

Chapter 7. Unsupervised Learning

The term *unsupervised learning* refers to statistical methods that extract meaning from data without training a model on labeled data (data where an outcome of interest is known). In Chapters 4 and 5, the goal is to build a model (set of rules) to predict a response from a set of predictor variables. Unsupervised learning also constructs a model of the data, but does not distinguish between a response variable and predictor variables.

Unsupervised learning can have different possible goals. In some cases, it can be used to create a predictive rule in the absence of a labeled response. *Clustering* methods can be used to identify meaningful groups of data. For example, using the web clicks and demographic data of a user on a website, we may be able to group together different types of users. The website could then be personalized to these different types.

In other cases, the goal may be to *reduce the dimension* of the data to a more manageable set of variables. This reduced set could then be used as input into a predictive model, such as regression or classification. For example, we may have thousands of sensors to monitor an industrial process. By reducing the data to a smaller set of features, we may be able to build a more powerful and interpretable model to predict process failure than by including data streams from thousands of sensors.

Finally, unsupervised learning can be viewed as an extension of the exploratory data analysis (see Chapter 1) to situations where you are confronted with a large number of variables and records. The aim is to gain insight into a set of data and how the different variables relate to each other. Unsupervised techniques give ways to sift through and analyze these variables and discover relationships.

UNSUPERVISED LEARNING AND PREDICTION

Unsupervised learning can play an important role for prediction, both for regression and classification problems. In some cases, we want to predict a category in the absence of any labeled data. For example, we might want to predict the type of vegetation in an area from a set of satellite sensory data. Since we don't have a response variable to train a model, clustering gives us a way to identify common patterns and categorize the regions.

Clustering is an especially important tool for the “cold-start problem.” In these types of problems, such as launching a new marketing campaign or identifying potential new types of fraud or spam, we initially may not have any response to train a model. Over time, as data is collected, we can learn more about the system and build a traditional predictive model. But clustering helps us start the learning process more quickly by identifying population segments.

Unsupervised learning is also important as a building block for regression and classification techniques. With big data, if a small subpopulation is not well represented in the overall population, the trained model may not perform well for that subpopulation. With clustering, it is possible to identify and label subpopulations. Separate models can then be fit to the different subpopulations. Alternatively, the subpopulation can be represented with its own feature, forcing the overall model to explicitly consider subpopulation identity as a predictor.

Principal Components Analysis

Often, variables will vary together (covary), and some of the variation in one is actually duplicated by variation in another. Principal components analysis (PCA) is a technique to discover the way in which numeric variables covary.¹

KEY TERMS FOR PRINCIPAL COMPONENTS ANALYSIS

Principal component

A linear combination of the predictor variables.

Loadings

The weights that transform the predictors into the components.

Synonym

Weights

Screeplot

A plot of the variances of the components, showing the relative importance of the components.

The idea in PCA is to combine multiple numeric predictor variables into a smaller set of variables, which are weighted linear combinations of the original set. The smaller set of variables, the *principal components*, “explains” most of the variability of the full set of variables, reducing the dimension of the data. The weights used to form the principal components reveal the relative contributions of the original variables to the new principal components.

PCA was first **proposed by Karl Pearson**. In what was perhaps the first paper on unsupervised learning, Pearson recognized that in many problems there is variability in the predictor variables, so he developed PCA as a technique to model this variability. PCA can be viewed as the unsupervised version of linear discriminant analysis; see **“Discriminant Analysis”**.

A Simple Example

For two variables, X_1 and X_2 , there are two principal components Z_i ($i = 1$ or 2):

$$Z_i = w_{i,1}X_1 + w_{i,2}X_2$$

The weights ($w_{i,1}$, $w_{i,2}$) are known as the component *loadings*. These transform the original variables into the principal components. The first principal component, Z_1 , is the linear combination that best explains the total variation. The second principal component, Z_2 , explains the remaining variation (it is also the linear combination that is the worst fit).

NOTE

It is also common to compute principal components on deviations from the means of the predictor variables, rather than on the values themselves.

You can compute principal components in R using the `princomp` function. The following performs a PCA on the stock price returns for Chevron (CVX) and ExxonMobil (XOM):

```
oil_px <- sp500_px[, c('CVX', 'XOM')]
pca <- princomp(oil_px)
pca$loadings

Loadings:
  Comp.1 Comp.2
CVX -0.747  0.665
XOM -0.665 -0.747
```

The weights for CVX and XOM for the first principal component are -0.747 and -0.665 and for the second principal component they are 0.665 and -0.747 . How to interpret this? The first principal component is essentially an average of CVX and XOM, reflecting the correlation between the two energy companies. The second principal component measures when the stock prices of CVX and XOM diverge.

It is instructive to plot the principal components with the data:

```
loadings <- pca$loadings
ggplot(data=oil_px, aes(x=CVX, y=XOM)) +
  geom_point(alpha=.3) +
  stat_ellipse(type='norm', level=.99) +
  geom_abline(intercept = 0, slope = loadings[2,1]/loadings[1,1]) +
  geom_abline(intercept = 0, slope = loadings[2,2]/loadings[1,2])
```

The result is shown in **Figure 7-1**.

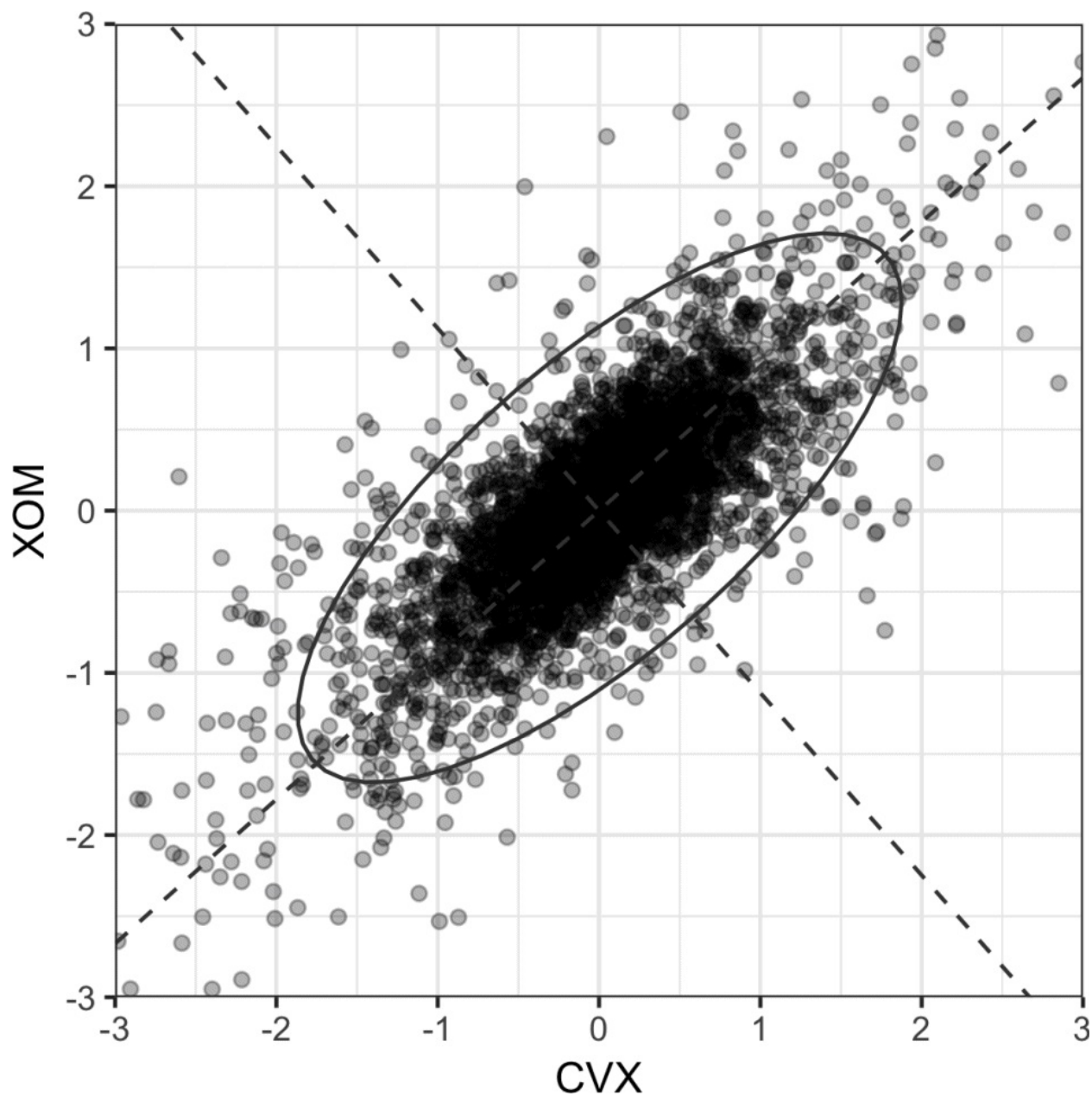


Figure 7-1. The principal components for the stock returns for Chevron and ExxonMobil

The solid dashed lines show the two principal components: the first one is along the long axis of the ellipse and the second one is along the short axis. You can see that a majority of the variability in the two stock returns is explained by the first principal component. This makes sense since energy stock prices tend to move as a group.

NOTE

The weights for the first principal component are both negative, but reversing the sign of all the weights does not change the principal component. For example, using weights of 0.747 and 0.665 for the first principal component is equivalent to the negative weights, just as an infinite line defined by the origin and 1,1 is the same as one defined by the origin and $-1, -1$.

Computing the Principal Components

Going from two variables to more variables is straightforward. For the first component, simply include the additional predictor variables in the linear combination, assigning weights that optimize the collection of the covariation from all the predictor variables into this first principal component (*covariance* is the statistical term; see “**Covariance Matrix**”). Calculation of principal components is a classic statistical method, relying on either the correlation matrix of the data or the covariance matrix, and it executes rapidly, not relying on iteration. As noted earlier, it works only with numeric variables, not categorical ones. The full process can be described as follows:

1. In creating the first principal component, PCA arrives at the linear combination of predictor variables that maximizes the percent of total variance explained.
2. This linear combination then becomes the first “new” predictor, Z_1 .
3. PCA repeats this process, using the same variables, with different weights to create a second new predictor, Z_2 . The weighting is done such that Z_1 and Z_2 are uncorrelated.
4. The process continues until you have as many new variables, or components, Z_i as original variables X_i .
5. Choose to retain as many components as are needed to account for most of the variance.
6. The result so far is a set of weights for each component. The final step is to convert the original data into new principal component scores by applying the weights to the original values. These new scores can then be used as the reduced set of predictor variables.

Interpreting Principal Components

The nature of the principal components often reveals information about the structure of the data. There are a couple of standard visualization displays to help you glean insight about the principal components. One such method is a *Screeplot* to visualize the relative importance of principal components (the name derives from the resemblance of the plot to a scree slope). The following is an example for a few top companies in the S&P 500:

```
syms <- c( 'AAPL', 'MSFT', 'CSCO', 'INTC', 'CVX', 'XOM',  
           'SLB', 'COP', 'JPM', 'WFC', 'USB', 'AXP', 'WMT', 'TGT', 'HD', 'COST')  
top_sp <- sp500_px[row.names(sp500_px) >= '2005-01-01', syms]  
sp_pca <- princomp(top_sp)  
screeplot(sp_pca)
```

As seen in **Figure 7-2**, the variance of the first principal component is quite large (as is often the case), but the other top principal components are significant.

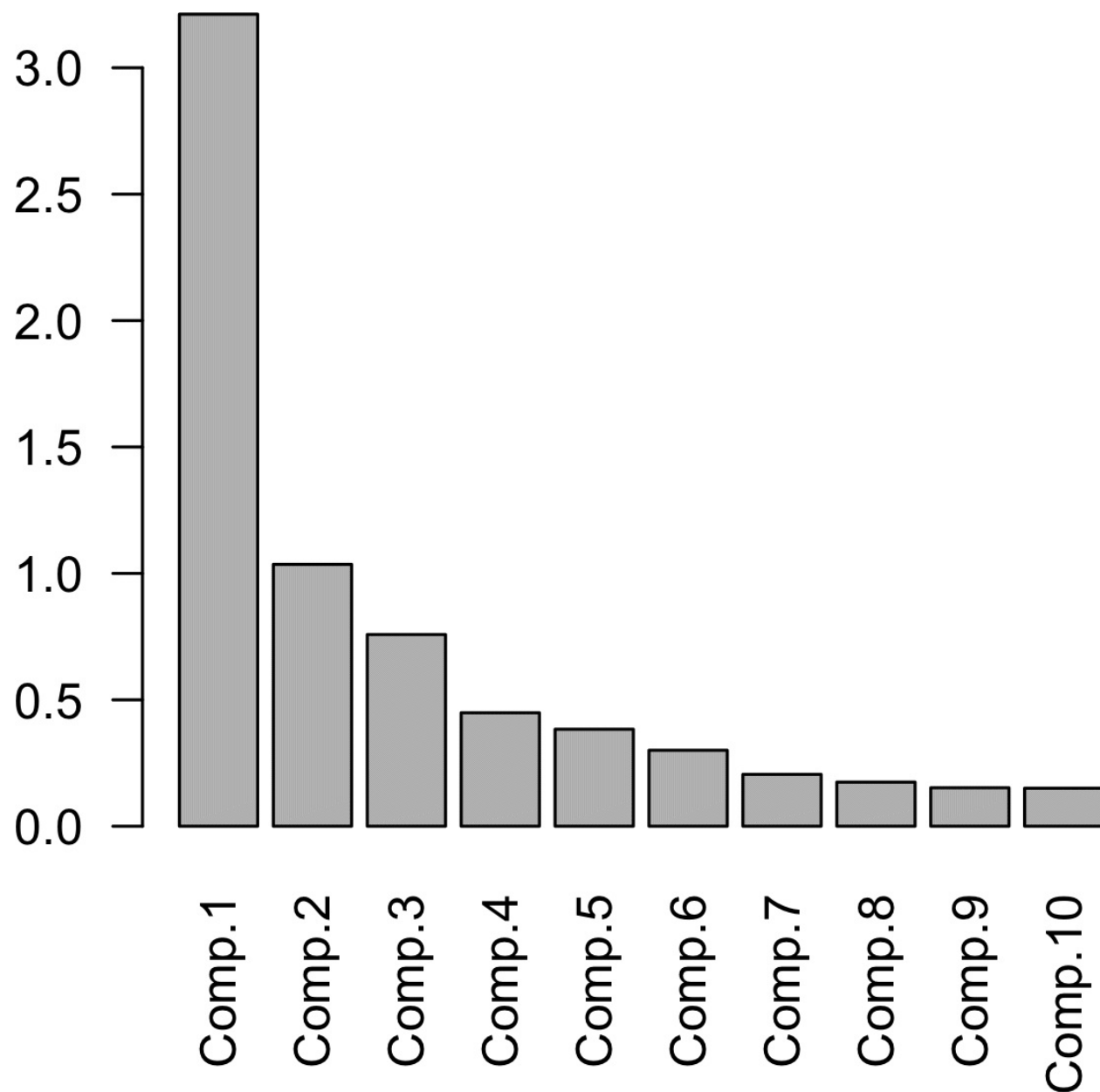


Figure 7-2. A screeplot for a PCA of top stocks from the S&P 500

It can be especially revealing to plot the weights of the top principal components. One way to do this is to use the `gather` function from the `tidyr` package in conjunction with `ggplot`:

```
library(tidyr)
loadings <- sp_pca$loadings[,1:5]
loadings$Symbol <- row.names(loadings)
loadings <- gather(loadings, "Component", "Weight", -Symbol)
ggplot(loadings, aes(x=Symbol, y=Weight)) +
  geom_bar(stat='identity') +
```

```
facet_grid(Component ~ ., scales='free_y')
```

The loadings for the top five components are shown in **Figure 7-3**. The loadings for the first principal component have the same sign: this is typical for data in which all the columns share a common factor (in this case, the overall stock market trend). The second component captures the price changes of energy stocks as compared to the other stocks. The third component is primarily a contrast in the movements of Apple and CostCo. The fourth component contrasts the movements of Schlumberger to the other energy stocks. Finally, the fifth component is mostly dominated by financial companies.

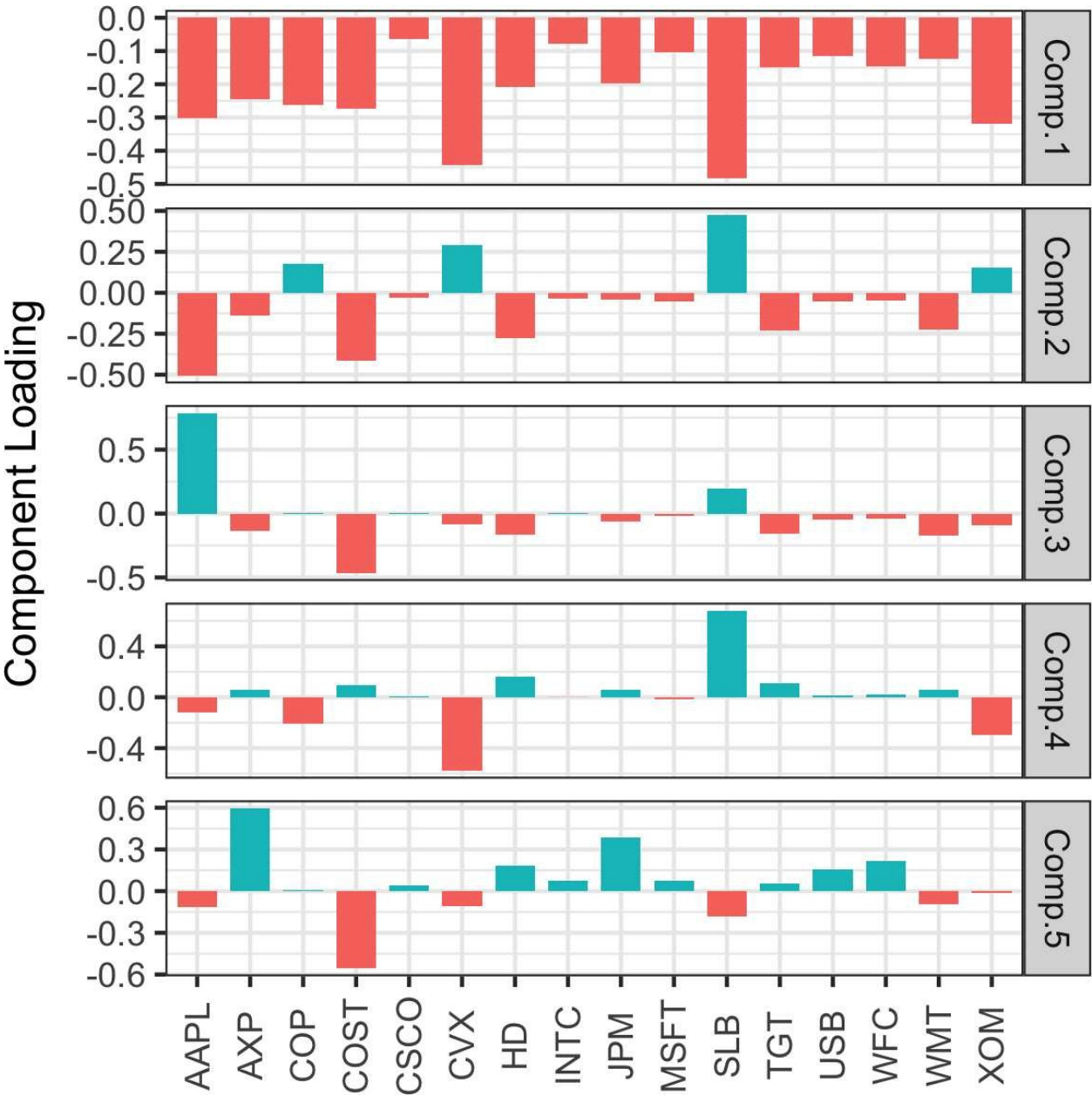


Figure 7-3. The loadings for the top five principal components of stock price returns

HOW MANY COMPONENTS TO CHOOSE?

