem; most people can understand  
3 why it is useful to study methods for supervised learning. Applications for  
4 unsupervised learning can be less obvious.

5 There can be distinct groups present in data sets. This grouping can pro-  
6 vide insight into the structure of the data, and guide further analyses (in  
7 the same way that stocks in different “sectors” may be best modeled sepa-  
8 rately).

9 Clustering can also help in the identification of outliers.

10 The following examples are from finance. Note that I am not making any  
11 claims as to the legitimacy and/or accuracy of their analyses; I am just  
12 providing examples of ways in which clustering techniques could be used.

1 **Pavlidis, et al. (2006):** “Financial Forecasting through Unsupervised Clus-  
2 tering and Neural Networks,” *Operational Research*.

3 The authors use clustering to divide the predictor space into regions of  
4 greater homogeneity, and fit distinct (neural network) models in each of  
5 these.

6 “Although global approximation methods can be applied to model and  
7 forecast time series . . . , it is reasonable to expect that forecasting accuracy  
8 can be improved if regions of the input space exhibiting similar dynamics  
9 are identified and subsequently a local model is constructed for each of  
10 them.”

- 1 **Lee, et al. (2010):** “An Effective Clustering Approach to Stock Market Pre-  
2 diction,” *PACIS 2010 Proceedings*
- 3 The authors use clustering to divide the financial reports into homoge-  
4 neous groups, and then take the cluster center as a “representative.”
- 5 “When a financial report is released, we will transform it into a feature  
6 vector . . . [and] we assign [it] to the nearest representative feature vector.  
7 Then, we predict the direction of the stock price movement according to  
8 the class label of the nearest representative feature vector.”
- 9 In this framework, the clustering serves as a “smoothing” operation on  
10 “nonstandard” data (the financial report).

Financial report		$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	$d_7$	$d_8$	$d_9$
Qualitative features	efficient	1	1	0	1	0	0	1	0	0
	growth	1	1	0	0	0	0	0	0	0
	advantage	1	1	1	0	0	0	0	0	0
	improvement	1	0	0	0	0	0	0	0	0
	deficient	0	0	0	1	1	0	0	0	0
	reorganize	0	0	0	1	1	0	0	1	1
	difficulty	0	0	0	1	1	1	1	0	0
Quantitative features	complaint	0	0	0	1	0	0	0	0	1
	operating margin	0.4	0.38	0.37	0.1	0.07	0.08	0.05	0.06	0.03
	ROE	0.3	0.28	0.27	0.01	0.04	0.04	0.07	0.02	0.05
	ROTA	0.25	0.23	0.22	0.02	0.05	0.07	0.04	0.01	0.04
	equity to capital	0.8	0.78	0.77	0.45	0.4	0.5	0.45	0.5	0.55
Class label		1	1	0	-1	-1	-1	0	-1	-1

Table 1. An example dataset.

Figure 1: From Lee, et al. (2010)

- 1 **Miceli and Susinno (2003): “Using Trees to Grow Money,” *Risk***
- 2 The authors cluster hedge fund performance time series in order to “visu-
- 3 alize a taxonomy embedded in synchronous historical data.” They state
- 4 that “[i]n a universe with low transparency from an investor’s point of
- 5 view, and where operating strategies are protected from full disclosure . . .
- 6 instruments of classification . . . would allow the investor to verify the de-
- 7 clared strategies in the statements issued by hedge funds. In particular,
- 8 nodes anomolous to a reference cluster help to identify funds that require
- 9 a deeper investigation.”

