

### 10.2.4 Other Uses for Principal Components

We saw in Section 6.3.1 that we can perform regression using the principal component score vectors as features. In fact, many statistical techniques, such as regression, classification, and clustering, can be easily adapted to use the  $n \times M$  matrix whose columns are the first  $M \ll p$  principal component score vectors, rather than using the full  $n \times p$  data matrix. This can lead to *less noisy* results, since it is often the case that the signal (as opposed to the noise) in a data set is concentrated in its first few principal components.

## 10.3 Clustering Methods

*Clustering* refers to a very broad set of techniques for **finding subgroups, or clusters, in a data set**. When we cluster the observations of a data set, we seek to partition them into distinct groups so that **the observations within each group are quite similar to each other**, while observations in different groups are quite different from each other. Of course, to make this concrete, we must define what it means for two or more observations to be *similar* or *different*. Indeed, this is often a domain-specific consideration that must be made based on knowledge of the data being studied. clustering

For instance, suppose that we have a set of  $n$  observations, each with  $p$  features. The  $n$  observations could correspond to tissue samples for patients with breast cancer, and the  $p$  features could correspond to measurements collected for each tissue sample; these could be clinical measurements, such as tumor stage or grade, or they could be gene expression measurements. We may have a reason to believe that **there is some heterogeneity among the  $n$  tissue samples**; for instance, perhaps there are a few different *unknown* subtypes of breast cancer. Clustering could be used to find these subgroups. **This is an unsupervised problem because we are trying to discover structure—in this case, distinct clusters—on the basis of a data set.** **The goal in supervised problems, on the other hand, is to try to predict some outcome vector such as survival time or response to drug treatment.**

Both clustering and PCA seek to simplify the data via a small number of summaries, but their mechanisms are different:

- **PCA looks to find a low-dimensional representation of the observations that explain a good fraction of the variance;**
- **Clustering looks to find homogeneous subgroups among the observations.**

Another application of clustering arises in marketing. We may have access to a large number of measurements (e.g. median household income, occupation, distance from nearest urban area, and so forth) for a large

number of people. Our goal is to perform *market segmentation* by identifying subgroups of people who might be more receptive to a particular form of advertising, or more likely to purchase a particular product. The task of performing market segmentation amounts to clustering the people in the data set.

Since clustering is popular in many fields, there exist a great number of clustering methods. In this section we focus on perhaps the two best-known clustering approaches: *K-means clustering* and *hierarchical clustering*. In *K-means clustering*, we seek to partition the observations into a pre-specified number of clusters. On the other hand, in *hierarchical clustering*, we do not know in advance how many clusters we want; in fact, we end up with a tree-like visual representation of the observations, called a *dendrogram*, that allows us to view at once the clusterings obtained for each possible number of clusters, from 1 to  $n$ . There are advantages and disadvantages to each of these clustering approaches, which we highlight in this chapter.

In general, we can cluster observations on the basis of the features in order to identify subgroups among the observations, or we can cluster features on the basis of the observations in order to discover subgroups among the features. In what follows, for simplicity we will discuss clustering observations on the basis of the features, though the converse can be performed by simply transposing the data matrix.

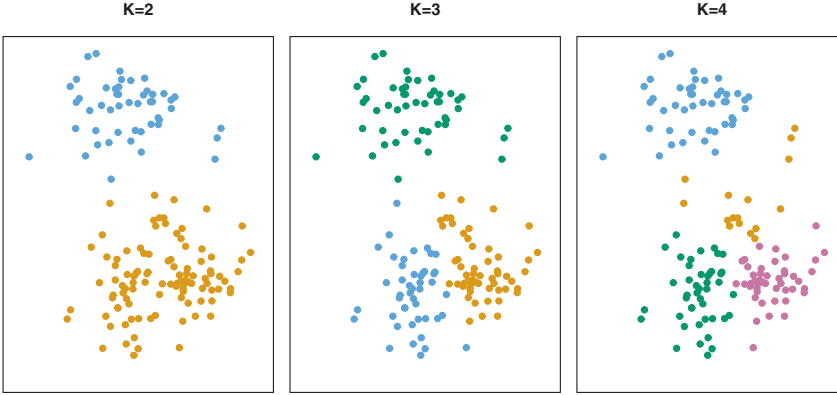
### 10.3.1 K-Means Clustering

*K-means clustering* is a simple and elegant approach for partitioning a data set into  $K$  distinct, non-overlapping clusters. To perform  $K$ -means clustering, we must first specify the desired number of clusters  $K$ ; then the  $K$ -means algorithm will assign each observation to exactly one of the  $K$  clusters. Figure 10.5 shows the results obtained from performing  $K$ -means clustering on a simulated example consisting of 150 observations in two dimensions, using three different values of  $K$ .