If your goal is to reduce the dimension of the data, you must decide how many principal components to select. The most common approach is to use an ad hoc rule to select the components that explain “most” of the variance. You can do this visually through the screeplot; for example, in [Figure 7-2](#), it would be natural to restrict the analysis to the top five components. Alternatively, you could select the top components such that the cumulative variance exceeds a threshold, such as 80%. Also, you can inspect the loadings to determine if the component has an intuitive interpretation. Cross-validation provides a more formal method to select the number of significant components (see [“Cross-Validation”](#) for more).

KEY IDEAS FOR PRINCIPAL COMPONENTS

- Principal components are linear combinations of the predictor variables (numeric data only).
- They are calculated so as to minimize correlation between components, reducing redundancy.
- A limited number of components will typically explain most of the variance in the outcome variable.
- The limited set of principal components can then be used in place of the (more numerous) original predictors, reducing dimensionality.

Further Reading

For a detailed look at the use of cross-validation in principal components, see Rasmus Bro, K. Kjeldahl, A.K. Smilde, and Henk A. L. Kiers, “Cross-Validation of Component Models: A Critical Look at Current Methods”, *Analytical and Bioanalytical Chemistry* 390, no. 5 (2008).

K-Means Clustering

Clustering is a technique to divide data into different groups, where the records in each group are similar to one another. A goal of clustering is to identify significant and meaningful groups of data. The groups can be used directly, analyzed in more depth, or passed as a feature or an outcome to a predictive regression or classification model. *K-means* is the first clustering method to be developed; it is still widely used, owing its popularity to the relative simplicity of the algorithm and its ability to scale to large data sets.

KEY TERMS FOR K-MEANS CLUSTERING

Cluster

A group of records that are similar.

Cluster mean

The vector of variable means for the records in a cluster.

K

The number of clusters.

K-means divides the data into *K* clusters by minimizing the sum of the squared distances of each record to the *mean* of its assigned cluster. This is referred to as the *within-cluster sum of squares* or *within-cluster SS*. *K*-means does not ensure the clusters will have the same size, but finds the clusters that are the best separated.

NORMALIZATION

It is typical to normalize (standardize) continuous variables by subtracting the mean and dividing by the standard deviation. Otherwise, variables with large scale will dominate the clustering process (see “[Standardization \(Normalization, Z-Scores\)](#)”).

A Simple Example

Start by considering a data set with n records and just two variables, x and y . Suppose we want to split the data into $K = 4$ clusters. This means assigning each record (x_i, y_i) to a cluster k . Given an assignment of n_k records to cluster k , the center of the cluster (\bar{x}_k, \bar{y}_k) is the mean of the points in the cluster:

$$\bar{x}_k = \frac{1}{n_k} \sum_{i \in \text{Cluster } k} x_i$$

$$\bar{y}_k = \frac{1}{n_k} \sum_{i \in \text{Cluster } k} y_i$$

CLUSTER MEAN

In clustering records with multiple variables (the typical case), the term *cluster mean* refers not to a single number, but to the vector of means of the variables.

The sum of squares within a cluster is given by:

$$SS_k = \sum_{i \in \text{Cluster } k} (x_i - \bar{x}_k)^2 + (y_i - \bar{y}_k)^2$$

K-means finds the assignment of records that minimizes within-cluster sum of squares across all four clusters $SS_1 + SS_2 + SS_3 + SS_4$.

$$\sum_{k=1}^4 SS_i$$

K-means clustering can be used to gain insight into how the price movements of stocks tend to cluster. Note that stock returns are reported in a fashion that is, in effect, standardized, so we do not need to normalize the data. In R, K-means clustering can be performed using the `kmeans` function. For example, the following finds four clusters based on two variables: the stock returns for ExxonMobil (XOM) and Chevron (CVX):

```
df <- sp500_px[row.names(sp500_px) >= '2011-01-01', c('XOM', 'CVX')]
km <- kmeans(df, centers=4)
```

The cluster assignment for each record is returned as the `cluster` component:

```
> df$cluster <- factor(km$cluster)
> head(df)
      XOM      CVX cluster
2011-01-03 0.73680496 0.2406809      2
```

```

2011-01-04 0.16866845 -0.5845157 1
2011-01-05 0.02663055 0.4469854 2
2011-01-06 0.24855834 -0.9197513 1
2011-01-07 0.33732892 0.1805111 2
2011-01-10 0.00000000 -0.4641675 1

```

The first six records are assigned to either cluster 1 or cluster 2. The means of the clusters are also returned:

```

> centers <- data.frame(cluster=factor(1:4), km$centers)
> centers
  cluster      XOM      CVX
1       1 -0.3284864 -0.5669135
2       2  0.2410159  0.3342130
3       3 -1.1439800 -1.7502975
4       4  0.9568628  1.3708892

```

Clusters 1 and 3 represent “down” markets, while clusters 2 and 4 represent “up markets.” In this example, with just two variables, it is straightforward to visualize the clusters and their means:

```

ggplot(data=df, aes(x=XOM, y=CVX, color=cluster, shape=cluster)) +
  geom_point(alpha=.3) +
  geom_point(data=centers, aes(x=XOM, y=CVX), size=3, stroke=2)

```

The resulting plot, given by [Figure 7-4](#), shows the cluster assignments and the cluster means.

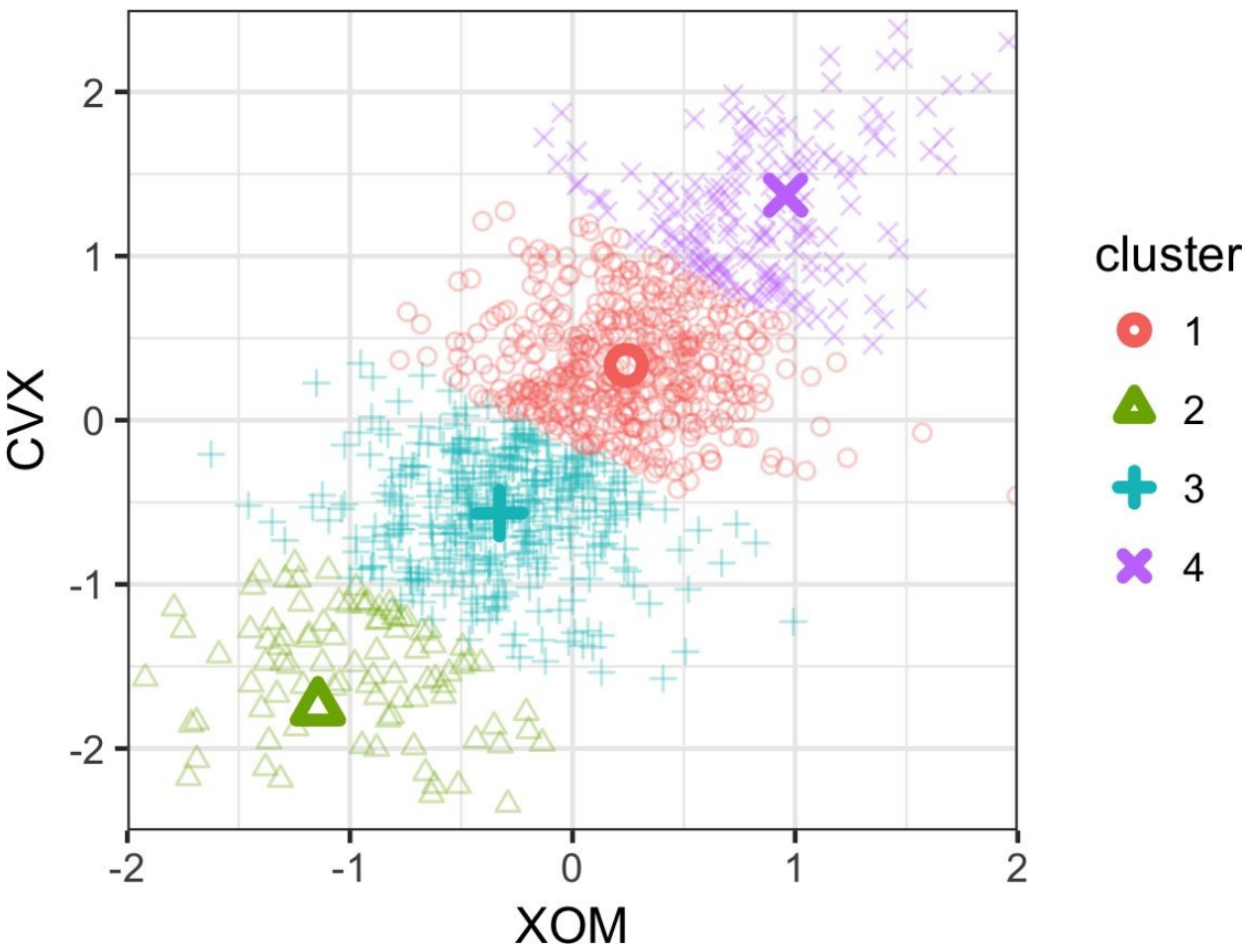


Figure 7-4. The clusters of K-means applied to stock price data for ExxonMobil and Chevron (the two cluster centers in the dense area are hard to distinguish)

K-Means Algorithm

In general, K -means can be applied to a data set with p variables X_1, \dots, X_p . While the exact solution to K -means is computationally very difficult, heuristic algorithms provide an efficient way to compute a locally optimal solution.

The algorithm starts with a user-specified K and an initial set of cluster means, then iterates the following steps:

1. Assign each record to the nearest cluster mean as measured by squared distance.
2. Compute the new cluster means based on the assignment of records.

The algorithm converges when the assignment of records to clusters does not change.

For the first iteration, you need to specify an initial set of cluster means. Usually you do this by randomly assigning each record to one of the K clusters, then finding the means of those clusters.

Since this algorithm isn't guaranteed to find the best possible solution, it is recommended to run the algorithm several times using different random samples to initialize the algorithm. When more than one set of iterations is used, the K -means result is given by the iteration that has the lowest within-cluster sum of squares.

The `nstart` parameter to the R function `kmeans` allows you to specify the number of random starts to try. For example, the following code runs K -means to find 5 clusters using 10 different starting cluster means:

```
syms <- c( 'AAPL', 'MSFT', 'CSCO', 'INTC', 'CVX', 'XOM', 'SLB', 'COP',
           'JPM', 'WFC', 'USB', 'AXP', 'WMT', 'TGT', 'HD', 'COST')
df <- sp500_px[row.names(sp500_px)>='2011-01-01', syms]
km <- kmeans(df, centers=5, nstart=10)
```

The function automatically returns the best solution out of the 10 different starting points. You can use the argument `iter.max` to set the maximum number of iterations the algorithm is allowed for each random start.

Interpreting the Clusters

An important part of cluster analysis can involve the interpretation of the clusters. The two most important outputs from `kmeans` are the sizes of the clusters and the cluster means. For the example in the previous subsection, the sizes of resulting clusters are given by this R command:

```
km$size
[1] 186 106 285 288 266
```

The cluster sizes are relatively balanced. Imbalanced clusters can result from distant outliers, or groups of records very distinct from the rest of the data — both may warrant further inspection.

You can plot the centers of the clusters using the `gather` function in conjunction with `ggplot`:

```
centers <- as.data.frame(t(centers))
names(centers) <- paste("Cluster", 1:5)
centers$Symbol <- row.names(centers)
centers <- gather(centers, "Cluster", "Mean", -Symbol)
centers$Color = centers$Mean > 0
ggplot(centers, aes(x=Symbol, y=Mean, fill=Color)) +
  geom_bar(stat='identity', position = "identity", width=.75) +
  facet_grid(Cluster ~ ., scales='free_y')
```

The resulting plot is shown in [Figure 7-5](#) and reveals the nature of each cluster. For example, clusters 1 and 2 correspond to days on which the market is down and up, respectively. Clusters 3 and 5 are characterized by up-market days for consumer stocks and down-market days for energy stocks, respectively. Finally, cluster 4 captures the days in which energy stocks were up and consumer stocks were down.

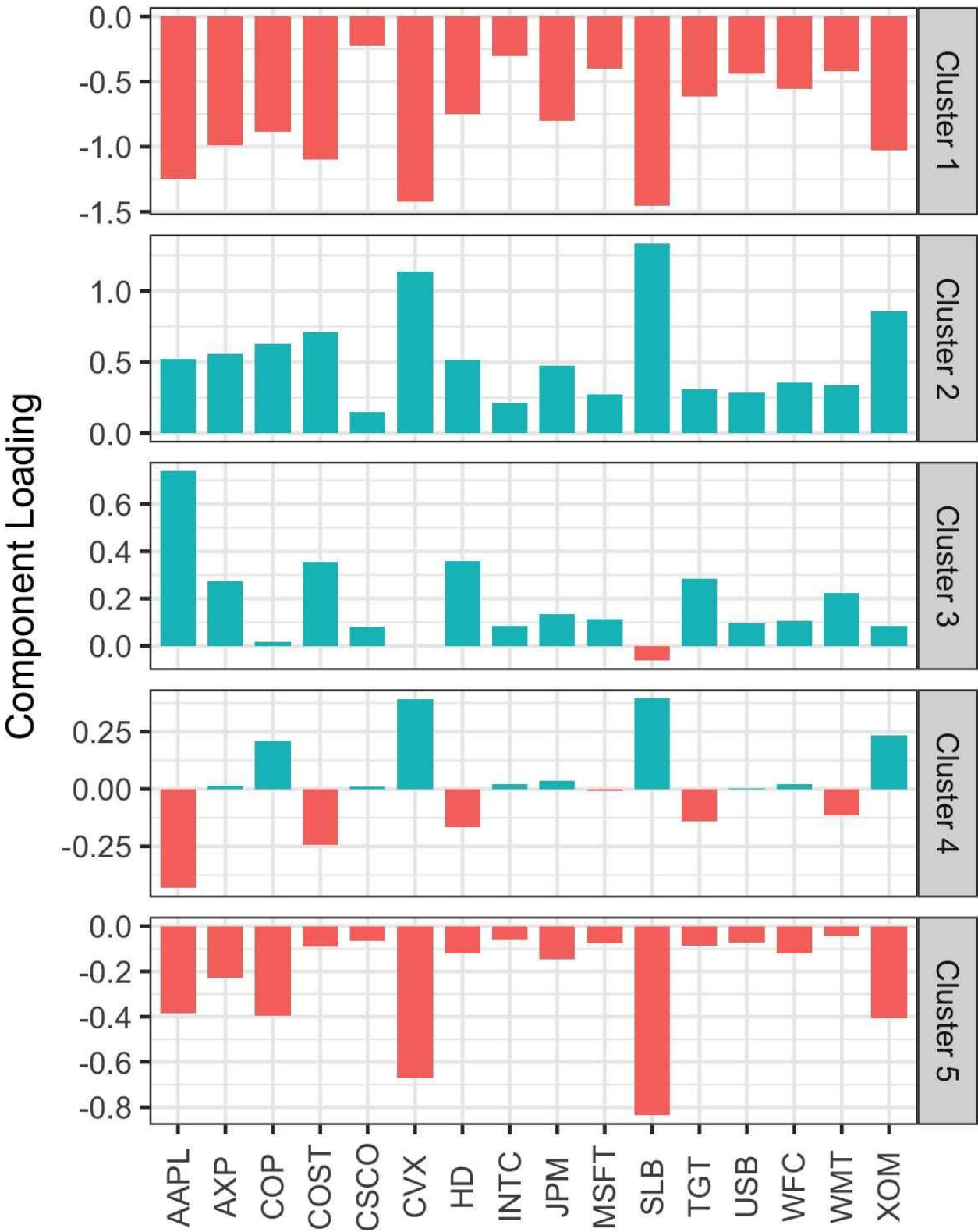


Figure 7-5. The means of the variables in each cluster ("cluster means")

CLUSTER ANALYSIS VERSUS PCA

The plot of cluster means is similar in spirit to looking at the loadings for principal component analysis (PCA); see “[Interpreting Principal Components](#)”. A major distinction is that unlike with PCA, the sign of the cluster means is meaningful. PCA identifies principal directions of variation, whereas cluster analysis finds groups of records located near one another.

Selecting the Number of Clusters

The K -means algorithm requires that you specify the number of clusters K . Sometimes the number of clusters is driven by the application. For example, a company managing a sales force might want to cluster customers into “personas” to focus and guide sales calls. In such a case, managerial considerations would dictate the number of desired customer segments — for example, two might not yield useful differentiation of customers, while eight might be too many to manage.

In the absence of a cluster number dictated by practical or managerial considerations, a statistical approach could be used. There is no single standard method to find the “best” number of clusters.

A common approach, called the *elbow method*, is to identify when the set of clusters explains “most” of the variance in the data. Adding new clusters beyond this set contributes relatively little incremental contribution in the variance explained. The elbow is the point where the cumulative variance explained flattens out after rising steeply, hence the name of the method.

Figure 7-6 shows the cumulative percent of variance explained for the default data for the number of clusters ranging from 2 to 15. Where is the elbow in this example? There is no obvious candidate, since the incremental increase in variance explained drops gradually. This is fairly typical in data that does not have well-defined clusters. This is perhaps a drawback of the elbow method, but it does reveal the nature of the data.