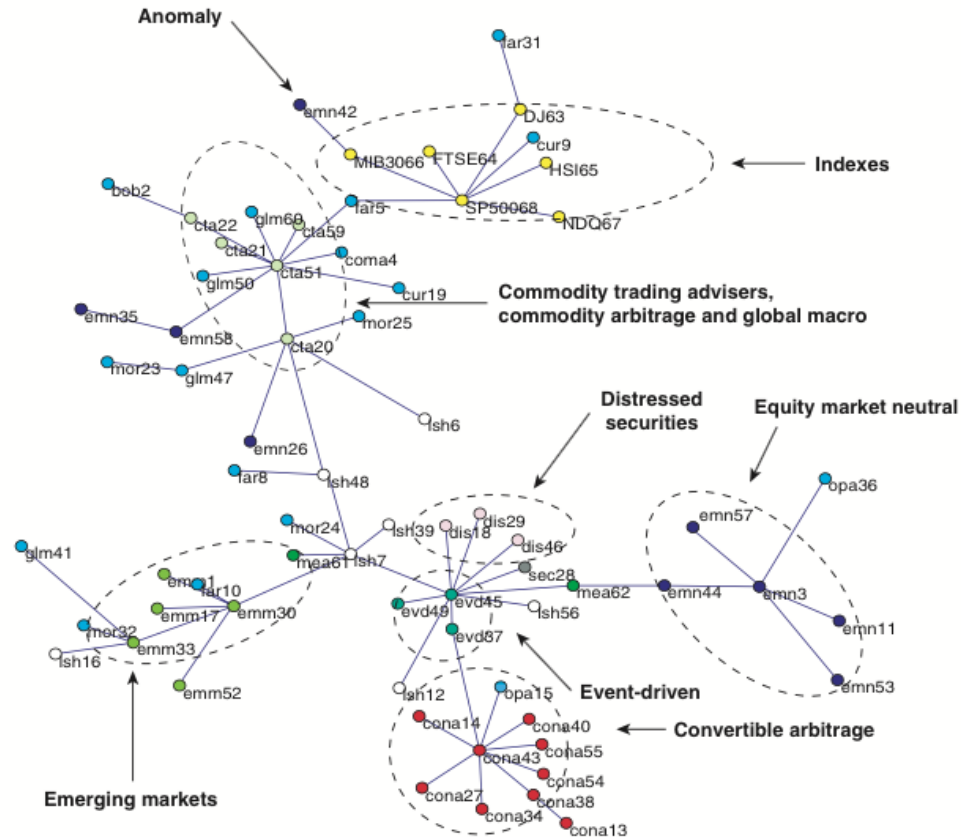


Figure 2: From Miceli and Susinno (2003)

1 **Das (2003):** “Hedge Fund Classification using K-means Clustering  
2 Method,” *9th International Conference on Computing in Economics and*  
3 *Finance*.

4 The author proposes using clustering as a means of categorizing hedge  
5 funds, hence improving on the existing “myriad of classifications, some  
6 overlapping and some mutually exclusive.”

7 **Craighead and Klemesrud (2002):** “Stock Selection Based on Cluster and  
8 Outlier Analysis,” Nationwide Financial, Columbus, OH.

9 The authors use clustering to identify outliers that could potentially create  
10 problems with portfolio selection algorithms.



## 1 K-Means Clustering

2 The **K-means clustering algorithm** is a relatively simple and intuitive ap-  
3 proach to dividing the sample  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  into  $K$  distinct groups.

4 Despite repeated extensions and enhancements of the method and the  
5 growth of other clustering methods, K-means remains popular and use-  
6 ful.

7 The algorithm is based on a very natural measure of within-cluster hetero-  
8 geneity. Observations are assigned to clusters in an effort to minimize this  
9 measure.

10 It is assumed that  $K$  is fixed by the user; the choice of  $K$  is a tricky issue.

11 I will adopt the notation used by ISL, Section 10.3.1.

1 First, define  $C_k$  to be the set of indices that belong to cluster  $k$ , for  $k =$   
2  $1, 2, \dots, K$ .

3 Each observation is assigned to exactly one cluster. In other words,  $C_k \subseteq$   
4  $\{1, 2, \dots, n\}$  with

5 
$$\bigcup_{k=1}^K C_k = \{1, 2, \dots, n\} \quad \text{and} \quad C_k, C_\ell \text{ disjoint for } k \neq \ell.$$

6 With K-means, the similarity within cluster  $k$  is measured by

7 
$$W(C_k) = \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2$$

8 where  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $\bar{\mathbf{x}}_k$  is the cluster **centroid**, calculated as the  
9 sample mean of the cluster members.

1 **Exercise:** Under what circumstance(s) will  $W(C_k)$  be small?

2

3

4

5

6

7

8

9

10

K-means clustering seeks to find the allocation of the observations that minimizes

$$\sum_{k=1}^K W(C_k)$$

Comparing every possible allocation is not realistic. (Note that there is no restriction on the sizes of the clusters.)

Instead, the following algorithm is utilized: (Algorithm 10.1 in ISL.)

Randomly assign a number from 1 to  $K$  to each of the observations. These are the initial cluster assignments. Then repeat until no changes are made to cluster assignments:

1. Find the centroid for each of the  $K$  clusters.
2. Re-assign each observation to the cluster whose centroid is closest as measured by Euclidean distance.