The  $K$ -means clustering procedure results from a simple and intuitive mathematical problem. We begin by defining some notation. Let  $C_1, \dots, C_K$  denote sets containing the indices of the observations in each cluster. These sets satisfy two properties:

1.  $C_1 \cup C_2 \cup \dots \cup C_K = \{1, \dots, n\}$ . In other words, each observation belongs to at least one of the  $K$  clusters.
2.  $C_k \cap C_{k'} = \emptyset$  for all  $k \neq k'$ . In other words, the clusters are non-overlapping; no observation belongs to more than one cluster.

For instance, if the  $i$ th observation is in the  $k$ th cluster, then  $i \in C_k$ . The idea behind  $K$ -means clustering is that a good clustering is one for which the within-cluster variation is as small as possible. The within-cluster variation



**FIGURE 10.5.** A simulated data set with 150 observations in two-dimensional space. Panels show the results of applying K-means clustering with different values of  $K$ , the number of clusters. The color of each observation indicates the cluster to which it was assigned using the K-means clustering algorithm. Note that there is no ordering of the clusters, so the cluster coloring is arbitrary. These cluster labels were not used in clustering; instead, they are the outputs of the clustering procedure.

for cluster  $C_k$  is a measure  $W(C_k)$  of the amount by which the observations within a cluster differ from each other. Hence we want to solve the problem

$$\text{minimize}_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K W(C_k) \right\}. \quad (10.9)$$

In words, this formula says that we want to partition the observations into  $K$  clusters such that the total within-cluster variation, summed over all  $K$  clusters, is as small as possible.

Solving (10.9) seems like a reasonable idea, but in order to make it actionable we need to define the within-cluster variation. There are many possible ways to define this concept, but by far the most common choice involves squared Euclidean distance. That is, we define

$$W(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2, \quad (10.10)$$

where  $|C_k|$  denotes the number of observations in the  $k$ th cluster. In other words, the within-cluster variation for the  $k$ th cluster is the sum of all of the pairwise squared Euclidean distances between the observations in the  $k$ th cluster, divided by the total number of observations in the  $k$ th cluster. Combining (10.9) and (10.10) gives the optimization problem that defines K-means clustering,

$$\text{minimize}_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}. \quad (10.11)$$

Now, we would like to find an algorithm to solve (10.11)—that is, a method to partition the observations into  $K$  clusters such that the objective of (10.11) is minimized. This is in fact a very difficult problem to solve precisely, since there are almost  $K^n$  ways to partition  $n$  observations into  $K$  clusters. This is a huge number unless  $K$  and  $n$  are tiny! Fortunately, a very simple algorithm can be shown to provide a local optimum—a *pretty good solution*—to the  $K$ -means optimization problem (10.11). This approach is laid out in Algorithm 10.1.

---

**Algorithm 10.1** *K-Means Clustering*

---

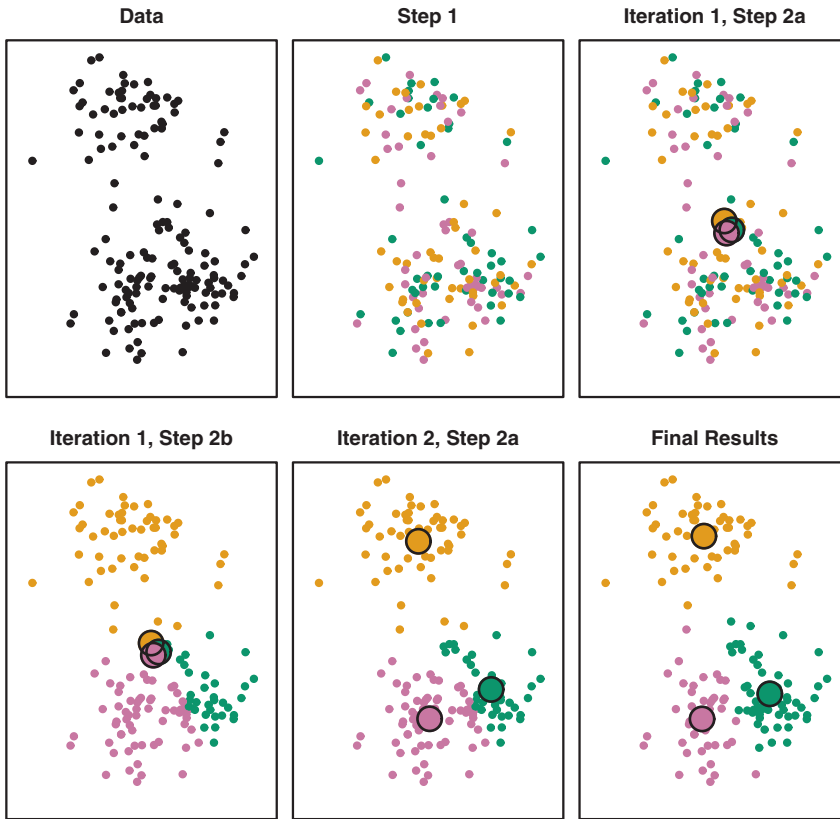
1. Randomly assign a number, from 1 to  $K$ , to each of the observations. These serve as initial cluster assignments for the observations.
  2. Iterate until the cluster assignments stop changing:
    - (a) For each of the  $K$  clusters, compute the *cluster centroid*. The *kth cluster centroid* is the *vector of the  $p$  feature means for the observations in the  $k$ th cluster*.
    - (b) Assign each observation to the cluster whose centroid is closest (where *closest* is defined using Euclidean distance).
- 