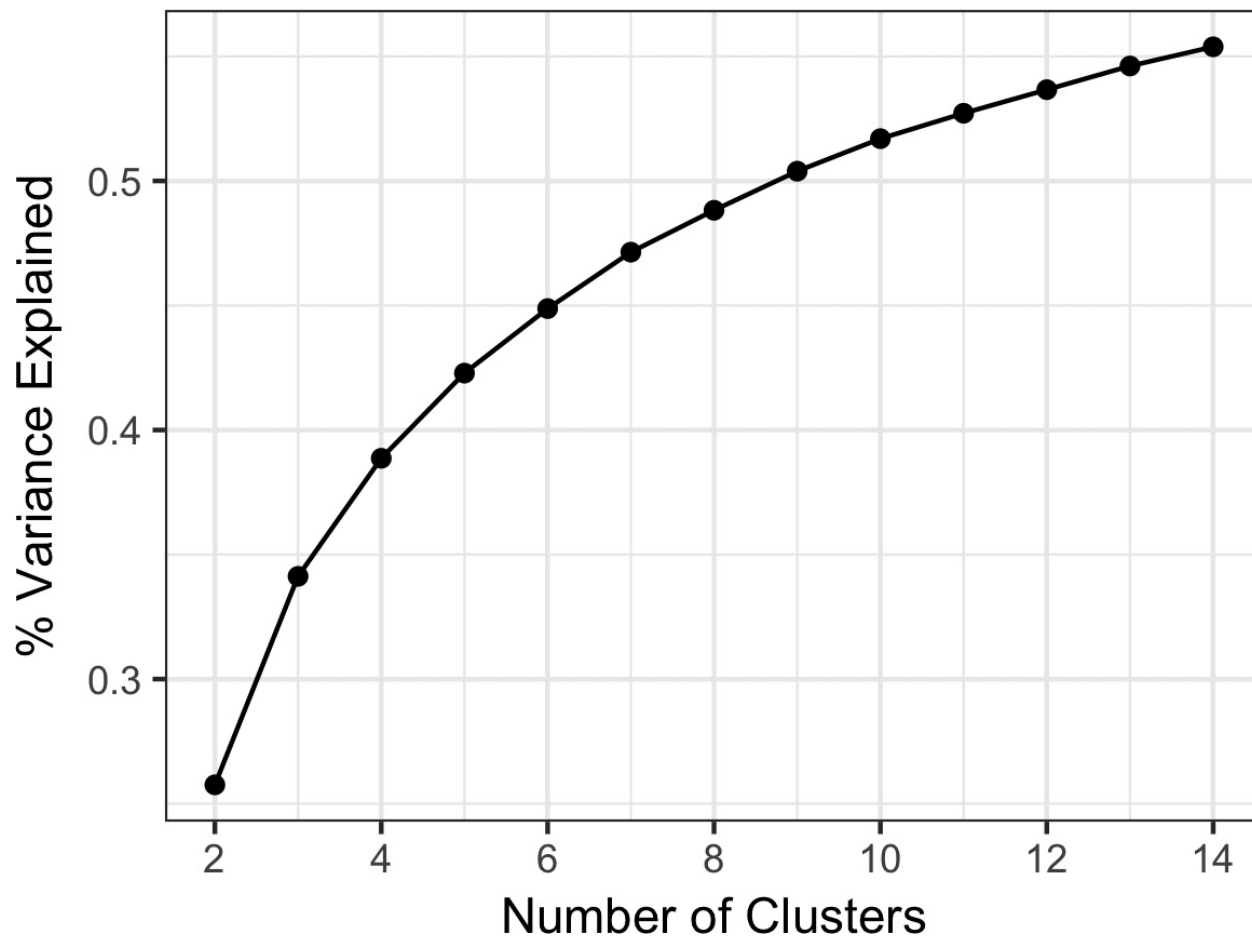


Figure 7-6. The elbow method applied to the stock data

In R, the `kmeans` function doesn't provide a single command for applying the elbow method, but it can be readily applied from the output of `kmeans` as shown here:

```
pct_var <- data.frame(pct_var = 0,
                      num_clusters=2:14)
totalss <- kmeans(df, centers=14, nstart=50, iter.max = 100)$totss
for(i in 2:14){
  pct_var[i-1, 'pct_var'] <- kmeans(df, centers=i, nstart=50, iter.max = 100)
    $betweenss/totalss
}
```

In evaluating how many clusters to retain, perhaps the most important test is this: how likely are the clusters to be replicated on new data? Are the clusters interpretable, and do they relate to a general characteristic of the data, or do they just reflect a specific instance? You can assess this, in part, using cross-

validation; see “Cross-Validation”.

In general, there is no single rule that will reliably guide how many clusters to produce.

NOTE

There are several more formal ways to determine the number of clusters based on statistical or information theory. For example, Robert Tibshirani, Guenther Walther, and Trevor Hastie (<http://www.stanford.edu/~hastie/Papers/gap.pdf>) propose a “gap” statistic based on statistical theory to identify the elbow. For most applications, a theoretical approach is probably not necessary, or even appropriate.

KEY IDEAS FOR K-MEANS CLUSTERING

- The number of desired clusters, K , is chosen by the user.
- The algorithm develops clusters by iteratively assigning records to the nearest cluster mean until cluster assignments do not change.
- Practical considerations usually dominate the choice of K ; there is no statistically determined optimal number of clusters.

Hierarchical Clustering

Hierarchical clustering is an alternative to *K*-means that can yield very different clusters. Hierarchical clustering is more flexible than *K*-means and more easily accommodates non-numerical variables. It is more sensitive in discovering outlying or aberrant groups or records. Hierarchical clustering also lends itself to an intuitive graphical display, leading to easier interpretation of the clusters.

KEY TERMS FOR HIERARCHICAL CLUSTERING

Dendrogram

A visual representation of the records and the hierarchy of clusters to which they belong.

Distance

A measure of how close one *record* is to another.

Dissimilarity

A measure of how close one *cluster* is to another.

Hierarchical clustering's flexibility comes with a cost, and hierarchical clustering does not scale well to large data sets with millions of records. For even modest-sized data with just tens of thousands of records, hierarchical clustering can require intensive computing resources. Indeed, most of the applications of hierarchical clustering are focused on relatively small data sets.

A Simple Example

Hierarchical clustering works on a data set with n records and p variables and is based on two basic building blocks:

- A distance metric $d_{i,j}$ to measure the distance between two records i and j .
- A dissimilarity metric $D_{A,B}$ to measure the difference between two clusters A and B based on the distances $d_{i,j}$ between the members of each cluster.

For applications involving numeric data, the most importance choice is the dissimilarity metric. Hierarchical clustering starts by setting each record as its own cluster and iterates to combine the least dissimilar clusters.

In R, the `hclust` function can be used to perform hierarchical clustering. One big difference with `hclust` versus `kmeans` is that it operates on the pairwise distances $d_{i,j}$ rather than the data itself. You can compute these using the `dist` function. For example, the following applies hierarchical clustering to the stock returns for a set of companies:

```
syms1 <- c('GOOGL', 'AMZN', 'AAPL', 'MSFT', 'CSCO', 'INTC', 'CVX',
           'XOM', 'SLB', 'COP', 'JPM', 'WFC', 'USB', 'AXP',
           'WMT', 'TGT', 'HD', 'COST')
# take transpose: to cluster companies, we need the stocks along the rows
df <- t(sp500_px[row.names(sp500_px) >= '2011-01-01', syms1])
d <- dist(df)
hcl <- hclust(d)
```

Clustering algorithms will cluster the records (rows) of a data frame. Since we want to cluster the companies, we need to *transpose* the data frame and put the stocks along the rows and the dates along the columns.

The Dendrogram

Hierarchical clustering lends itself to a natural graphical display as a tree, referred to as a *dendrogram*. The name comes from the Greek words *dendro* (tree) and *gramma* (drawing). In R, you can easily produce this using the `plot` command:

```
plot(hcl)
```

The result is shown in **Figure 7-7**. The leaves of the tree correspond to the records. The length of the branch in the tree indicates the degree of dissimilarity between corresponding clusters. The returns for Google and Amazon are quite dissimilar to the returns for the other stocks. The other stocks fall into natural groups: energy stocks, financial stocks, and consumer stocks are all separated into their own subtrees.

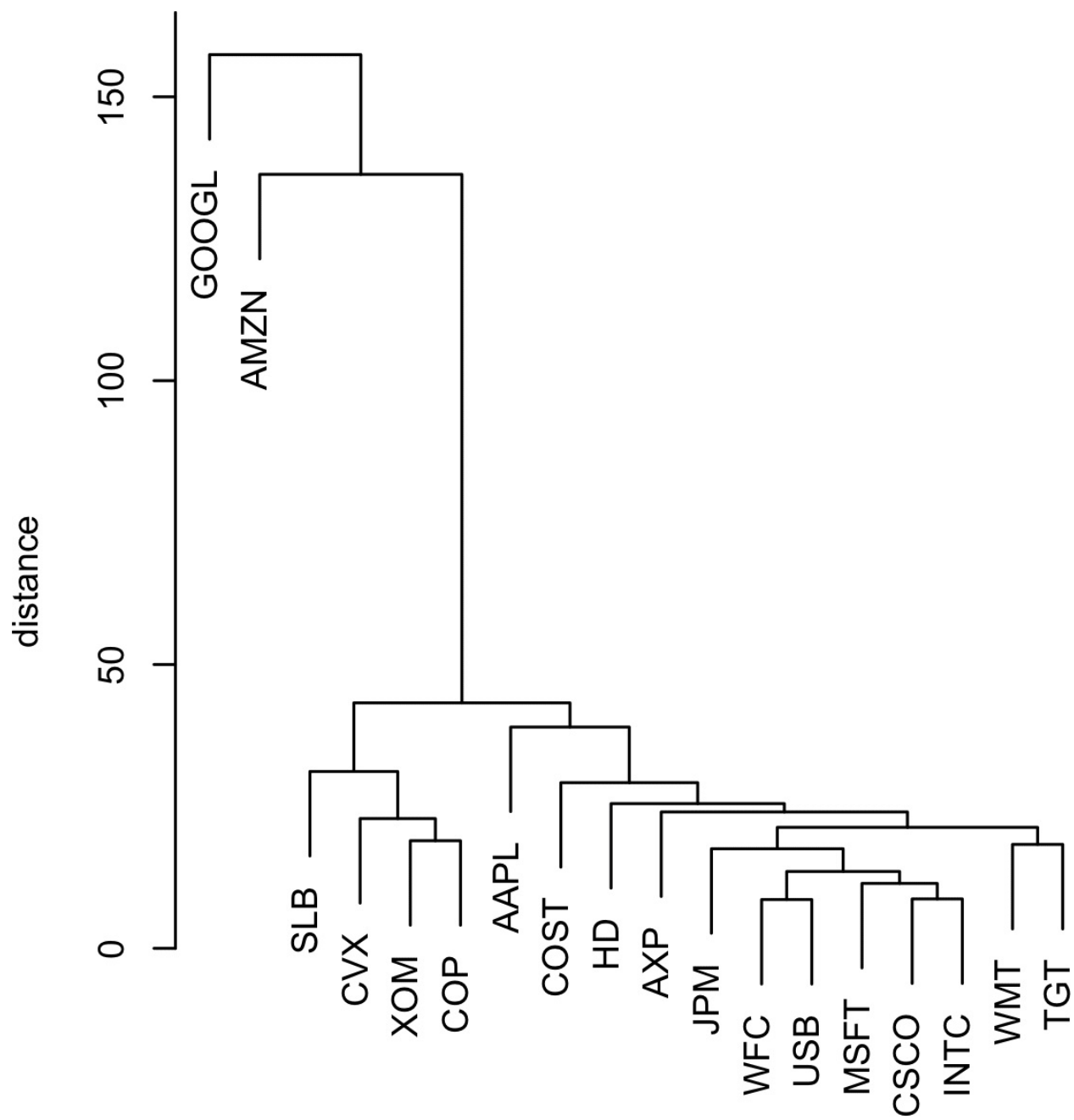


Figure 7-7. A dendrogram of stocks

In contrast to *K*-means, it is not necessary to prespecify the number of clusters. To extract a specific number of clusters, you can use the `cutree` function:

```
cutree(hc1, k=4)
```

GOOGL	AMZN	AAPL	MSFT	CSCO	INTC	CVX	XOM	SLB	COP	JPM	WFC
1	2	3	3	3	3	4	4	4	4	3	3
USB	AXP	WMT	TGT	HD	COST						
3	3	3	3	3	3						

The number of clusters to extract is set to 4, and you can see that Google and Amazon each belong to their own cluster. The oil stocks (XOM, CVS, SLB, COP) all belong to another cluster. The remaining stocks are in the fourth cluster.

The Agglomerative Algorithm

The main algorithm for hierarchical clustering is the *agglomerative* algorithm, which iteratively merges similar clusters. The agglomerative algorithm begins with each record constituting its own single-record cluster, then builds up larger and larger clusters. The first step is to calculate distances between all pairs of records.

For each pair of records (x_1, x_2, \dots, x_p) and (y_1, y_2, \dots, y_p) , we measure the distance between the two records, $d_{x,y}$, using a distance metric (see “Distance Metrics”). For example, we can use Euclidian distance:

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2}$$

We now turn to inter-cluster distance. Consider two clusters A and B , each with a distinctive set of records, $A = (a_1, a_2, \dots, a_m)$ and $B = (b_1, b_2, \dots, b_q)$.

We can measure the dissimilarity between the clusters $D(A, B)$ by using the distances between the members of A and the members of B .

One measure of dissimilarity is the *complete-linkage* method, which is the maximum distance across all pairs of records between A and B :

$$D(A, B) = \max d(a_i, b_j) \text{ for all pairs } i, j$$

This defines the dissimilarity as the biggest difference between all pairs.

The main steps of the agglomerative algorithm are:

1. Create an initial set of clusters with each cluster consisting of a single record for all records in the data.
2. Compute the dissimilarity $D(C_k, C_\ell)$ between all pairs of clusters k, ℓ .
3. Merge the two clusters C_k and C_ℓ that are least dissimilar as measured by $D(C_k, C_\ell)$.

4. If we have more than one cluster remaining, return to step 2. Otherwise, we are done.

Measures of Dissimilarity

There are four common measures of dissimilarity: *complete linkage*, *single linkage*, *average linkage*, and *minimum variance*. These (plus other measures) are all supported by most hierarchical clustering software, including `hclust`. The complete linkage method defined earlier tends to produce clusters with members that are similar. The single linkage method is the minimum distance between the records in two clusters:

$$D(A, B) = \min d(a_i, b_j) \text{ for all pairs } i, j$$

This is a “greedy” method and produces clusters that can contain quite disparate elements. The average linkage method is the average of all distance pairs and represents a compromise between the single and complete linkage methods. Finally, the minimum variance method, also referred to as *Ward’s method*, is similar to *K-means* since it minimizes the within-cluster sum of squares (see “**K-Means Clustering**”).

Figure 7-8 applies hierarchical clustering using the four measures to the ExxonMobil and Chevron stock returns. For each measure, four clusters are retained.

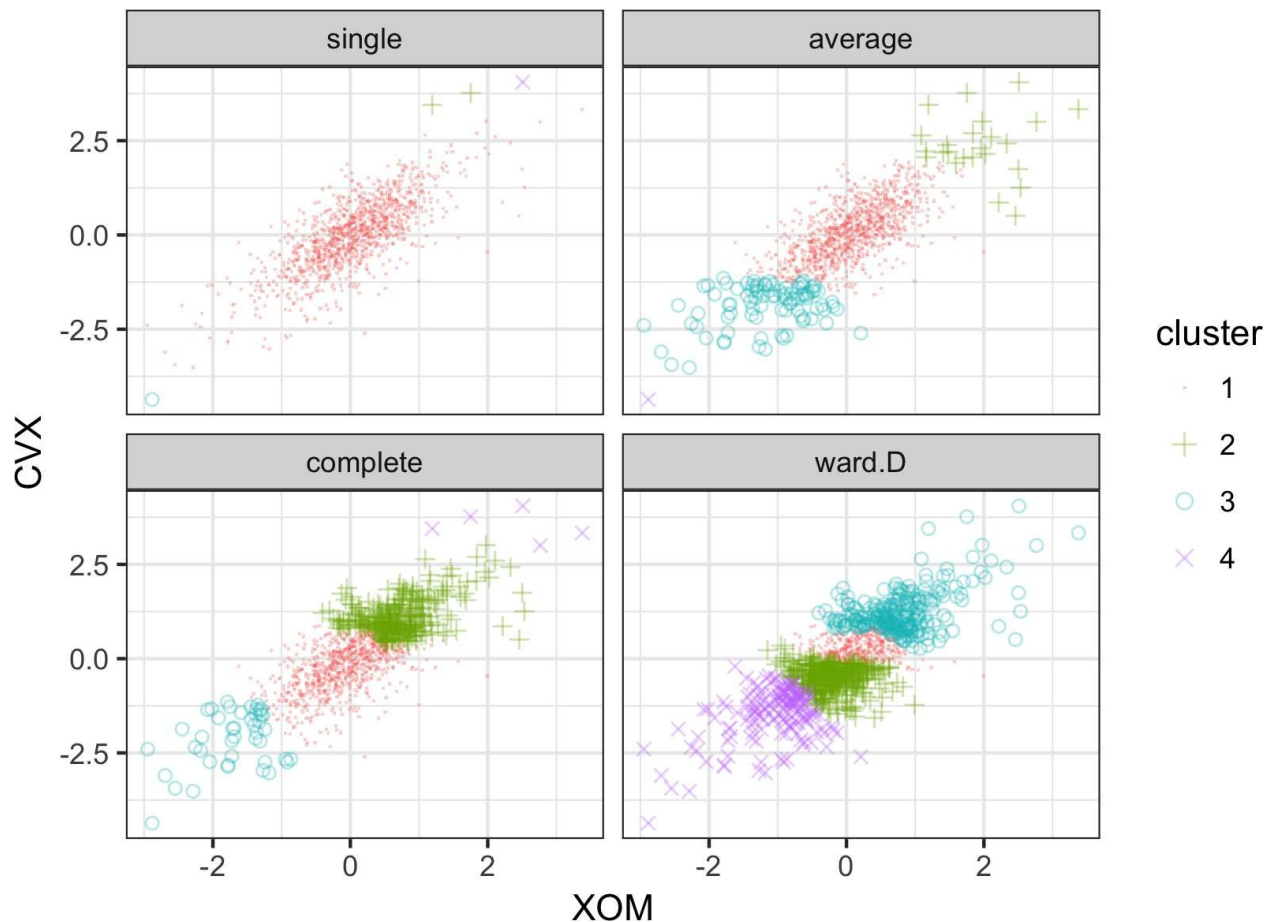


Figure 7-8. A comparison of measures of dissimilarity applied to stock data

The results are strikingly different: the single linkage measure assigns almost all of the points to a single cluster. Except for the minimum variance method (ward.D), all measures end up with at least one cluster with just a few outlying points. The minimum variance method is most similar to the *K*-means cluster; compare with [Figure 7-4](#).

KEY IDEAS FOR HIERARCHICAL CLUSTERING

- Start with every record in its own cluster.
- Progressively, clusters are joined to nearby clusters until all records belong to a single cluster (the agglomerative algorithm).
- The agglomeration history is retained and plotted, and the user (without specifying the number of clusters beforehand) can visualize the number and structure of clusters at different stages.
- Inter-cluster distances are computed in different ways, all relying on the set of all inter-record distances.