- 1 **Comment:** Note that the criterion is guaranteed to not increase with each  
2 iteration of the algorithm, but the search could get caught in a local mini-  
3 mum (and not the global minimum).
- 4 To address this problem, it is standard to repeat the algorithm several times  
5 with different (randomly chosen) starting allocations. If you keep track  
6 of the criterion achieved from each repetition, you can easily determine  
7 which allocation is the best.

## 1 **K-means in R**

2 K-means is implemented in R using the function `kmeans()`.

3 The syntax is simple:

```
> kmout = kmeans(datamat, centers = 5, nstart = 10)
```

4 Here I have specified that I want  $K = 5$  clusters, and that I want the basic  
5 algorithm repeated ten times in order to increase the chances of finding the  
6 global minimum.

7 Note that `datamat` must be formatted so that each of the vectors to be  
8 clustered are stored in its `rows`. Use the transpose function `t()`, if needed.

## 1 Example: Clustering Stock Time Series

2 In this example, we will cluster our data set of random-chosen NYSE  
3 stocks, with data taken from August 18, 2017 to September 29, 2017. This  
4 is a total of 30 trading days.

5 The data can be read in using the command

```
> stocksample = read.table("stocksample.txt", header=T,  
+                          sep="\t", comment.char="")  
>  
> stocksamplescl = apply(stocksample[,5:34], 1, scale)
```

6 Note that the data are formatted such that each row is a different stock; the  
7 columns are the different days. As before, each time series has been scaled  
8 so that it has a mean of zero and variance of one. It would be difficult to  
9 compare the series without this transformation.

1 The call to `kmeans()` is below. The matrix needs to be transposed:

```
> kmout = kmeans(t(stocksamplescl), centers = 10,  
+               nstart = 10)
```

2 Some comments on the output of `kmeans()`:

3 1. The rows of `kmout$centers` holds the centers of the  $K$  clusters. In  
4 our example, each center is a vector of length 30.

5 2. `kmout$cluster` specifies the cluster that each of the observations  
6 has been assigned to.

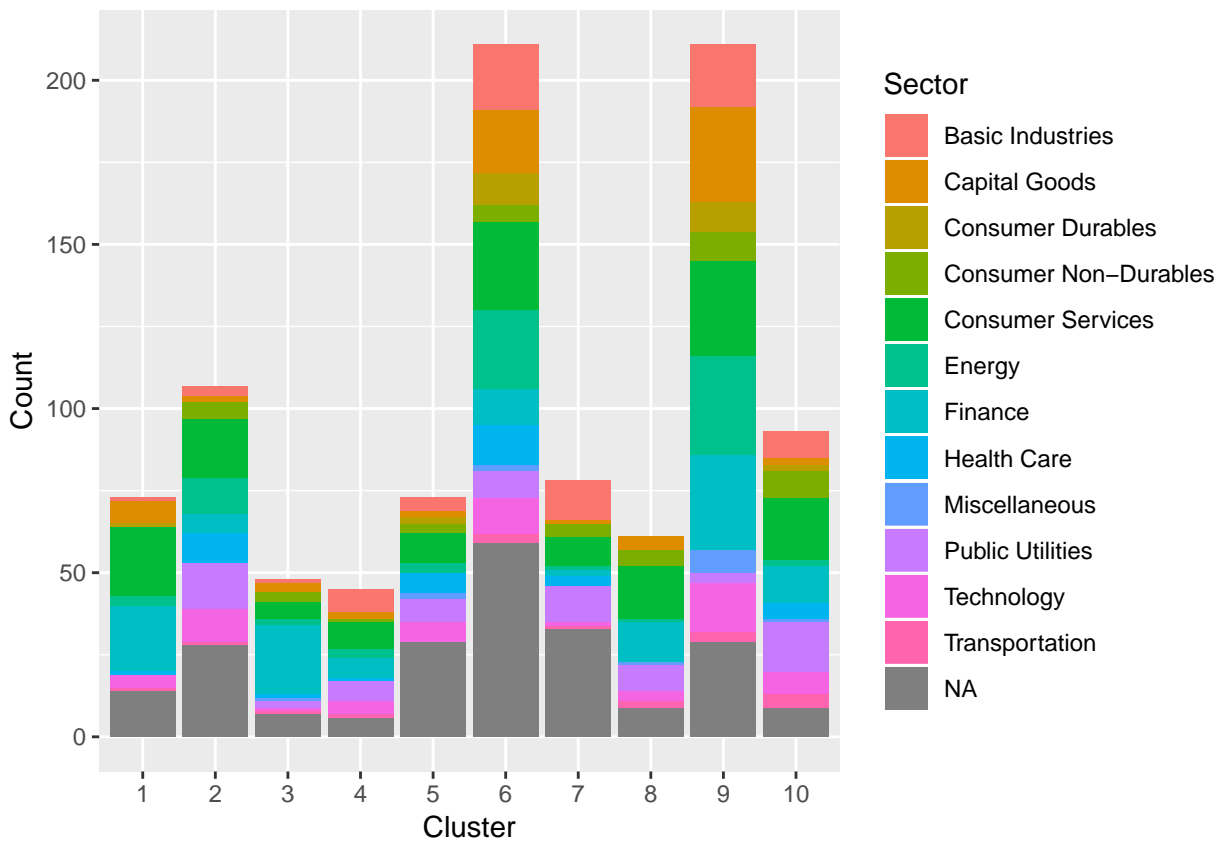
7 3. The value of  $W(C_k)$  for each of the clusters can be found from  
8 `kmout$withinss`.