Algorithm 10.1 is guaranteed to decrease the value of the objective (10.11) at each step. To understand why, the following identity is illuminating:

$$\frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2, \quad (10.12)$$

where  $\bar{x}_{kj} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{ij}$  is the mean for feature  $j$  in cluster  $C_k$ . In Step 2(a) the cluster means for each feature are the constants that minimize the sum-of-squared deviations, and in Step 2(b), reallocating the observations can only improve (10.12). This means that as the algorithm is run, the clustering obtained will continually improve until the result no longer changes; the objective of (10.11) will never increase. When the result no longer changes, a *local optimum* has been reached. Figure 10.6 shows the progression of the algorithm on the toy example from Figure 10.5.  $K$ -means clustering derives its name from the fact that in Step 2(a), the cluster centroids are computed as the mean of the observations assigned to each cluster.

Because the  $K$ -means algorithm finds a local rather than a global optimum, the results obtained will depend on the initial (random) cluster assignment of each observation in Step 1 of Algorithm 10.1. For this reason, it is important to run the algorithm multiple times from different random



**FIGURE 10.6.** The progress of the  $K$ -means algorithm on the example of Figure 10.5 with  $K=3$ . Top left: the observations are shown. Top center: in Step 1 of the algorithm, each observation is randomly assigned to a cluster. Top right: in Step 2(a), the cluster centroids are computed. These are shown as large colored disks. Initially the centroids are almost completely overlapping because the initial cluster assignments were chosen at random. Bottom left: in Step 2(b), each observation is assigned to the nearest centroid. Bottom center: Step 2(a) is once again performed, leading to new cluster centroids. Bottom right: the results obtained after ten iterations.

initial configurations. Then one selects the *best* solution, i.e. that for which the objective (10.11) is smallest. Figure 10.7 shows the local optima obtained by running  $K$ -means clustering six times using six different initial cluster assignments, using the toy data from Figure 10.5. In this case, the best clustering is the one with an objective value of 235.8.

As we have seen, to perform  $K$ -means clustering, we must decide how many clusters we expect in the data. The problem of selecting  $K$  is far from simple. This issue, along with other practical considerations that arise in performing  $K$ -means clustering, is addressed in Section 10.3.3.



**FIGURE 10.7.** *K*-means clustering performed six times on the data from Figure 10.5 with  $K = 3$ , each time with a different random assignment of the observations in Step 1 of the *K*-means algorithm. Above each plot is the value of the objective (10.11). Three different local optima were obtained, one of which resulted in a smaller value of the objective and provides better separation between the clusters. Those labeled in red all achieved the same best solution, with an objective value of 235.8.

### 10.3.2 Hierarchical Clustering

One potential disadvantage of *K*-means clustering is that it requires us to pre-specify the number of clusters  $K$ . *Hierarchical clustering* is an alternative approach which does not require that we commit to a particular choice of  $K$ . Hierarchical clustering has an added advantage over *K*-means clustering in that it results in an attractive tree-based representation of the observations, called a *dendrogram*.

In this section, we describe *bottom-up* or *agglomerative* clustering. This is the most common type of hierarchical clustering, and refers to the fact that a dendrogram (generally depicted as an upside-down tree; see

bottom-up  
agglomerative