Model-Based Clustering

Clustering methods such as hierarchical clustering and K -means are based on heuristics and rely primarily on finding clusters whose members are close to one another, as measured directly with the data (no probability model involved). In the past 20 years, significant effort has been devoted to developing *model-based clustering* methods. Adrian Raftery and other researchers at the University of Washington made critical contributions to model-based clustering, including both theory and software. The techniques are grounded in statistical theory and provide more rigorous ways to determine the nature and number of clusters. They could be used, for example, in cases where there might be one group of records that are similar to one another but not necessarily close to one another (e.g., tech stocks with high variance of returns), and another group of records that is similar, and also close (e.g., utility stocks with low variance).

Multivariate Normal Distribution

The most widely used model-based clustering methods rest on the *multivariate normal* distribution. The multivariate normal distribution is a generalization of the normal distribution to set of p variables X_1, X_2, \dots, X_p . The distribution is defined by a set of means $\mu = \mu_1, \mu_2, \dots, \mu_p$ and a covariance matrix Σ . The covariance matrix is a measure of how the variables correlate with each other (see “**Covariance Matrix**” for details on the covariance). The covariance matrix Σ consists of p variances $\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2$ and covariances $\sigma_{i,j}$ for all pairs of variables $i \neq j$. With the variables put along the rows and duplicated along the columns, the matrix looks like this:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} & \dots & \sigma_{1,p} \\ \sigma_{2,1} & \sigma_2^2 & \dots & \sigma_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p,1} & \sigma_{p,2} & \dots & \sigma_p^2 \end{bmatrix}$$

Since a covariance matrix is symmetric, and $\sigma_{i,j} = \sigma_{j,i}$, there are only $p \times (p - 1) - p$ covariance terms. In total, the covariance matrix has $p \times (p - 1)$ parameters. The distribution is denoted by:

$$(X_1, X_2, \dots, X_p) \widetilde{N}_p(\mu, \Sigma)$$

This is a symbolic way of saying that the variables are all normally distributed, and the overall distribution is fully described by the vector of variable means and the covariance matrix.

Figure 7-9 shows the probability contours for a multivariate normal distribution for two variables X and Y (the 0.5 probability contour, for example, contains 50% of the distribution).

The means are $\mu_x = 0.5$ and $\mu_y = -0.5$ and the covariance matrix is:

$$\Sigma = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

Since the covariance σ_{xy} is positive, X and Y are positively correlated.

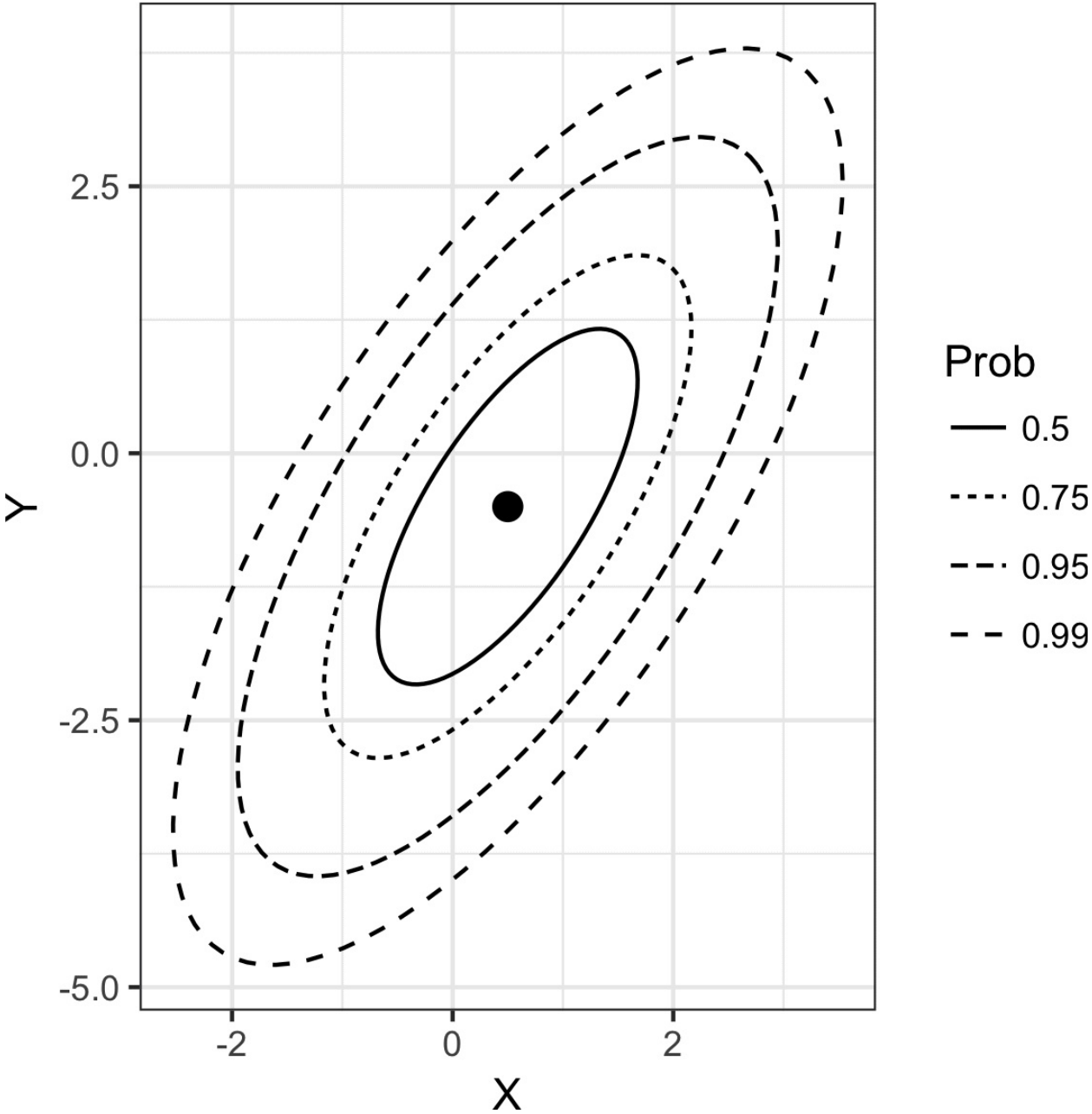


Figure 7-9. Probability contours for a two-dimensional normal distribution

Mixtures of Normals

The key idea behind model-based clustering is that each record is assumed to be distributed as one of K multivariate-normal distributions, where K is the number of clusters. Each distribution has a different mean μ and covariance matrix Σ .

For example, if you have two variables, X and Y , then each row (X_i, Y_i) is modeled as having been sampled from one of K distributions $N_1(\mu_1, \Sigma_1), N_1(\mu_2, \Sigma_2), \dots, N_1(\mu_K, \Sigma_K)$.

R has a very rich package for model-based clustering called `mclust`, originally developed by Chris Fraley and Adrian Raftery. With this package, we can apply model-based clustering to the stock return data we previously analyzed using K -means and hierarchical clustering:

```
> library(mclust)
> df <- sp500_px[row.names(sp500_px)>='2011-01-01', c('XOM', 'CVX')]
> mcl <- Mclust(df)
> summary(mcl)
Mclust VEE (ellipsoidal, equal shape and orientation) model with 2 components:

log.likelihood    n df      BIC      ICL
      -2255.134 1131  9 -4573.546 -5076.856

Clustering table:
  1  2
963 168
```

If you execute this code, you will notice that the computation takes significantly longer than other procedures. Extracting the cluster assignments using the `predict` function, we can visualize the clusters:

```
cluster <- factor(predict(mcl)$classification)
ggplot(data=df, aes(x=XOM, y=CVX, color=cluster, shape=cluster)) +
  geom_point(alpha=.8)
```

The resulting plot is shown in [Figure 7-10](#). There are two clusters: one cluster in the middle of the data, and a second cluster in the outer edge of the data. This is very different from the clusters obtained using K -means ([Figure 7-4](#)) and hierarchical clustering ([Figure 7-8](#)), which find clusters that are compact.

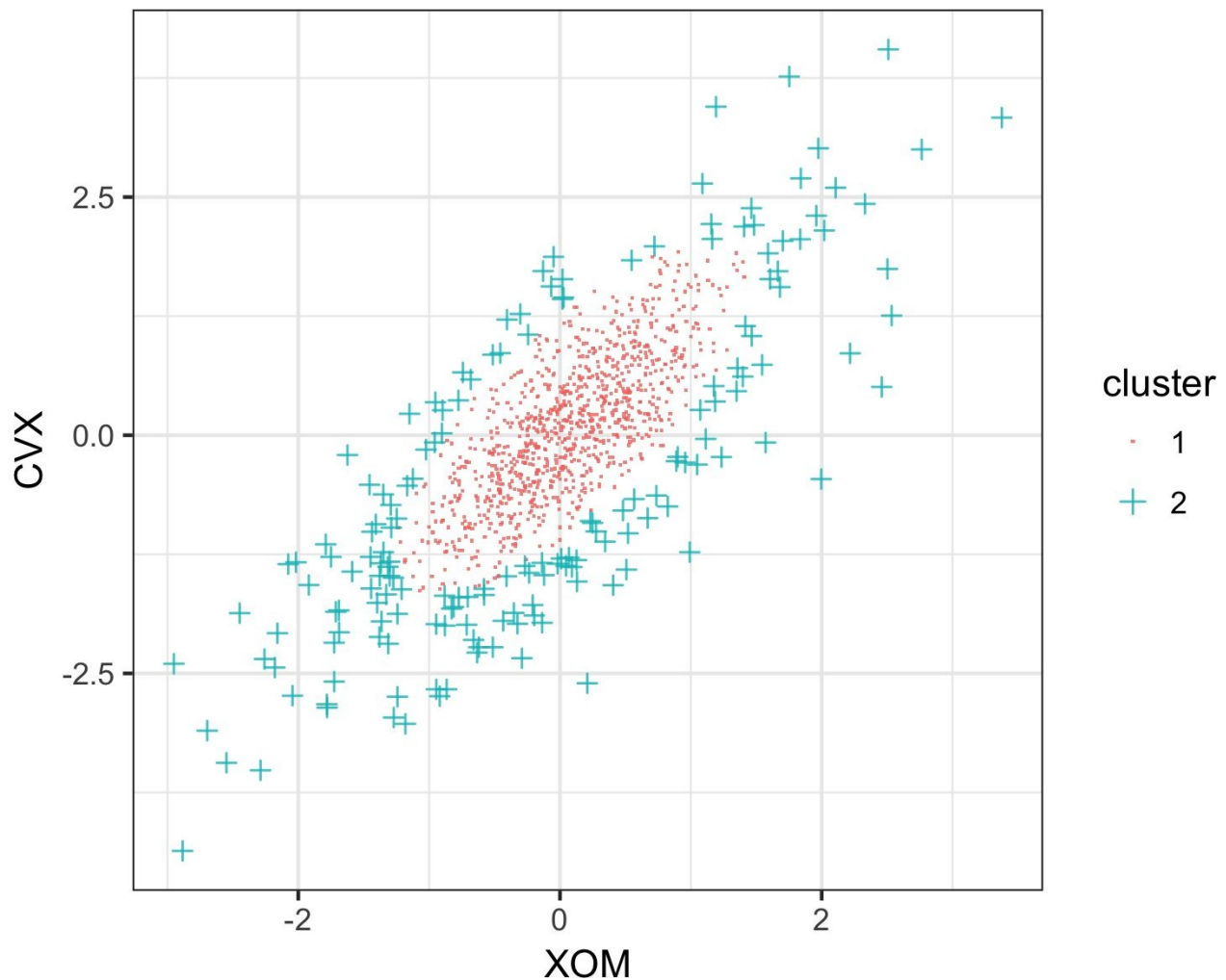


Figure 7-10. Two clusters are obtained for stock return data using *mclust*

You can extract the parameters to the normal distributions using the `summary` function:

```
> summary(mcl, parameters=TRUE)$mean
      [,1]      [,2]
XOM 0.05783847 -0.04374944
CVX 0.07363239 -0.21175715
> summary(mcl, parameters=TRUE)$variance
, , 1
      XOM      CVX
XOM 0.3002049 0.3060989
CVX 0.3060989 0.5496727
, , 2
      XOM      CVX
XOM 1.046318 1.066860
CVX 1.066860 1.915799
```

The distributions have similar means and correlations, but the second distribution has much larger variances and covariances.

The clusters from `mclust` may seem surprising, but in fact, they illustrate the statistical nature of the method. The goal of model-based clustering is to find the best-fitting set of multivariate normal distributions. The stock data appears to have a normal-looking shape: see the contours of [Figure 7-9](#). In fact, though, stock returns have a longer-tailed distribution than a normal distribution. To handle this, `mclust` fits a distribution to the bulk of the data, but then fits a second distribution with a bigger variance.

Selecting the Number of Clusters

Unlike K -means and hierarchical clustering, `mclust` automatically selects the number of clusters (in this case, two). It does this by choosing the number of clusters for which the *Bayesian Information Criteria* (*BIC*) has the largest value. BIC (similar to AIC) is a general tool to find the best model amongst a candidate set of models. For example, AIC (or BIC) is commonly used to select a model in stepwise regression; see “[Model Selection and Stepwise Regression](#)”. BIC works by selecting the best-fitting model with a penalty for the number of parameters in the model. In the case of model-based clustering, adding more clusters will always improve the fit at the expense of introducing additional parameters in the model.

You can plot the BIC value for each cluster size using a function in `hclust`:

```
plot(mcl, what='BIC', ask=FALSE)
```

The number of clusters — or number of different multivariate normal models (components) — is shown on the x-axis (see [Figure 7-11](#)).

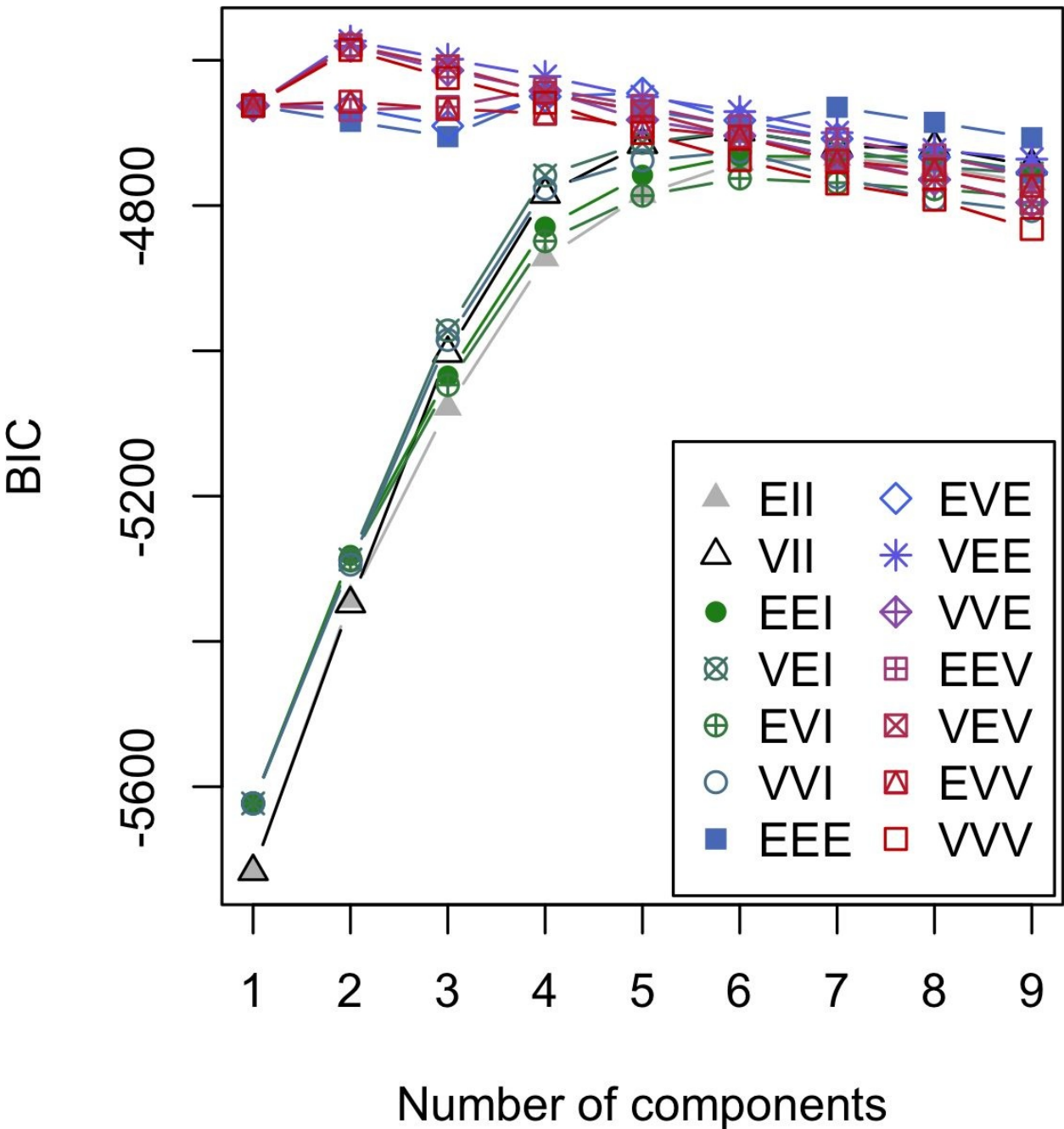


Figure 7-11. BIC scores for the stock return data for different numbers of clusters (components)

This plot is similar to the elbow plot used to identify the number of clusters to choose for *K*-means, except the value being plotted is BIC instead of percent of variance explained (see Figure 7-6). One big difference is that instead of one line, *mc1ust* shows 14 different lines! This is because *mc1ust* is actually fitting 14 different models for each cluster size, and ultimately it chooses the best-fitting

model.

Why does `mclust` fit so many models to determine the best set of multivariate normals? It's because there are different ways to parameterize the covariance matrix Σ for fitting a model. For the most part, you do not need to worry about the details of the models and can simply use the model chosen by `mclust`. In this example, according to BIC, three different models (called VEE, VEV, and VVE) give the best fit using two components.

NOTE

Model-based clustering is a rich and rapidly developing area of study, and the coverage in this text only spans a small part of the field. Indeed, the `mclust` help file is currently 154 pages long. Navigating the nuances of model-based clustering is probably more effort than is needed for most problems encountered by data scientists.

Model-based clustering techniques do have some limitations. The methods require an underlying assumption of a model for the data, and the cluster results are very dependent on that assumption. The computations requirements are higher than even hierarchical clustering, making it difficult to scale to large data. Finally, the algorithm is more sophisticated and less accessible than that of other methods.

KEY IDEAS FOR MODEL-BASED CLUSTERING

- Clusters are assumed to derive from different data-generating processes with different probability distributions.
- Different models are fit, assuming different numbers of (typically normal) distributions.
- The method chooses the model (and the associated number of clusters) that fits the data well without using too many parameters (i.e., overfitting).

Further Reading

For more detail on model-based clustering, see the [mclust documentation](#).

Scaling and Categorical Variables

Unsupervised learning techniques generally require that the data be appropriately scaled. This is different from many of the techniques for regression and classification in which scaling is not important (an exception is *K*-nearest neighbors; see “*K-Nearest Neighbors*”).