9 4. The size of each cluster is found from `kmout$size`.



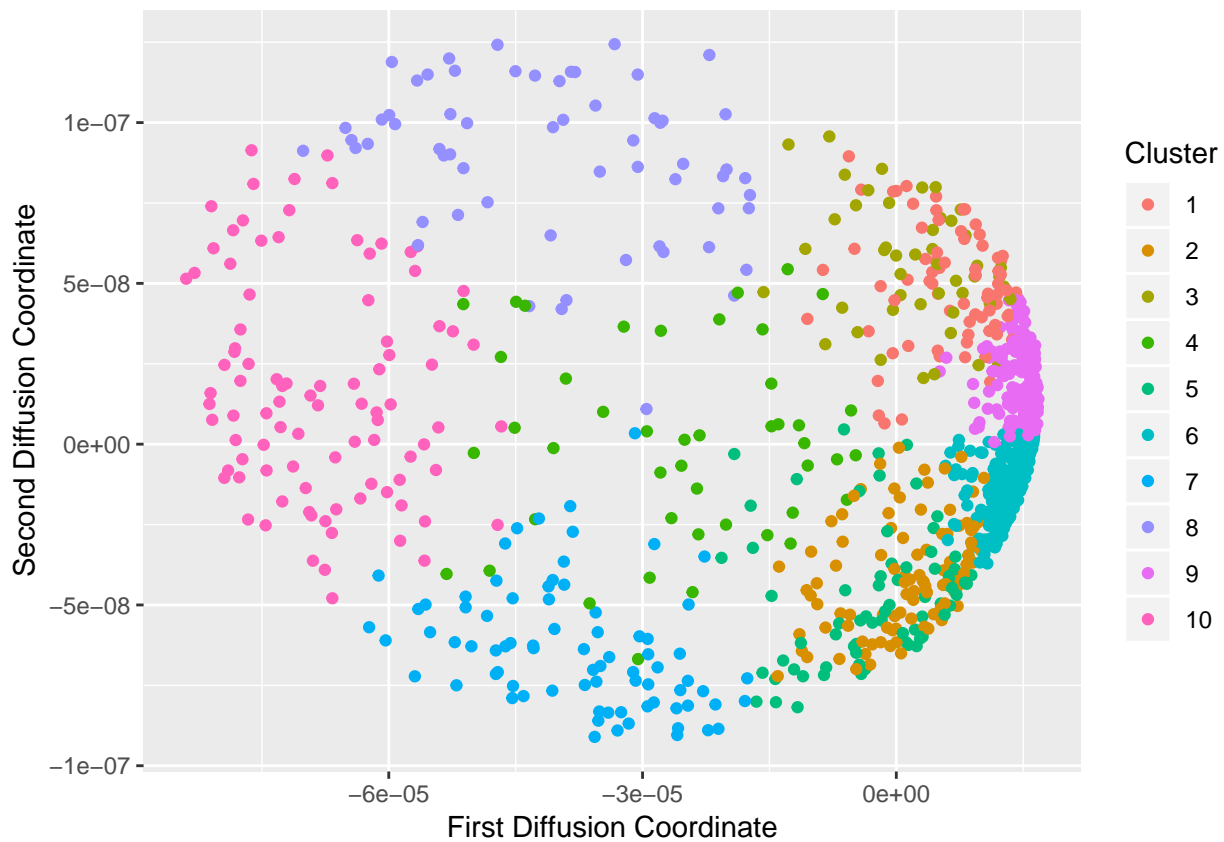
- 1 We can inspect how the clusters vary with respect to sector:

```
> stocksample$kmclust = factor(kmout$cluster)
>
> ggplot(stocksample, aes(x=kmclust, fill=sector)) +
+   geom_bar() +
+   labs(x="Cluster", y="Count", fill="Sector")
```