KEY TERMS FOR SCALING DATA

Scaling

Squashing or expanding data, usually to bring multiple variables to the same scale.

Normalization

One method of scaling — subtracting the mean and dividing by the standard deviation.

Synonym

Standardization

Gower’s distance

A scaling algorithm applied to mixed numeric and categorical data to bring all variables to a 0–1 range.

For example, with the personal loan data, the variables have widely different units and magnitude. Some variables have relatively small values (e.g., number of years employed), while others have very large values (e.g., loan amount in dollars). If the data is not scaled, then the PCA, *K*-means, and other clustering methods will be dominated by the variables with large values and ignore the variables with small values.

Categorical data can pose a special problem for some clustering procedures. As with *K*-nearest neighbors, unordered factor variables are generally converted to a set of binary (0/1) variables using one hot encoding (see “*One Hot Encoder*”). Not only are the binary variables likely on a different scale from other data, the fact that binary variables have only two values can prove problematic with techniques such as PCA and *K*-means.

Scaling the Variables

Variables with very different scale and units need to be normalized appropriately before you apply a clustering procedure. For example, let's apply kmeans to a set of data of loan defaults without normalizing:

```
df <- defaults[, c('loan_amnt', 'annual_inc', 'revol_bal', 'open_acc',
                   'dti', 'revol_util')]
km <- kmeans(df, centers=4, nstart=10)
centers <- data.frame(size=km$size, km$centers)
round(centers, digits=2)
  size loan_amnt annual_inc revol_bal open_acc dti revol_util
1   55  23157.27  491522.49  83471.07   13.35  6.89    58.74
2  1218  21900.96  165748.53  38299.44   12.58 13.43    63.58
3  7686  18311.55   83504.68  19685.28   11.68 16.80    62.18
4 14177  10610.43   42539.36  10277.97    9.60 17.73    58.05
```

The variables `annual_inc` and `revol_bal` dominate the clusters, and the clusters have very different sizes. Cluster 1 has only 55 members with comparatively high income and revolving credit balance.

A common approach to scaling the variables is to convert them to z-scores by subtracting the mean and dividing by the standard deviation. This is termed standardization or normalization (see “**Standardization (Normalization, Z-Scores)**” for more discussion about using z-scores):

$$z = \frac{x - \bar{x}}{s}$$

See what happens to the clusters when kmeans is applied to the normalized data:

```
df0 <- scale(df)
km0 <- kmeans(df0, centers=4, nstart=10)
centers0 <- scale(km0$centers, center=FALSE,
                  scale=1/attr(df0, 'scaled:scale'))
centers0 <- scale(centers0, center=-attr(df0, 'scaled:center'), scale=F)
data.frame(size=km0$size, centers0)
  size loan_amnt annual_inc revol_bal open_acc dti revol_util
1  5429  10393.60   53689.54   6077.77    8.69 11.35    30.69
2  6396  13310.43   55522.76  16310.95   14.25 24.27    59.57
3  7493  10482.19   51216.95  11530.17    7.48 15.79    77.68
4  3818  25933.01  116144.63  32617.81   12.44 16.25    66.01
```

The cluster sizes are more balanced, and the clusters are not just dominated by

`annual_inc` and `revol_bal`, revealing more interesting structure in the data. Note that the centers are rescaled to the original units in the preceding code. If we had left them unscaled, the resulting values would be in terms of z-scores, and therefore less interpretable.

NOTE

Scaling is also important for PCA. Using the z-scores is equivalent to using the correlation matrix (see “**Correlation**”) instead of the covariance matrix in computing the principal components. Software to compute PCA usually has an option to use the correlation matrix (in R, the `princomp` function has the argument `cor`).

Dominant Variables

Even in cases where the variables are measured on the same scale and accurately reflect relative importance (e.g., movement to stock prices), it can sometimes be useful to rescale the variables.

Suppose we add Alphabet (GOOGL) and Amazon (AMZN) to the analysis in “Interpreting Principal Components”.

```
syms <- c('AMZN', 'GOOGL', 'AAPL', 'MSFT', 'CSCO', 'INTC', 'CVX', 'XOM',
          'SLB', 'COP', 'JPM', 'WFC', 'USB', 'AXP', 'WMT', 'TGT', 'HD', 'COST')
top_sp1 <- sp500_px[row.names(sp500_px) >= '2005-01-01', syms]
sp_pca1 <- princomp(top_sp1)
screepplot(sp_pca1)
```

The screeplot displays the variances for the top principal components. In this case, the screeplot in [Figure 7-12](#) reveals that the variances of the first and second components are much larger than the others. This often indicates that one or two variables dominate the loadings. This is, indeed, the case in this example:

```
round(sp_pca1$loadings[,1:2], 3)
      Comp.1 Comp.2
GOOGL  0.781  0.609
AMZN   0.593 -0.792
AAPL   0.078  0.004
MSFT   0.029  0.002
CSCO   0.017 -0.001
INTC   0.020 -0.001
CVX    0.068 -0.021
XOM    0.053 -0.005
...
```

The first two principal components are almost completely dominated by GOOGL and AMZN. This is because the stock price movements of GOOGL and AMZN dominate the variability.

To handle this situation, you can either include them as is, rescale the variables (see “[Scaling the Variables](#)”), or exclude the dominant variables from the analysis and handle them separately. There is no “correct” approach, and the treatment depends on the application.

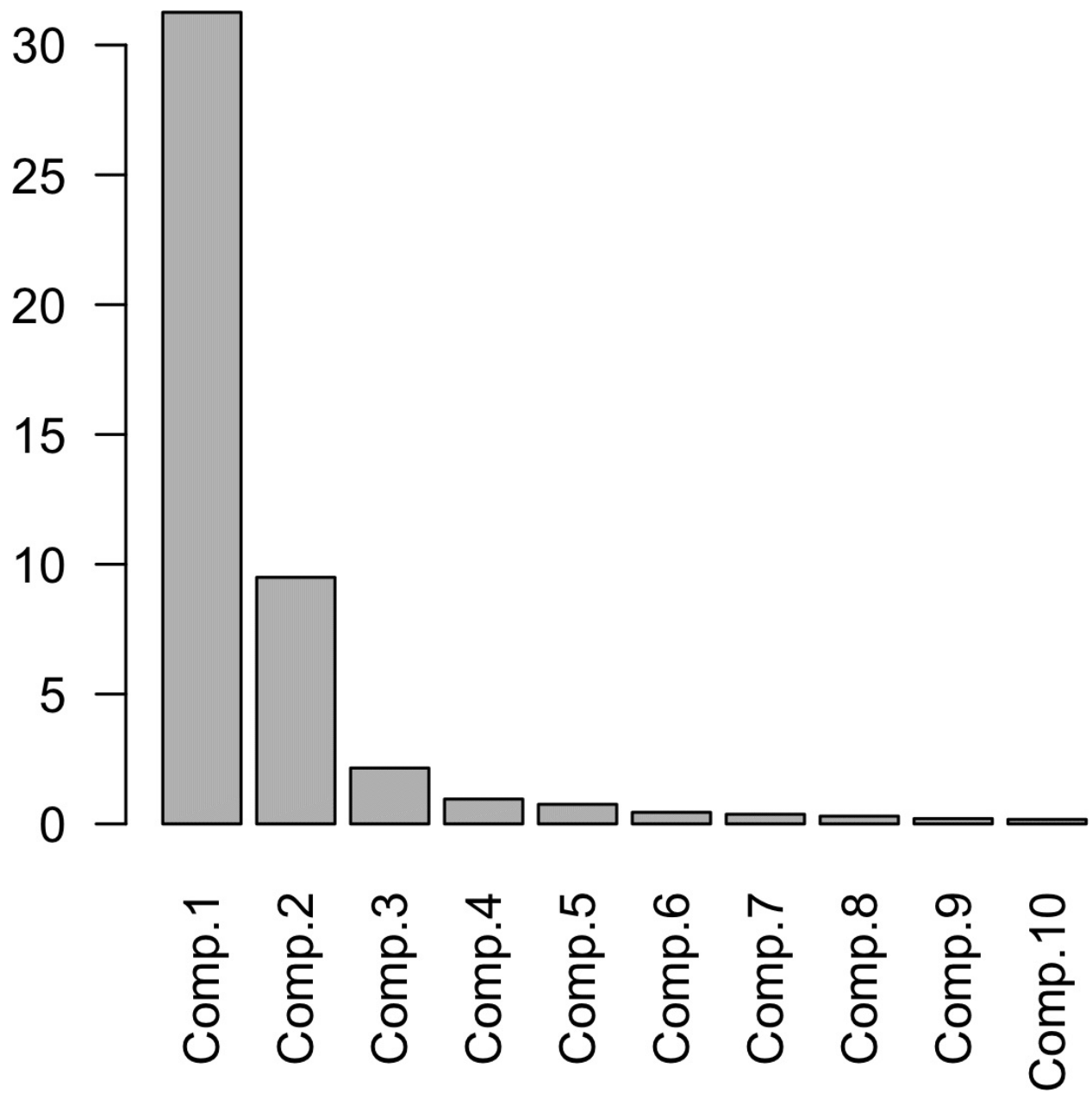


Figure 7-12. A screeplot for a PCA of top stocks from the S&P 500 including GOOGL and AMZN

Categorical Data and Gower's Distance

In the case of categorical data, you must convert it to numeric data, either by ranking (for an ordered factor) or by encoding as a set of binary (dummy) variables. If the data consists of mixed continuous and binary variables, you will usually want to scale the variables so that the ranges are similar; see “[Scaling the Variables](#)”. One popular method is to use *Gower's distance*.

The basic idea behind Gower's distance is to apply a different distance metric to each variable depending on the type of data:

- For numeric variables and ordered factors, distance is calculated as the absolute value of the difference between two records (*Manhattan distance*).
- For categorical variables, the distance is 1 if the categories between two records are different and the distance is 0 if the categories are the same.

Gower's distance is computed as follows:

1. Compute the distance $d_{i,j}$ for all pairs of variables i and j for each record.
2. Scale each pair $d_{i,j}$ so the minimum is 0 and the maximum is 1.
3. Add the pairwise scaled distances between variables together, either using a simple or weighted mean, to create the distance matrix.

To illustrate Gower's distance, take a few rows from the loan data:

```
> x = defaults[1:5, c('dti', 'payment_inc_ratio', 'home', 'purpose')]
> x
# A tibble: 5 × 4
  dti payment_inc_ratio home      purpose
<dbl>      <dbl> <fctr>      <fctr>
1  1.00      2.39320 RENT        car
2  5.55      4.57170 OWN         small_business
3 18.08      9.71600 RENT        other
4 10.08     12.21520 RENT debt_consolidation
5  7.06      3.90888 RENT        other
```

The function `daisy` in the `cluster` package can be used to compute Gower's distance:

```
> library(cluster)
> daisy(x, metric='gower')
Dissimilarities :
      1      2      3      4
2 0.6220479
3 0.6863877 0.8143398
4 0.6329040 0.7608561 0.4307083
5 0.3772789 0.5389727 0.3091088 0.5056250
```

All distances are between 0 and 1. The pair of records with the biggest distance is 2 and 3: neither has the same values for home or purpose and they have very different levels of `dti` (debt-to-income) and `payment_inc_ratio`. Records 3 and 5 have the smallest distance because they share the same values for home or purpose.

You can apply hierarchical clustering (see “[Hierarchical Clustering](#)”) to the resulting distance matrix using `hclust` to the output from `daisy`:

```
df <- defaults[sample(nrow(defaults), 250),
                  c('dti', 'payment_inc_ratio', 'home', 'purpose')]
d = daisy(df, metric='gower')
hcl <- hclust(d)
dnd <- as.dendrogram(hcl)
plot(dnd, leaflab='none')
```

The resulting dendrogram is shown in [Figure 7-13](#). The individual records are not distinguishable on the x-axis, but we can examine the records in one of the subtrees (on the left, using a “cut” of 0.5), with this code:

```
> df[labels(dnd_cut$lower[[1]]),]
# A tibble: 9 × 4
  dti payment_inc_ratio home purpose
<dbl>      <dbl> <fctr> <fctr>
1 24.57      0.83550 RENT other
2 34.95      5.02763 RENT other
3 1.51       2.97784 RENT other
4 8.73      14.42070 RENT other
5 12.05      9.96750 RENT other
6 10.15     11.43180 RENT other
7 19.61     14.04420 RENT other
8 20.92      6.90123 RENT other
9 22.49      9.36000 RENT other
```

This subtree entirely consists of renters with a loan purpose labeled as “other.” While strict separation is not true of all subtrees, this illustrates that the categorical variables tend to be grouped together in the clusters.

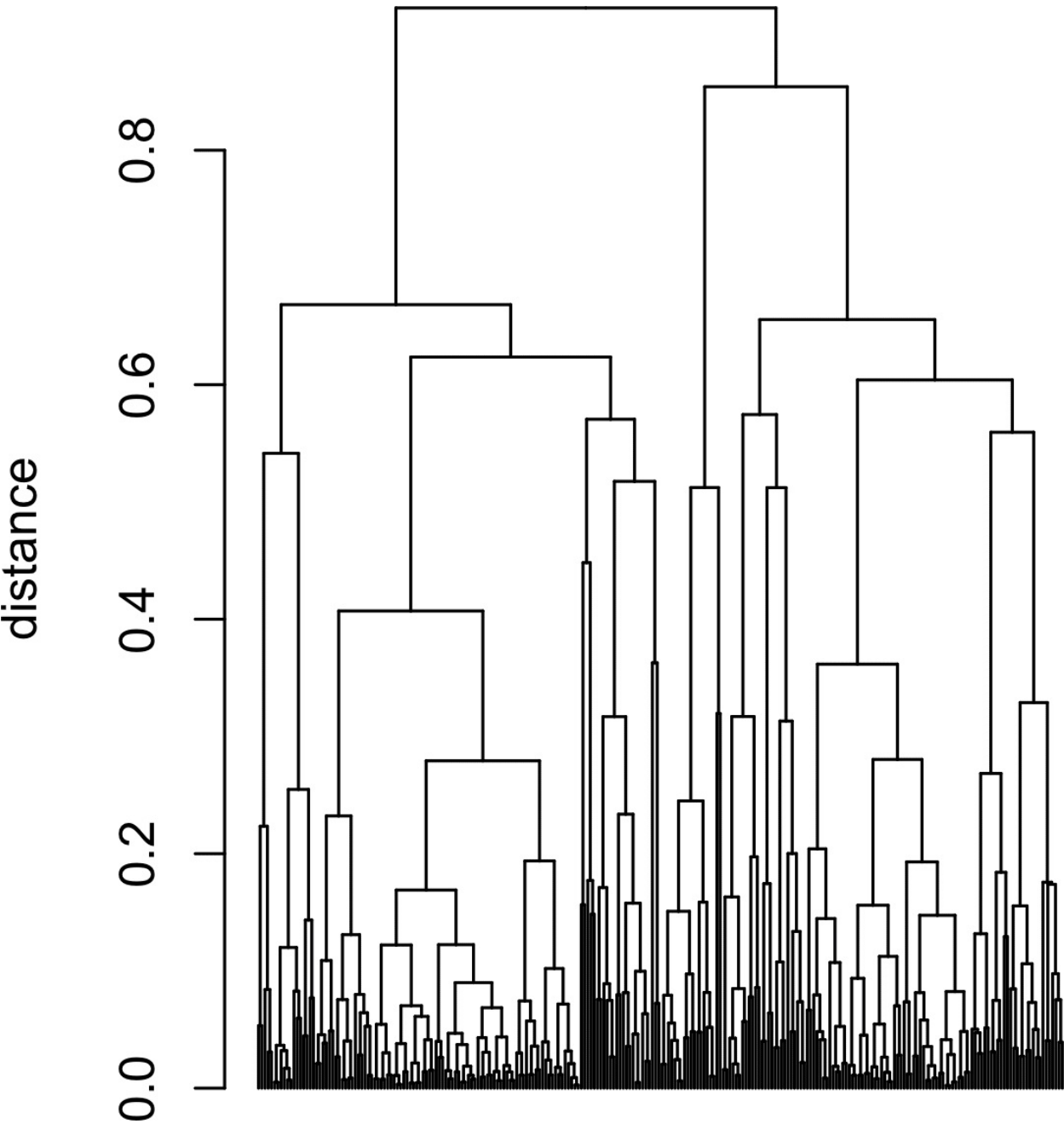


Figure 7-13. A dendrogram of hclust applied to a sample of loan default data with mixed variable types

Problems with Clustering Mixed Data

K-means and PCA are most appropriate for continuous variables. For smaller data sets, it is better to use hierarchical clustering with Gower's distance. In principle, there is no reason why K-means can't be applied to binary or categorical data.

You would usually use the “one hot encoder” representation (see “**One Hot Encoder**”) to convert the categorical data to numeric values. In practice, however, using K-means and PCA with binary data can be difficult.

If the standard z-scores are used, the binary variables will dominate the definition of the clusters. This is because 0/1 variables take on only two values and K-means can obtain a small within-cluster sum-of-squares by assigning all the records with a 0 or 1 to a single cluster. For example, apply kmeans to loan default data including factor variables home and pub_rec_zero:

```
df <- model.matrix(~ -1 + dti + payment_inc_ratio + home + pub_rec_zero,
                  data=defaults)
df0 <- scale(df)
km0 <- kmeans(df0, centers=4, nstart=10)
centers0 <- scale(km0$centers, center=FALSE,
                 scale=1/attr(df0, 'scaled:scale'))
scale(centers0, center=-attr(df0, 'scaled:center'), scale=F)
      dti payment_inc_ratio homeMORTGAGE homeOWN homeRENT pub_rec_zero
1 17.02           9.10         0.00      0         1.00         1.00
2 17.47           8.43         1.00      0         0.00         1.00
3 17.23           9.28         0.00      1         0.00         0.92
4 16.50           8.09         0.52      0         0.48         0.00
```

The top four clusters are essentially proxies for the different levels of the factor variables. To avoid this behavior, you could scale the binary variables to have a smaller variance than other variables. Alternatively, for very large data sets, you could apply clustering to different subsets of data taking on specific categorical values. For example, you could apply clustering separately to those loans made to someone who has a mortgage, owns a home outright, or rents.

KEY IDEAS FOR SCALING DATA

- Variables measured on different scales need to be transformed to similar scales, so that their impact on algorithms is not determined mainly by their scale.
- A common scaling method is normalization (standardization) — subtracting the mean and dividing by the standard deviation.
- Another method is Gower's distance, which scales all variables to the 0–1 range (it is often used

with mixed numeric and categorical data).

Summary

For dimension reduction of numeric data, the main tools are either principal components analysis or *K*-means clustering. Both require attention to proper scaling of the data to ensure meaningful data reduction.

For clustering with highly structured data in which the clusters are well separated, all methods will likely produce a similar result. Each method offers its own advantage. *K*-means scales to very large data and is easily understood.

Hierarchical clustering can be applied to mixed data types — numeric and categorical — and lends itself to an intuitive display (the dendrogram). Model-based clustering is founded on statistical theory and provides a more rigorous approach, as opposed to the heuristic methods. For very large data, however, *K*-means is the main method used.

With noisy data, such as the loan and stock data (and much of the data that a data scientist will face), the choice is more stark. *K*-means, hierarchical clustering, and especially model-based clustering all produce very different solutions. How should a data scientist proceed? Unfortunately, there is no simple rule of thumb to guide the choice. Ultimately, the method used will depend on the data size and the goal of the application.

¹ This and subsequent sections in this chapter © 2017 Datastats, LLC, Peter Bruce and Andrew Bruce, used by permission.

Bibliography

- [bokeh] Bokeh Development Team. “Bokeh: Python library for interactive visualization” (2014). <http://www.bokeh.pydata.org>.
- [Deng-Wickham-2011] Deng, H. and Wickham, H. “Density estimation in R” (2011). <http://vita.had.co.nz/papers/density-estimation.pdf>.
- [Wikipedia-2016] “Diving.” Wikipedia: The Free Encyclopedia. Wikimedia Foundation, Inc. 10 Mar 2016. Web. 19 Mar 2016.
- [Donoho-2015] Donoho, David. “50 Years of Data Science” (2015). <http://courses.csail.mit.edu/18.337/2015/docs/50YearsDataScience.pdf>.
- [Duong-2001] Duang, Tarn. “An introduction to kernel density estimation” (2001). <http://www.mvstat.net/tduong/research/seminars/seminar-2001-05.pdf>.
- [Few-2007] Few, Stephen. “Save the Pies for Dessert.” Visual Intelligence Newsletter, Perceptual Edge (2007). https://www.perceptualedge.com/articles/visual_business_intelligence/save
- [Hintze-Nelson-1998] Hintze, J. and Nelson, R. “Violin Plots: A Box Plot-Density Trace Synergism.” *The American Statistician* 52.2 (May 1998): 181–184.
- [Galton-1886] Galton, Francis. “Regression towards mediocrity in Hereditary stature.” *The Journal of the Anthropological Institute of Great Britain and Ireland*, 15:246-273. JSTOR 2841583.
- [ggplot2] Wickham, Hadley. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York (2009). ISBN: 978-0-387-98140-6. <http://had.co.nz/ggplot2/book>.
- [Hyndman-Fan-1996] Hyndman, R. J. and Fan, Y. “Sample quantiles in statistical packages,” *American Statistician* 50, (1996) 361–365.

- [lattice] Sarkar, Deepayan. *Lattice: Multivariate Data Visualization with R*. Springer (2008). ISBN 978-0-387-75968-5. <http://lmdvr.r-forge.r-project.org>.
- [Legendre] Legendre, Adrien-Marie. *Nouvelle methodes pour la determination des orbites des cometes*. F. Didot, Paris (1805).
- [NIST-Handbook-2012] NIST/SEMATECH *e-Handbook of Statistical Methods*, <http://www.itl.nist.gov/div898/handbook/eda/section3/eda35b.htm> (2012).
- [R-base-2015] R Core Team. “R: A Language and Environment for Statistical Computing,” R Foundation for Statistical Computing (2015). <http://www.R-project.org/>.
- [seaborn] Wasdom, Michael. “Seaborn: statistical data visualization” (2015). <http://stanford.edu/~mwaskom/software/seaborn/#>.
- [Stigler-Gauss] Stigler, Stephen M. “Gauss and the Invention of Least Squares.” *Ann. Stat.* 9(3), 465–474 (1981).
- [Trellis-Graphics] Becker, R., Cleveland, W, Shyu, M. and Kaluzny, S. “A Tour of Trellis Graphics” (1996). http://polisci.msu.edu/jacoby/icpsr/graphics/manuscripts/Trellis_tour.pdf.
- [Tukey-1962] Tukey, John W. “The Future of Data Analysis.” *Ann. Math. Statist.* 33 (1962), no. 1, 1–67. https://projecteuclid.org/download/pdf_1/euclid.aoms/1177704711
- [Tukey-1977] Tukey, John W. *Exploratory Data Analysis*. Pearson (1977). ISBN: 978-0-201-07616-5.
- [Tukey-1987] Tukey, John W. Edited by Jones, L. V. *The collected works of John W. Tukey: Philosophy and Principles of Data Analysis 1965–1986*, Volume IV. Chapman and Hall/CRC (1987). ISBN: 978-0-534-05101-3.
- [UCLA] “R Library: Contrast Coding Systems for Categorical Variables.” UCLA: Statistical Consulting Group. http://www.ats.ucla.edu/stat/r/library/contrast_coding.htm. Accessed June 2016.

- [Zhang-Wang-2007] Zhang, Qi and Wang, Wei. 19th International Conference on Scientific and Statistical Database Management, IEEE Computer Society (2007).

Index

A

A/B testing, **A/B Testing-For Further Reading**

control group, advantages of using, **Why Have a Control Group?**

epsilon-greedy algorithm, **Multi-Arm Bandit Algorithm**

importance of permissions, **Why Just A/B? Why Not C, D...?**

traditional, shortcoming of, **Multi-Arm Bandit Algorithm**

accuracy, **Evaluating Classification Models**

improving in random forests, **Random Forest**

Adaboost, **Boosting**

boosting algorithm, **The Boosting Algorithm**

adjusted R-squared, **Assessing the Model**

adjustment of p-values, **Multiple Testing, Multiple Testing**

agglomerative algorithm, **The Agglomerative Algorithm**

AIC (Akaike's Information Criteria), **Model Selection and Stepwise Regression, Selecting the Number of Clusters**

variants of, **Model Selection and Stepwise Regression**

Akaike, Hirotugu, **Model Selection and Stepwise Regression**

all subset regression, **Model Selection and Stepwise Regression**

alpha, **Statistical Significance and P-Values, Alpha**

dividing up in multiple testing, **Multiple Testing**

alternative hypothesis, **Hypothesis Tests, Alternative Hypothesis**

American Statistical Association (ASA), statement on p-values, **Value of the p-value**

anomaly detection, **Outliers, Regression and Prediction**

ANOVA (analysis of variance

statistical test based on F-statistic, **F-Statistic**

ANOVA (analysis of variance), **ANOVA-Further Reading**

computing ANOVA table in R, **F-Statistic**

decomposition of variance, **F-Statistic**

two-way, **Two-Way ANOVA**

arms (multi-arm bandits), **Multi-Arm Bandit Algorithm**

AUC (area under the ROC curve), **AUC**

average linkage, **Measures of Dissimilarity**

B

backward elimination, **Model Selection and Stepwise Regression**

backward selection, **Model Selection and Stepwise Regression**

bagging, **The Bootstrap, Resampling, Statistical Machine Learning, Bagging**

better predictive performance than single trees, **How Trees Are Used**

boosting vs., **Boosting**

using with random forests, **Random Forest**

bandit algorithms, **Multi-Arm Bandit Algorithm**

(see also multi-arm bandits)

bar charts, **Exploring Binary and Categorical Data**

Bayesian classification, **Naive Bayes**

(see also naive Bayes algorithm)

impracticality of exact Bayesian classification, **Why Exact Bayesian Classification Is Impractical**

Bayesian information criteria (BIC), **Model Selection and Stepwise Regression, Selecting the Number of Clusters**

beta distribution, **Multi-Arm Bandit Algorithm**

bias, **Bias**

selection bias, **Selection Bias-Further Reading**

bias-variance tradeoff, **Choosing K**

biased estimates, **Standard Deviation and Related Estimates**

from naive Bayes classifier, **The Naive Solution**

BIC (Bayesian information criteria), **Model Selection and Stepwise Regression, Selecting the Number of Clusters**

bidirectional alternative hypothesis, **One-Way, Two-Way Hypothesis Test**

big data

and outliers in regression, **Outliers**

use of regression in, **Prediction versus Explanation (Profiling)**

value of, **Size versus Quality: When Does Size Matter?**

binary data, **Elements of Structured Data**

exploring, **Exploring Binary and Categorical Data-Correlation**

binomial, **Binomial Distribution**

binomial distribution, **Binomial Distribution-Further Reading**

binomial trials, **Binomial Distribution**

bins

hexagonal binning, **Hexagonal Binning and Contours (Plotting Numeric versus Numeric Data)**

in frequency tables, **Frequency Table and Histograms**

in histograms, **Frequency Table and Histograms**

bivariate analysis, **Exploring Two or More Variables**

black swan theory, **Long-Tailed Distributions**

blind studies, **Why Have a Control Group?**

boosting, **Statistical Machine Learning, Tree Models, Boosting-Summary**

bagging vs., **Boosting**

boosting algorithm, **The Boosting Algorithm**

hyperparameters and cross-validation, **Hyperparameters and Cross-**

Validation

overfitting, avoiding using regularization, **Regularization: Avoiding Overfitting**

XGBoost, **XGBoost**

bootstrap, **The Bootstrap-Further Reading, Resampling**

confidence interval generation, **Confidence Intervals, Confidence and Prediction Intervals**

permutation tests, **Exhaustive and Bootstrap Permutation Test**

resampling vs. bootstrapping, **Resampling versus Bootstrapping**

using with random forests, **Random Forest**

bootstrap sample, **The Bootstrap**

boxplots, **Exploring the Data Distribution**

combining with a violin plot, example, **Categorical and Numeric Data**

example, percent of airline delays by carrier, **Categorical and Numeric Data**

outliers in, **Outliers**

percentiles and, **Percentiles and Boxplots**

Breiman, Leo, **Statistical Machine Learning**

bubble plots, **Influential Values**

C

categorical data, **Elements of Structured Data**

exploring, **Exploring Binary and Categorical Data-Correlation**

expected value, **Expected Value**

mode, **Mode**

numerical data as categorical data, **Exploring Binary and Categorical Data**

exploring two categorical variables, **Two Categorical Variables**

importance of the concept, **Elements of Structured Data**

numeric variable grouped by categorical variable, **Categorical and Numeric Data**

scaling and categorical variables, **Scaling and Categorical Variables-Summary**

dominant variables, **Dominant Variables**

Gower's distance, **Categorical Data and Gower's Distance**

scaling the variables, **Scaling the Variables**

categorical variables, **Factor Variables in Regression**

(see also factor variables)

causation, regression and, **Prediction versus Explanation (Profiling)**

central limit theorem, **Sampling Distribution of a Statistic, Central Limit Theorem, Student's t-Distribution**

data science and, **Student's t-Distribution**

chi-square distribution, **Chi-Squared Test: Statistical Theory**

chi-square statistic, [Chi-Square Test](#)

chi-square test, [Chi-Square Test-Further Reading](#)

detecting scientific fraud, [Fisher's Exact Test](#)

Fisher's exact test, [Fisher's Exact Test](#)

relevance for data science, [Relevance for Data Science](#)

resampling approach, [Chi-Square Test: A Resampling Approach](#)

statistical theory, [Chi-Squared Test: Statistical Theory](#)

class purity, [Measuring Homogeneity or Impurity](#)

classification, [Classification-Summary](#)

discriminant analysis, [Discriminant Analysis-Further Reading](#)

covariance matrix, [Covariance Matrix](#)

Fisher's linear discriminant, [Fisher's Linear Discriminant](#)

simple example, [A Simple Example](#)

evaluating models, [Evaluating Classification Models-Further Reading](#)

AUC metric, [AUC](#)

confusion matrix, [Confusion Matrix](#)

lift, [Lift](#)

precision, recall, and specificity, [Precision, Recall, and Specificity](#)

rare class problem, [The Rare Class Problem](#)

ROC curve, [ROC Curve](#)

K-Nearest Neighbors, **K-Nearest Neighbors**

logistic regression, **Logistic Regression-Further Reading**

and the GLM, **Logistic Regression and the GLM**

assessing the model, **Assessing the Model**

comparison to linear regression, **Linear and Logistic Regression:
Similarities and Differences**

interpreting coefficients and odds ratios, **Interpreting the Coefficients
and Odds Ratios**

logistic response function and logit, **Logistic Response Function and
Logit**

predicted values from, **Predicted Values from Logistic Regression**

more than two possible outcomes, **Classification**

naive Bayes algorithm, **Naive Bayes-Further Reading**

impracticality of exact Bayesian classification, **Why Exact Bayesian
Classification Is Impractical**

using numeric predictor variables, **Numeric Predictor Variables**

strategies for imbalanced data, **Strategies for Imbalanced Data-Further
Reading**

cost-based classification, **Cost-Based Classification**

data generation, **Data Generation**

exploring the predictions, **Exploring the Predictions**

oversampling and up/down weighting, **Oversampling and Up/Down**

Weighting

undersampling, **Undersampling**

unsupervised learning as building block, **Unsupervised Learning**

cluster mean, **K-Means Clustering, A Simple Example, Interpreting the Clusters**

clustering, **Unsupervised Learning**

application to cold-start problems, **Unsupervised Learning**

cluster analysis vs. PCA, **Interpreting the Clusters**

hierarchical, **Hierarchical Clustering-Measures of Dissimilarity, Categorical Data and Gower's Distance**

agglomerative algorithm, **The Agglomerative Algorithm**

dendrogram, **The Dendrogram**

dissimilarity measures, **Measures of Dissimilarity**

simple example, **A Simple Example**

K-means, **K-Means Clustering-Selecting the Number of Clusters, Scaling the Variables**

interpreting the clusters, **Interpreting the Clusters**

K-means algorithm, **K-Means Algorithm**

selecting the number of clusters, **Selecting the Number of Clusters**

simple example, **A Simple Example-K-Means Algorithm**

model-based, **Model-Based Clustering-Further Reading**

mixtures of normals, **Mixtures of Normals**

selecting the number of clusters, **Selecting the Number of Clusters**

problems with mixed data, **Problems with Clustering Mixed Data**

standardizing data, **Standardization (Normalization, Z-Scores)**

clusters, **K-Means Clustering**

coefficient of determination, **Assessing the Model**

coefficients

in logistic regression, **Interpreting the Coefficients and Odds Ratios**

in simple linear regression, **The Regression Equation**

estimates vs. known, **Fitted Values and Residuals**

interpretation in multiple linear regression, **Example: King County Housing Data**

complete linkage, **The Agglomerative Algorithm**

complexity parameter (cp), **Stopping the Tree from Growing**

conditional probabilities, **Naive Bayes**

conditioning variables, **Visualizing Multiple Variables**

confidence intervals, **Confidence Intervals-Further Reading, Confidence and Prediction Intervals**

generating with bootstrap, **Confidence Intervals**

level of confidence, **Confidence Intervals**

prediction intervals vs., **Confidence and Prediction Intervals**

confidence level, **Confidence Intervals-Confidence Intervals**

confounding variables, **Interpreting the Regression Equation, Confounding Variables**

confusion matrix, **Evaluating Classification Models-Confusion Matrix**

contingency tables, **Exploring Two or More Variables**

example, loan grade and status, **Two Categorical Variables**

continuous data, **Elements of Structured Data**

continuous variable as test metric, **A/B Testing**

predicting continuous value with a tree, **Predicting a Continuous Value**

contour plots, **Exploring Two or More Variables**

using with hexagonal binning, **Hexagonal Binning and Contours (Plotting Numeric versus Numeric Data)**

contrast coding systems, **Dummy Variables Representation**

control group, **A/B Testing**

advantages of using, **Why Have a Control Group?**

Cook's distance, **Influential Values**

correlated variables, **Interpreting the Regression Equation**

multicollinearity, **Multicollinearity**

predictor variables, **Correlated Predictors**

correlation, **Correlation-Further Reading**

key terms for, **Correlation**

regression vs., **Simple Linear Regression**

scatterplots, **Scatterplots**

correlation coefficient, **Correlation**

computing Pearson's correlation coefficient, **Correlation**

key concepts, **Scatterplots**

other types of, **Correlation**

correlation matrix, **Correlation**

example, correlation between telecommunication stock returns,
Correlation

cost-based classification, **Cost-Based Classification**

count data

as test metric, **A/B Testing**

Fisher's exact test for, **Fisher's Exact Test**

covariance, **Discriminant Analysis, Covariance Matrix, Computing the Principal Components**

covariance matrix

in discriminant analysis, **Covariance Matrix**

in model-based clustering, **Multivariate Normal Distribution**

using to compute Mahalanobis distance, **Distance Metrics**

cross-validation, **Cross-Validation, Choosing K**

for selection of principal components, **Interpreting Principal Components**

using for hyperparameters in boosting, **Hyperparameters and Cross-Validation**

using to estimate value of complexity parameter, **Stopping the Tree from Growing**

cumulative gains charts, **Lift**

D

d.f. (degrees of freedom), **Degrees of Freedom, Chi-Square Test**

(see also degrees of freedom)

data analysis, **Exploratory Data Analysis**

(see also exploratory data analysis)

data distribution, **Exploring the Data Distribution-Further Reading, Sampling Distribution of a Statistic**

frequency tables and histograms, **Frequency Table and Histograms**

key terms for, **Exploring the Data Distribution**

percentiles and boxplots, **Percentiles and Boxplots**

sampling distribution vs., **Sampling Distribution of a Statistic**

data frames, **Rectangular Data**

and indexes, **Data Frames and Indexes**

typical data format, **Rectangular Data**

data generation, **Strategies for Imbalanced Data, Data Generation**

data snooping, **Selection Bias**

data types

key terms for, **Elements of Structured Data**

resources for further reading, **Further Reading**

database normalization, **Standardization (Normalization, Z-Scores)**

decile gains charts, **Lift**

decision trees, **The Bootstrap, Statistical Machine Learning**

meaning in operations research, **Tree Models**

recursive partitioning algorithm, **Random Forest**

decomposition of variance, **ANOVA, F-Statistic**

degrees of freedom, **Standard Deviation and Related Estimates, Student's t-Distribution, Degrees of Freedom-Further Reading**

in chi-square test, **Chi-Squared Test: Statistical Theory**

dendrograms, **Hierarchical Clustering**

example, dendrogram of stocks, **The Dendrogram**

hierarchical clustering with mixed variable types, **Categorical Data and Gower's Distance**

density plots, **Exploring the Data Distribution, Density Estimates**

example, density of state murder rates, **Density Estimates**

dependent variable, **The Regression Equation**

(see also response)

deviation coding, **Factor Variables in Regression, Dummy Variables Representation**

deviations, **Estimates of Variability**

standard deviation and related estimates, **Standard Deviation and Related Estimates**

directional alternative hypothesis, **One-Way, Two-Way Hypothesis Test**

discrete data, **Elements of Structured Data**

discriminant analysis, **Discriminant Analysis-Further Reading**

covariance matrix, **Covariance Matrix**

extensions of, **A Simple Example**

Fisher's linear discriminant, **Fisher's Linear Discriminant**

simple example, **A Simple Example-A Simple Example**

discriminant function, **Discriminant Analysis**

discriminant weights, **Discriminant Analysis**

dispersion, **Estimates of Variability**

(see also variability, estimates of)

dissimilarity, **Hierarchical Clustering**

common measures of, **Measures of Dissimilarity**

measuring with, complete-linkage method, **The Agglomerative Algorithm**

metric in hierarchical clustering, **A Simple Example**

distance metrics, **K-Nearest Neighbors, Hierarchical Clustering**

Gower's distance and categorical data, **Categorical Data and Gower's Distance**

in hierarchical clustering, **A Simple Example, The Agglomerative Algorithm**
in K-Nearest Neighbors, **Distance Metrics**