- 1 It is also interesting to compare K-means with the output of a dimension
- 2 reduction technique seen previously:

```
> stockdistmat = dist(t(stocksamplescl))
> stockdiffmap = diffuse(stockdistmat, eps.val=50, t=10)
> stocksample$dmap1 = stockdiffmap$X[,1]
> stocksample$dmap2 = stockdiffmap$X[,2]
>
> ggplot(stocksample, aes(x=dmap1, y=dmap2, color=kmclust)) +
+   geom_point() +
+   labs(x="First Diffusion Coordinate",
+        y="Second Diffusion Coordinate", color="Cluster")
```



- 1 **Exercise:** Would it seem sensible to run K-means directly on the low-
- 2 dimensional representation, e.g., on the diffusion coordinates?

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

6 \_\_\_\_\_

7 \_\_\_\_\_

8 \_\_\_\_\_

9 \_\_\_\_\_

10 \_\_\_\_\_

## 1 Hierarchical Clustering

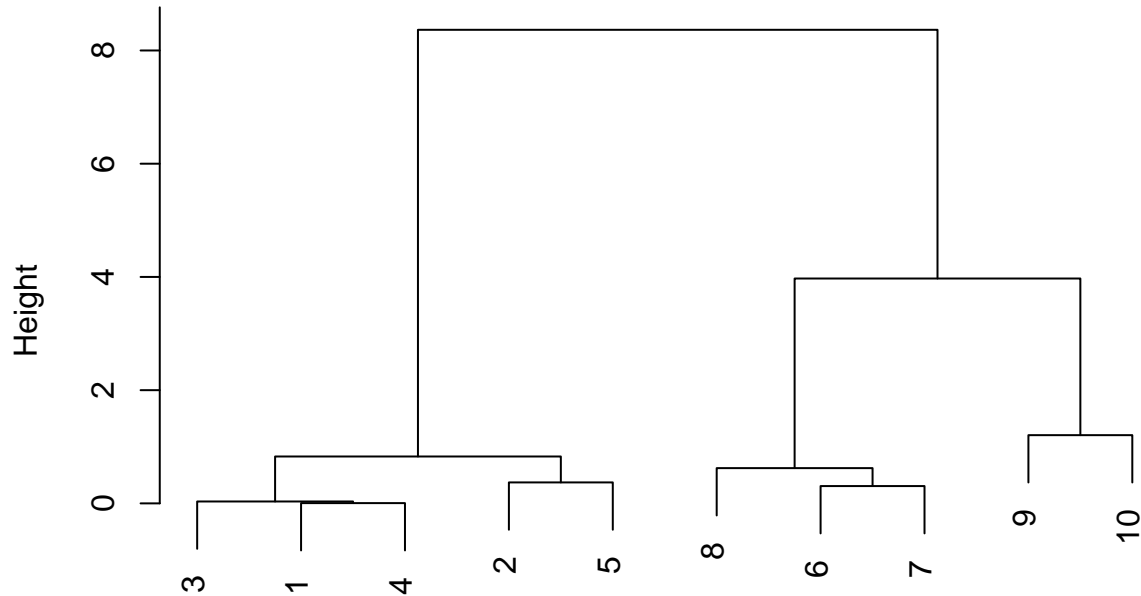
2 Another class of clustering procedures, called **hierarchical clustering**, are  
3 characterized by the distinctive **dendrograms** they create to depict the re-  
4 lationships between observations.

5 An appealing feature of these methods is that the creation of the dendro-  
6 gram does not require the specification of the number of clusters. Some-  
7 times, the dendrogram may reveal an “obvious” division among the obser-  
8 vations, and hence suggest the number of clusters needed.

9 This clustering method allows the user to specify a measure of dissimilar-  
10 ity. Often, Euclidean distance is used (just as in K-means).