Donoho, David, **Exploratory Data Analysis**

double blind studies, **Why Have a Control Group?**

dummy variables, **Factor Variables in Regression**

representation of factor variables in regression, **Dummy Variables Representation**

representing string factor data as numbers, **One Hot Encoder**

Durbin-Watson statistic, **Heteroskedasticity, Non-Normality and Correlated Errors**

E

EDA (see exploratory data analysis)

effect size, **Power and Sample Size, Sample Size**

elbow method, **Selecting the Number of Clusters**

ensemble learning, **Statistical Machine Learning**

staged used of K-Nearest Neighbors, **KNN as a Feature Engine**

ensemble models, **Boosting**

entropy, **Measuring Homogeneity or Impurity**

epsilon-greedy algorithm, **Multi-Arm Bandit Algorithm**

errors, **Normal Distribution**

estimates, **Estimates of Location**

indicated by hat notation, **Fitted Values and Residuals**

Euclidean distance, **Distance Metrics**

exact tests, **Exhaustive and Bootstrap Permutation Test**

Excel, pivot tables, **Two Categorical Variables**

exhaustive permutation tests, **Exhaustive and Bootstrap Permutation Test**

expectation or expected, **Chi-Square Test**

expected value, **Exploring Binary and Categorical Data, Expected Value**
calculating, **Expected Value**

explanation vs. prediction (in regression), **Prediction versus Explanation**
(Profiling)

exploratory data analysis, **Exploratory Data Analysis-Summary**

binary and categorical data, **Exploring Binary and Categorical Data-**
Correlation

correlation, **Correlation-Further Reading**

data distribution, **Exploring the Data Distribution-Further Reading**

estimates of location, **Estimates of Location-Further Reading**

estimates of variability, **Estimates of Variability-Further Reading**

exploring two or more variables, **Exploring Two or More Variables-**
Summary

rectangular data, **Rectangular Data-Estimates of Location**

Exploratory Data Analysis (Tukey), **Exploratory Data Analysis**

exponential distribution, **Poisson and Related Distributions**

calculating, **Exponential Distribution**

extrapolation

dangers of, **The Dangers of Extrapolation**

definition of, **Prediction Using Regression**

F

F-statistic, **ANOVA, F-Statistic, Assessing the Model**

facets, **Visualizing Multiple Variables**

factor variables, **Factor Variables in Regression-Ordered Factor Variables**

different codings, **Dummy Variables Representation**

dummy variables representation, **Dummy Variables Representation**

handling in logistic regression, **Fitting the model**

in naive Bayes algorithm, **Naive Bayes**

ordered, **Ordered Factor Variables**

reference coding, **Interactions and Main Effects**

with many levels, **Factor Variables with Many Levels**

factors, conversion of text columns to, **Elements of Structured Data**

failure rate, estimating, **Estimating the Failure Rate**

false discovery rate, **Multiple Testing, Multiple Testing**

false positive rate, **AUC**

feature selection

chi-square tests in, **Relevance for Data Science**

using discriminant analysis, **A Simple Example**

features, **Rectangular Data**

terminology differences, **Data Frames and Indexes**

field view (spatial data), **Nonrectangular Data Structures**

Fisher's exact test, **Fisher's Exact Test**

Fisher's linear discriminant, **Fisher's Linear Discriminant**

Fisher's scoring, **Fitting the model**

Fisher, R.A., **Fisher's Exact Test, Discriminant Analysis**

fitted values, **Simple Linear Regression, Fitted Values and Residuals**

folds, **Cross-Validation, Hyperparameters and Cross-Validation**

forward selection and backward selection, **Model Selection and Stepwise Regression**

frequency tables, **Exploring the Data Distribution**

example, population by state, **Frequency Table and Histograms**

Friedman, Jerome H. (Jerry), **Statistical Machine Learning**

G

gains, **Lift**

(see also lift)

Gallup Poll, **Random Sampling and Sample Bias**

Gallup, George, **Random Sampling and Sample Bias, Random Selection**

Galton, Francis, **Regression to the Mean**

GAM (see generalized additive models)

Gaussian distribution, **Normal Distribution**

(see also normal distribution)

generalized additive models, **Polynomial and Spline Regression, Generalized Additive Models, Exploring the Predictions**

generalized linear model (GLM), **Logistic Regression and the GLM**

Gini coefficient, **Measuring Homogeneity or Impurity**

Gini impurity, **Measuring Homogeneity or Impurity**

GLM (see generalized linear model)

Gossett, W.S., **Student's t-Distribution**

Gower's distance, **Scaling and Categorical Variables**

categorical data and, **Categorical Data and Gower's Distance**

gradient boosted trees, **Interactions and Main Effects**

gradient boosting, **The Boosting Algorithm**

definition of, **Boosting**

graphs, **Nonrectangular Data Structures**

computer science versus statistics, **Nonrectangular Data Structures**

lesson on misleading graphs, **Further Reading**

greedy algorithms, **Multi-Arm Bandit Algorithm**

H

hat notation, **Fitted Values and Residuals**

hat-value, **Testing the Assumptions: Regression Diagnostics, Influential Values**

heat maps, **Hexagonal Binning and Contours (Plotting Numeric versus Numeric Data)**

heteroskedastic errors, **Heteroskedasticity, Non-Normality and Correlated Errors**

heteroskedasticity, **Testing the Assumptions: Regression Diagnostics, Heteroskedasticity, Non-Normality and Correlated Errors**

hexagonal binning, **Exploring Two or More Variables**

example, using with contour plot, **Hexagonal Binning and Contours (Plotting Numeric versus Numeric Data)**

hierarchical clustering, **Hierarchical Clustering-Measures of Dissimilarity, Categorical Data and Gower's Distance**

agglomerative algorithm, **The Agglomerative Algorithm**

measures of dissimilarity, **Measures of Dissimilarity**

simple example, **A Simple Example**

histograms, **Exploring the Data Distribution**

example, population by state, **Frequency Table and Histograms**

homogeneity, measuring, **Measuring Homogeneity or Impurity**

hyperparameters

and cross-validation in boosting, **Hyperparameters and Cross-Validation**

for HGBost, **Hyperparameters and Cross-Validation**

in random forests, **Hyperparameters**

hypothesis tests, **Hypothesis Tests-Further Reading**

alternative hypothesis, **Alternative Hypothesis**

false discovery rate, **Multiple Testing**

null hypothesis, **The Null Hypothesis**

one-way and two-way tests, **One-Way, Two-Way Hypothesis Test**

I

impurity, **Tree Models**

measuring, **Measuring Homogeneity or Impurity**

in-sample methods to assess and tune models, **Model Selection and Stepwise Regression**

independent variables, **Simple Linear Regression, The Regression Equation**

main effects, **Interactions and Main Effects**

indexes, data frames and, **Data Frames and Indexes**

indicator variables, **Factor Variables in Regression**

inference, **Exploratory Data Analysis, Statistical Experiments and Significance Testing**

influence plots, **Influential Values**

influential values, **Testing the Assumptions: Regression Diagnostics, Influential Values**

information, **Measuring Homogeneity or Impurity**

interactions, **Interpreting the Regression Equation**

and main effects, **Interactions and Main Effects**

deciding which interaction terms to include in the model, **Interactions and Main Effects**

intercepts, **Simple Linear Regression**

in cotton exposure and lung capacity example, **The Regression Equation**

Internet of Things (IoT), **Elements of Structured Data**

interquantile range (IQR), **Estimates of Variability, Estimates Based on Percentiles**

interval endpoints, **Confidence Intervals**

K

K (in K-Nearest Neighbors), **K-Nearest Neighbors**

k-fold cross-validation, **Cross-Validation**

K-means clustering, **K-Means Clustering-Selecting the Number of Clusters**

interpreting the clusters, **Interpreting the Clusters**

K-means algorithm, **K-Means Algorithm**

selecting the number of clusters, **Selecting the Number of Clusters**

simple example, **A Simple Example-K-Means Algorithm**

using on unnormalized and normalized variables, **Scaling the Variables**

K-Nearest Neighbors, **Predicted Values from Logistic Regression, K-Nearest Neighbors-KNN as a Feature Engine**

as a feature engine, **KNN as a Feature Engine**

choosing K, **Choosing K**

distance metrics, **Distance Metrics**

example, predicting loan default, **A Small Example: Predicting Loan Default**

one hot encoder, **One Hot Encoder**

standardization, **Standardization (Normalization, Z-Scores)**

kernal density estimates, **Density Estimates**

KernSmooth package, **Density Estimates**

KNN (see K-Nearest Neighbors)

knots, **Polynomial and Spline Regression, Splines**

kurtosis, **Frequency Table and Histograms**

L

lambda, in Poisson and related distributions, **Poisson and Related Distributions**

Lasso regression, **Model Selection and Stepwise Regression, Regularization: Avoiding Overfitting**

Latent Dirichlet Allocation (LDA), **Discriminant Analysis**

leaf, **Tree Models**

least squares, **Simple Linear Regression, Least Squares**

leverage, **Testing the Assumptions: Regression Diagnostics**

influential values in regression, **Influential Values**

lift, **Evaluating Classification Models, Lift**

lift curve, **Lift**

linear discriminant analysis (LDA), **Discriminant Analysis, Exploring the Predictions**

linear regression, **Simple Linear Regression-Weighted Regression**

comparison to logistic regression, **Linear and Logistic Regression: Similarities and Differences**

fitted values and residuals, **Fitted Values and Residuals**

generalized linear model (GLM), **Logistic Regression and the GLM**

least squares, **Least Squares**

multiple, **Multiple Linear Regression-Weighted Regression**

assessing the model, **Assessing the Model**

cross-validation, **Cross-Validation**

example, King County housing data, **Example: King County Housing**

Data

model selection and stepwise regression, **Model Selection and Stepwise Regression**

weighted regression, **Weighted Regression**

prediction vs. explanation, **Prediction versus Explanation (Profiling)**

regression equation, **The Regression Equation**

Literary Digest poll of 1936, **Random Sampling and Sample Bias, Random Selection**

loadings, **Principal Components Analysis, A Simple Example**

for top five components (example), **Interpreting Principal Components**

log odds, **Logistic Regression**

log-odds function (see logit function)

log-odds ratio, **Interpreting the Coefficients and Odds Ratios**

logistic regression, **Logistic Regression-Further Reading, Exploring the Predictions**

and the generalized linear model (GLM), **Logistic Regression and the GLM**

assessing the model, **Assessing the Model**

comparison to linear regression, **Linear and Logistic Regression: Similarities and Differences**

interpreting the coefficients and odds ratios, **Interpreting the Coefficients and Odds Ratios**

logistic response function and logit, **Logistic Response Function and Logit**

predicted values from, **Predicted Values from Logistic Regression**

logit function, **Logistic Regression, Logistic Response Function and Logit**

long-tail distributions, **Long-Tailed Distributions-Further Reading**

loss, **Tree Models**

loss function, **Oversampling and Up/Down Weighting**

M

machine learning

statistics vs., **Statistical Machine Learning**

machine learning, **Statistical Machine Learning**

(see also statistical machine learning)

Mahalanobis distance, **Covariance Matrix, Distance Metrics**

main effects, **Interpreting the Regression Equation**

interactions and, **Interactions and Main Effects**

Mallows Cp, **Model Selection and Stepwise Regression**

Manhattan distance, **Distance Metrics, Regularization: Avoiding Overfitting, Categorical Data and Gower's Distance**

maximum likelihood estimation (MLE), **Fitting the model**

mean, **Estimates of Location**

formula for, **Mean**

regression to, **Regression to the Mean**

sample mean vs. population mean, **Sample Mean versus Population Mean**

trimmed mean, **Mean**

weighted mean, **Mean**

mean absolute deviation, **Estimates of Variability, A/B Testing**

formula for calculating, **Standard Deviation and Related Estimates**

mean absolute deviation from the median (MAD), **Standard Deviation and Related Estimates**

median, **Estimates of Location**

and robust estimates, **Median and Robust Estimates**

median absolute deviation, **Estimates of Variability**

metrics, **Estimates of Location**

minimum variance, **Measures of Dissimilarity**

MLE (see maximum likelihood estimation)

mode, **Exploring Binary and Categorical Data**

examples in categorical data, **Mode**

model-based clustering, **Model-Based Clustering-Further Reading**

limitations, **Selecting the Number of Clusters**

mixtures of normals, **Mixtures of Normals**

multivariate normal distribution, **Multivariate Normal Distribution**

selecting the number of clusters, **Selecting the Number of Clusters**

moments, **Frequency Table and Histograms**

multi-arm bandits, **Why Just A/B? Why Not C, D...?, Multi-Arm Bandit Algorithm-Further Reading**

definition of, **Multi-Arm Bandit Algorithm**

multicollinearity, **Interpreting the Regression Equation, Multicollinearity problems with one hot encoding, One Hot Encoder**

multicollinearity errors, **Degrees of Freedom, Dummy Variables Representation**

multiple linear regression (see linear regression)

multiple testing, **Multiple Testing-Further Reading**

bottom line for data scientists, **Multiple Testing**

Multivariate analysis, **Exploring Two or More Variables**

Multivariate normal distribution, **Multivariate Normal Distribution**

N

n (sample size), **Student's t-Distribution**

n or sample size, **Degrees of Freedom**

Naive Bayes algorithm, **Naive Bayes-Further Reading**

applying to numeric predictor variables, **Numeric Predictor Variables**

neighbors, **K-Nearest Neighbors**

network data structures, **Nonrectangular Data Structures**

nodes, **Tree Models**

non-normal residuals, **Testing the Assumptions: Regression Diagnostics**

nonlinear regression, **Polynomial and Spline Regression-Further Reading**
definition of, **Polynomial and Spline Regression**

nonrectangular data structures, **Nonrectangular Data Structures**

normal distribution, **Normal Distribution-Standard Normal and QQ-Plots**
key concepts, **Standard Normal and QQ-Plots**

standard normal and QQ-Plots, **Standard Normal and QQ-Plots**

normalization, **Standard Normal and QQ-Plots, Standardization**
(**Normalization, Z-Scores**), **K-Means Clustering**

categorical variables before clustering, **Scaling the Variables**

data distribution and, **Standardization (Normalization, Z-Scores)**

in statistics vs. database context, **Standardization (Normalization, Z-Scores)**

null hypothesis, **Hypothesis Tests, The Null Hypothesis**

numeric variables

grouped according to a categorical variable, **Categorical and Numeric Data**

numeric predictor variables for naive Bayes, **Numeric Predictor Variables**

numerical data as categorical data, **Exploring Binary and Categorical Data**

O

object representation (spatial data), **Nonrectangular Data Structures**

Occam's razor, **Model Selection and Stepwise Regression**

odds, **Logistic Regression, Logistic Response Function and Logit**

odds ratios, **Interpreting the Coefficients and Odds Ratios**

log-odds ratio and, **Interpreting the Coefficients and Odds Ratios**

omnibus tests, **ANOVA**

one hot encoder, **Factor Variables in Regression, One Hot Encoder**

one hot encoding, **Dummy Variables Representation**

one-way tests, **Hypothesis Tests, One-Way, Two-Way Hypothesis Test**

order statistics, **Estimates of Variability, Estimates Based on Percentiles**

ordered factor variables, **Ordered Factor Variables**

ordinal data, **Elements of Structured Data**

importance of the concept, **Elements of Structured Data**

ordinary least squares (OLS), **Least Squares, Heteroskedasticity, Non-Normality and Correlated Errors**

(see also least squares)

out-of-bag (OOB) estimate of error, **Random Forest**

outcome, **Rectangular Data**

outliers, **Estimates of Location, Outliers, Testing the Assumptions: Regression Diagnostics**

in regression, **Outliers**

sensitivity of correlation coefficient to, **Correlation**

sensitivity of least squares to, **Least Squares**

variance, standard deviation, mean absolute deviation and, **Standard Deviation and Related Estimates**

overfitting, **Multiple Testing**

avoiding in boosting using regularization, **Regularization: Avoiding Overfitting**

in linear regression, **Model Selection and Stepwise Regression**

oversampling, **Strategies for Imbalanced Data, Oversampling and Up/Down Weighting**

P

p-values, **Statistical Significance and P-Values, P-Value**

adjusting, **Multiple Testing**

data science and, **Data Science and P-Values**

t-statistic and, **Assessing the Model**

value of, **Value of the p-value**

pairwise comparisons, **ANOVA**

partial residual plots, **Testing the Assumptions: Regression Diagnostics, Partial Residual Plots and Nonlinearity**

in logistic regression, **Assessing the Model**

PCA (see principal components analysis)