- 1 It is useful to begin by understanding how to interpret the dendrogram.
- 2 An example is given below.

**Cluster Dendrogram**



Observations which are similar are “linked” at a low height in the dendrogram. For instance, to link observations 6 and 7, you need only go up to a height of approximately 0.5. These two observations are quite similar. But, in order to link observations 2 and 10, you need to go to a height of over 8. These are quite dissimilar.

The units for “height” match the units of the dissimilarity metric that the user provides.

**Exercise:** Which pair are more similar: 8 and 6, or 9 and 10?

---

---



- 1 Clusters can be formed by “cutting” the dendrogram at a user-chosen
- 2 height. It may seem “clear” by looking at the dendrogram that cutting at
- 3 around a height of 3 would be a good choice. The results is three clusters.
- 4 Note that in this case I generated the ten observations from the normal
- 5 distribution with a standard deviation of 0.5. The mean for observations 1
- 6 thru 5 was 0, for observations 6 thru 8 was 5, and for observations 9 and
- 7 10 was 7. This explains the structure in the dendrogram.

## 1 Building the Dendrogram

2 The process of constructing the dendrogram is quite intuitive.

3 First, the two observations which are the most similar (as measured by the  
4 chosen dissimilarity metric) and joined together, at a height equal to their  
5 dissimilarity. In our example above, observations 1 and 4 were the first to  
6 be joined, at a height very close to zero.

7 This process then continues, with the most similar objects joined together  
8 at the appropriate height. **But:** Once observations are joined, they are  
9 **treated as a unit** from that point on. For instance, in the second step in  
10 the construction of the dendrogram above, observation 3 was joined with  
11 the **unit consisting of 1 and 4**.

12 This process continues until all observations are joined.

The obvious question is thus: **How is the dissimilarity of collections of observations calculated?** There are a few different approaches. Two of the standard ones are as follows:

With **complete linkage**, the distance between sets of observations,  $S_i$  and  $S_j$ , is defined to be the **maximal** dissimilarity between all pairs  $x \in S_i$  and  $y \in S_j$ .

With **single linkage**, the distance between sets of observations,  $S_i$  and  $S_j$ , is defined to be the **minimal** dissimilarity between all pairs  $x \in S_i$  and  $y \in S_j$ .

Less-commonly used options are **centroid linkage**, in which the distance between  $S_i$  and  $S_j$  is calculated as the dissimilarity between the centroids of the two groups, and **average linkage**, calculated as the average distance between all pairs of elements in  $S_i$  and  $S_j$ .

## 1 Hierarchical Clustering in R

2 The function `hclust()` performs hierarchical clustering in R.

3 Remember that this method is built around a user-chosen measure of dis-  
4 similarity. Hence, the primary input argument to `hclust()` is a symmet-  
5 ric  $n$  by  $n$  distance matrix. The diagonal of this matrix should consist of  
6 zeros.

- 1 The syntax for complete linkage clustering is as follows:

```
> hcout = hclust(dist(t(stocksamplescl[,1:100])),  
+               method="complete")
```

- 2 (Note that only the first 100 stocks are used, since the dendrogram becomes  
3 difficult to read otherwise. There is no such practical limit in using this  
4 method, however.)

- 5 Other choices for method include single, centroid, and average.

- 6 The `plot()` function, when applied to the output of `hclust()`, produces  
7 the dendrogram:

```
> plot(hcout, labels=stocksample$name[1:100],  
+       cex=0.35, sub="", xlab="")
```



- 1 Another useful function is `cutree()`. This will create return cluster mem-
- 2 bership for each of the observation, based on a desired number of clusters,
- 3 or on a desired height.
- 4 For instance, to cut our dendrogram into six clusters, we specify `k=6`:

```
> cutree(hcout, k=6) [1:10]  
[1] 1 2 3 2 2 1 1 1 1 2
```

- 5 And to cut our our dendrogram at a height of 10, specify `h=10`:

```
> cutree(hcout, h=10) [1:10]  
[1] 1 2 3 2 2 1 1 1 1 2
```

- 1 We can again inspect how the clusters vary with respect to sector:

```
> stocksample$hcclust = factor(cutree(hcout,10))  
>  
> ggplot(stocksample, aes(x=hcclust, fill=sector)) +  
+   geom_bar() +  
+   labs(x="Cluster", y="Count", fill="Sector")
```