Pearson residuals, **Chi-Square Test: A Resampling Approach**

Pearson's chi-squared test, **Chi-Squared Test: Statistical Theory**

Pearson's correlation coefficient, **Correlation**

Pearson, Karl, **Chi-Square Test, Principal Components Analysis**

penalized regression, **Model Selection and Stepwise Regression**

percentiles, **Estimates of Variability**

and boxplots, **Percentiles and Boxplots**

estimates based on, **Estimates Based on Percentiles**

precise definition of, **Estimates Based on Percentiles**

permission, obtaining for human subject testing, **Why Just A/B? Why Not C, D...?**

permutation tests, **Resampling**

exhaustive and bootstrap, **Exhaustive and Bootstrap Permutation Test**

for ANOVA, **ANOVA**

value for data science, **Permutation Tests: The Bottom Line for Data Science**

web stickiness example, **Example: Web Stickiness**

pertinent records (in searches), **Size versus Quality: When Does Size Matter?**

physical networks, **Nonrectangular Data Structures**

pie charts, **Exploring Binary and Categorical Data**

pivot tables (Excel), **Two Categorical Variables**

point estimates, **Confidence Intervals**

Poisson distributions, **Poisson and Related Distributions, Generalized Linear Models**

calculating, **Poisson Distributions**

polynomial coding, **Dummy Variables Representation**

polynomial regression, **Polynomial and Spline Regression, Polynomial**

population, **Random Sampling and Sample Bias**

sample mean vs. population mean, **Sample Mean versus Population Mean**

posterior probability, **Naive Bayes, The Naive Solution**

power and sample size, **Power and Sample Size-Further Reading**

precision, **Evaluating Classification Models**

in classification models, **Precision, Recall, and Specificity**

predicted values, **Fitted Values and Residuals**

(see also fitted values)

prediction

explanation vs., in linear regression, **Prediction versus Explanation (Profiling)**

harnessing results from multiple trees, **How Trees Are Used**

K-Nearest Neighbors, **K-Nearest Neighbors**

using as first stage, **KNN as a Feature Engine**

predicted values from logistic regression, **Predicted Values from Logistic Regression**

unsupervised learning and, **Unsupervised Learning**

using regression, **Prediction Using Regression-Factor Variables in Regression**

confidence and prediction intervals, **Confidence and Prediction Intervals**

dangers of extrapolation, **The Dangers of Extrapolation**

prediction intervals, **Prediction Using Regression**

confidence intervals vs., **Confidence and Prediction Intervals**

predictor variables, **Data Frames and Indexes, The Regression Equation**
(see also independent variables)

correlated, **Correlated Predictors**

in linear discriminant analysis, more than two, **A Simple Example**

in naive Bayes algorithm, **Naive Bayes**

main effects, **Interactions and Main Effects**

numeric, applying naive Bayes to, **Numeric Predictor Variables**

relationship between response and, **Partial Residual Plots and Nonlinearity**

principal components, **Principal Components Analysis**

principal components analysis, **Principal Components Analysis-Further Reading**

cluster analysis vs., **Interpreting the Clusters**

computing the principal components, **Computing the Principal Components**

interpreting principal components, **Interpreting Principal Components**

scaling the variables, **Scaling the Variables**

simple example, **A Simple Example-A Simple Example**

standardizing data, **Standardization (Normalization, Z-Scores)**

probability theory, **Exploratory Data Analysis**

profiling vs. explanation, **Prediction versus Explanation (Profiling)**

propensity score, **Classification**

proxy variables, **Example: Web Stickiness**

pruning, **Tree Models, Stopping the Tree from Growing**

pseudo-residuals, **The Boosting Algorithm**

Q

QQ-Plots, **Normal Distribution**

example, returns for Netflix, **Long-Tailed Distributions**

standard normal and, **Standard Normal and QQ-Plots**

quadratic discriminant analysis, **A Simple Example**

quantiles, **Estimates Based on Percentiles**

R function, quantile, **Estimates Based on Percentiles**

R

R-squared, **Multiple Linear Regression, Assessing the Model**

random forests, **Interactions and Main Effects, Tree Models, Random Forest-Hyperparameters**

better predictive performance than single trees, **How Trees Are Used**

determining variable importance, **Variable Importance**

hyperparameters, **Hyperparameters**

random sampling, **Random Sampling and Sample Bias-Further Reading**

bias, **Bias**

key terms for, **Random Sampling and Sample Bias**

random selection, **Random Selection**

sample mean vs. population mean, **Sample Mean versus Population Mean**

size versus quality, **Size versus Quality: When Does Size Matter?**

random subset of variables, **Random Forest**

randomization, **A/B Testing**

randomization tests, **Resampling**

(see also permutation tests)

randomness, misinterpreting, **Hypothesis Tests**

range, **Estimates of Variability, Estimates Based on Percentiles**

rare class problem, **The Rare Class Problem**

recall, **Evaluating Classification Models, Precision, Recall, and Specificity**

receiver operating characteristics (see ROC curve)

records, **Rectangular Data, Simple Linear Regression**

rectangular data, **Rectangular Data-Estimates of Location**

terminology differences, **Data Frames and Indexes**

recursive partitioning, **Tree Models, The Recursive Partitioning Algorithm, Random Forest**

reference coding, **Factor Variables in Regression-Dummy Variables Representation, Interactions and Main Effects, Logistic Regression and the GLM**

regression, **Regression and Prediction-Summary**

causation and, **Prediction versus Explanation (Profiling)**

diagnostics, **Testing the Assumptions: Regression Diagnostics-Polynomial and Spline Regression**

heteroskedasticity, non-normality, and correlated errors,
Heteroskedasticity, Non-Normality and Correlated Errors

influential values, **Influential Values**

outliers, **Outliers**

parial residual plots and nonlinearity, **Partial Residual Plots and Nonlinearity**

using scatterplot smoothers, **Heteroskedasticity, Non-Normality and Correlated Errors**

different meanings of the term, **Least Squares**

factor variables in, **Factor Variables in Regression-Ordered Factor Variables**

ordered factor variables, **Ordered Factor Variables**

with many levels, **Factor Variables with Many Levels**

interpreting the regression equation, **Interpreting the Regression Equation-Interactions and Main Effects**

confounding variables, **Confounding Variables**

correlated predictors, **Correlated Predictors**

interactions and main effects, **Interactions and Main Effects**

multicollinearity, **Multicollinearity**

KNN (K-Nearest Neighbors), **KNN as a Feature Engine**

logistic regression, **Logistic Regression-Further Reading**

comparison to linear regression, **Linear and Logistic Regression: Similarities and Differences**

multiple linear regression, **Multiple Linear Regression-Weighted Regression**

polynomial and spline regression, **Polynomial and Spline Regression-Summary**

generalized additive models, **Generalized Additive Models**

polynomial regression, **Polynomial**

splines, **Splines**

prediction with, **Prediction Using Regression-Factor Variables in Regression**

confidence and prediction intervals, **Confidence and Prediction Intervals**

dangers of extrapolation, **The Dangers of Extrapolation**

ridge regression, **Regularization: Avoiding Overfitting**

simple linear regression, **Simple Linear Regression-Further Reading**

fitted values and residuals, **Fitted Values and Residuals**

least squares, **Least Squares**

prediction vs. explanation, **Prediction versus Explanation (Profiling)**

regression equation, **The Regression Equation**

unsupervised learning as building block, **Unsupervised Learning**

with a tree, **Predicting a Continuous Value**

regression coefficient, **Simple Linear Regression**

in cotton exposure and lung capacity example, **The Regression Equation**

regression to the mean, **Regression to the Mean**

regularization, **Boosting**

avoiding overfitting with, **Regularization: Avoiding Overfitting**

replacement (in sampling), **Random Sampling and Sample Bias**

bootstrap, **The Bootstrap**

representativeness, **Random Sampling and Sample Bias**

resampling, **The Bootstrap, Resampling-For Further Reading**

bootstrapping vs., **Resampling versus Bootstrapping**

permutation tests, **Permutation Test**

exhaustive and bootstrap tests, **Exhaustive and Bootstrap Permutation Test**

value for data science, **Permutation Tests: The Bottom Line for Data Science**

web stickiness example, **Example: Web Stickiness**

using in chi-square test, **Chi-Square Test: A Resampling Approach**

residual standard error, **Multiple Linear Regression, Assessing the Model**

residual sum of squares, **Least Squares**

(see also least squares)

residuals, **Simple Linear Regression, Fitted Values and Residuals**

computing, **Fitted Values and Residuals**

distribution of, **Heteroskedasticity, Non-Normality and Correlated Errors**

standardized, **Testing the Assumptions: Regression Diagnostics**

response, **Simple Linear Regression, The Regression Equation**

relationship between predictor variable and, **Partial Residual Plots and Nonlinearity**

ridge regression, **Model Selection and Stepwise Regression, Regularization: Avoiding Overfitting**

robust, **Estimates of Location**

robust estimates of location

example, population and murder rate by state, **Example: Location Estimates of Population and Murder Rates**

mean absolute deviation from the median, **Standard Deviation and Related Estimates**

median, **Median and Robust Estimates**

outliers and, **Outliers**

ROC curve, **ROC Curve**

root mean squared error (RMSE), **Multiple Linear Regression, Assessing the Model, Predicting a Continuous Value**

RSE (see residual standard error)

RSS (residual sum of squares), **Least Squares**

(see also least squares)

S

sample bias, **Random Sampling and Sample Bias, Random Sampling and Sample Bias**

sample statistic, **Sampling Distribution of a Statistic**

samples

definition of, **Random Sampling and Sample Bias**

sample size, power and, **Power and Sample Size-Further Reading**

terminology differences, **Data Frames and Indexes**

sampling, **Data and Sampling Distributions-Summary**

binomial distribution, **Binomial Distribution-Further Reading**

bootstrap, **The Bootstrap-Further Reading**

confidence intervals, **Confidence Intervals-Further Reading**

long-tail distributions, **Long-Tailed Distributions-Further Reading**

normal distribution, **Normal Distribution-Standard Normal and QQ-Plots**

oversampling imbalanced data, **Oversampling and Up/Down Weighting**

Poisson and related distributions, **Poisson and Related Distributions-Summary**

estimating failure rate, **Estimating the Failure Rate**

exponential distribution, **Exponential Distribution**

Poisson distribution, **Poisson Distributions**

Weibull distribution, **Weibull Distribution**

population versus sample, **Data and Sampling Distributions**

random sampling and sample bias, **Random Sampling and Sample Bias-Further Reading**

sampling distribution of a statistic, **Sampling Distribution of a Statistic-Further Reading**

selection bias, **Selection Bias-Further Reading**

Student's t-distribution, **Student's t-Distribution-Further Reading**

Thompson's sampling, **Multi-Arm Bandit Algorithm**

undersampling imbalanced data, **Undersampling**

with and without replacement, **Random Sampling and Sample Bias, The Bootstrap, Resampling**

sampling distribution, **Sampling Distribution of a Statistic-Further Reading**

central limit theorem, **Central Limit Theorem**

data distribution vs., **Sampling Distribution of a Statistic**

standard error, **Standard Error**

scale parameter (Weibull distribution), **Weibull Distribution**

scaling data and categorical variables, **Scaling and Categorical Variables-Summary**

dominant variables, **Dominant Variables**

Gower's distance and categorical data, **Categorical Data and Gower's Distance**

problems clustering mixed data, **Problems with Clustering Mixed Data**

scaling the variables, **Scaling the Variables**

scatterplot smoothers, **Heteroskedasticity, Non-Normality and Correlated Errors**

scatterplots, **Correlation**

example, returns for ATT and Verizon, **Scatterplots**

scientific fraud, detecting, **Fisher's Exact Test**

screeplots, **Principal Components Analysis, Interpreting Principal Components**

for PCA of top stocks, **Dominant Variables**

searches

search queries on Google, **Size versus Quality: When Does Size Matter?**

vast search effect, **Selection Bias**

selection bias, **Selection Bias-Further Reading**

regression to the mean, **Regression to the Mean**

self-selection sampling bias, **Random Sampling and Sample Bias**

sensitivity, **Evaluating Classification Models, Precision, Recall, and Specificity**

shape parameter (Weibull distribution), **Weibull Distribution**

signal-to-noise ratio, **Choosing K**

significance level, **Power and Sample Size, Sample Size**

significance tests, **Hypothesis Tests, Data Science and P-Values**
(see also hypothesis tests)

simple random sample, **Random Sampling and Sample Bias**

single linkage, **Measures of Dissimilarity**

skew, **Long-Tailed Distributions**

skewness, **Frequency Table and Histograms**

slope, **Simple Linear Regression**
(see also regression coefficient)

in regression equation, **The Regression Equation**

SMOTE algorithm, **Data Generation**

spatial data structures, **Nonrectangular Data Structures**

specificity, **Evaluating Classification Models, Precision, Recall, and Specificity**

spline regression, **Polynomial and Spline Regression, Splines**

splines, **Splines**

split value, **Tree Models**

square-root of n rule, **Standard Error**

SS (sum of squares), **ANOVA**

withing cluster sum of squares, **K-Means Clustering**

standard deviation, **Estimates of Variability**

and related estimates, **Standard Deviation and Related Estimates**

covariance matrix and, **Covariance Matrix**

in statistical testing output, **A/B Testing**

sensitivity to outliers, **Standard Deviation and Related Estimates**

standard error vs., **Standard Error**

standard error, **Sampling Distribution of a Statistic**

formula for calculating, **Standard Error**

standard deviation vs., **Standard Error**

standard normal distribution, **Normal Distribution, Standard Normal and QQ-Plots**

standardization, **Standard Normal and QQ-Plots, K-Nearest Neighbors, K-Means Clustering**

in K-Nearest Neighbors, **Standardization (Normalization, Z-Scores)**

standardized residuals, **Testing the Assumptions: Regression Diagnostics**

examining to detect outliers, **Outliers**

statistical experiments and significance testing, **Statistical Experiments and Significance Testing-Summary**

A/B testing, **A/B Testing-For Further Reading**

chi-square test, **Chi-Square Test-Further Reading**

degrees of freedom, **Degrees of Freedom-Further Reading**

hypothesis tests, **Hypothesis Tests-Further Reading**

multi-arm bandit algorithm, **Multi-Arm Bandit Algorithm-Further Reading**

multiple tests, **Multiple Testing-Further Reading**

power and sample size, **Power and Sample Size-Further Reading**

resampling, **Resampling-Statistical Significance and P-Values**

statistical significance and p-values, **Statistical Significance and P-Values-Further Reading**

alpha, **Alpha**

data science and p-values, **Data Science and P-Values**

p-values, **P-Value**

type 1 and type 2 errors, **Type 1 and Type 2 Errors**

value of p-values, **Value of the p-value**

t-tests, **t-Tests-Further Reading**

statistical inference, classical inference pipeline, **Statistical Experiments and Significance Testing**

statistical machine learning, **Statistical Machine Learning-Summary**

bagging and the random forest, **Bagging and the Random Forest-Hyperparameters**

boosting, **Boosting-Summary**

avoiding overfitting using regularization, **Regularization: Avoiding Overfitting**

hyperparameters and cross-validation, **Hyperparameters and Cross-Validation**

XGBoost, **XGBoost**

K-Nearest Neighbors, **K-Nearest Neighbors-KNN as a Feature Engine as a feature engine, KNN as a Feature Engine**

choosing K, **Choosing K**

distance metrics, **Distance Metrics**

example, predicting loan default, **A Small Example: Predicting Loan Default**

one hot encoder, **One Hot Encoder**

standardization, **Standardization (Normalization, Z-Scores)**

tree models, **Tree Models-Further Reading**

measuring homogeneity or impurity, **Measuring Homogeneity or**

Impurity

predicting a continuous value, **Predicting a Continuous Value**

recursive partitioning algorithm, **The Recursive Partitioning Algorithm**

simple example, **A Simple Example**

stopping tree growth, **Stopping the Tree from Growing**

uses of trees, **How Trees Are Used**

statistical moments, **Frequency Table and Histograms**

statistical significance, **Permutation Test**

statistics vs. machine learning, **Statistical Machine Learning**

stepwise regression, **Model Selection and Stepwise Regression**

stochastic gradient boosting, **The Boosting Algorithm**

definition of, **Boosting**

XGBoost implementation, **XGBoost-Hyperparameters and Cross-Validation**

stratified sampling, **Random Sampling and Sample Bias, Random Selection**

structured data, **Elements of Structured Data-Further Reading**

Student's t-distribution, **Student's t-Distribution-Further Reading**

subjects, **A/B Testing**

success, **Binomial Distribution**

sum contrasts, **Dummy Variables Representation**

T

t-distributions, **Student's t-Distribution-Further Reading, t-Tests**

data science and, **Student's t-Distribution**

t-statistic, **t-Tests, Multiple Linear Regression, Assessing the Model**

t-tests, **t-Tests-Further Reading**

tail, **Long-Tailed Distributions**

target shuffling, **Selection Bias**

test sample, **Evaluating Classification Models**

test statistic, **A/B Testing, t-Tests**

selecting before the experiment, **Why Have a Control Group?**

Thompson sampling, **Multi-Arm Bandit Algorithm**

time series data, **Nonrectangular Data Structures**

time-to-failure analysis, **Weibull Distribution**

treatment, **A/B Testing**

treatment group, **A/B Testing**

tree models, **Interactions and Main Effects, Exploring the Predictions, Tree Models**

how trees are used, **How Trees Are Used**

measuring homogeneity or impurity, **Measuring Homogeneity or Impurity**

predicting a continuous value, **Predicting a Continuous Value**

recursive partitioning algorithm, **The Recursive Partitioning Algorithm**

simple example, **A Simple Example**

stopping tree growth, **Stopping the Tree from Growing**

Trellis graphics, **Visualizing Multiple Variables**

trials, **Binomial Distribution**

trimmed mean, **Estimates of Location**

formula for, **Mean**

Tukey, John Wilder, **Exploratory Data Analysis**

two-way tests, **Hypothesis Tests, One-Way, Two-Way Hypothesis Test**

type 1 errors, **Statistical Significance and P-Values, Type 1 and Type 2 Errors, Multiple Testing**

type 2 errors, **Statistical Significance and P-Values, Type 1 and Type 2 Errors**

U

unbiased estimates, **Standard Deviation and Related Estimates**

undersampling, **Undersampling**

uniform random distribution, **Fisher's Exact Test**

univariate analysis, **Exploring Two or More Variables**

unsupervised learning, **Unsupervised Learning-Summary**

and prediction, **Unsupervised Learning**

hierarchical clustering, **Hierarchical Clustering-Measures of Dissimilarity**

agglomerative algorithm, **The Agglomerative Algorithm**

dendrogram, **The Dendrogram**

dissimilarity measures, **Measures of Dissimilarity**

simple example, **A Simple Example**

K-means clustering, **K-Means Clustering-Selecting the Number of Clusters**

interpreting the clusters, **Interpreting the Clusters**

K-means algorithm, **K-Means Algorithm**

selecting the number of customers, **Selecting the Number of Clusters**

simple example, **A Simple Example-K-Means Algorithm**

model-based clustering, **Model-Based Clustering-Further Reading**

mixtures of normals, **Mixtures of Normals**

multivariate normal distribution, **Multivariate Normal Distribution**

selecting the number of clusters, **Selecting the Number of Clusters**

principal components analysis, **Principal Components Analysis-Further Reading**

computing the principal components, **Computing the Principal Components**

interpreting principal components, **Interpreting Principal Components**

simple example, **A Simple Example-A Simple Example**

scaling and categorical variables, **Scaling and Categorical Variables-**

Summary

dominant variables, **Dominant Variables**

Gower's distance and categorical data, **Categorical Data and Gower's Distance**

problems clustering mixed data, **Problems with Clustering Mixed Data**

scaling the variables, **Scaling the Variables**

up weight or down weight, **Strategies for Imbalanced Data, Oversampling and Up/Down Weighting**

uplift vs. lift, **Lift**

V

validation sample, **Evaluating Classification Models**

variability

variability, estimates of, **Estimates of Variability-Further Reading**

example, murder rate by state population, **Example: Variability Estimates of State Population**

key terminology, **Estimates of Variability**

percentiles, **Estimates Based on Percentiles**

standard deviation and related estimates, **Standard Deviation and Related Estimates**

variables

exploring two or more, **Exploring Two or More Variables-Summary**

categorical and numeric data, **Categorical and Numeric Data**

hexagonal binning and contours, **Hexagonal Binning and Contours (Plotting Numeric versus Numeric Data)**

key concepts, **Visualizing Multiple Variables**

visualizing multiple variables, **Visualizing Multiple Variables**

importance of, determining in random forests, **Variable Importance**

rescaling with z-scores, **Standardization (Normalization, Z-Scores)**

variance, **Estimates of Variability**

analysis of (ANOVA), **ANOVA**

formula for calculating, **Standard Deviation and Related Estimates**

sensitivity to outliers, **Standard Deviation and Related Estimates**

vast search effect, **Selection Bias**

violin plots, **Exploring Two or More Variables**

combining with a boxplot, example, **Categorical and Numeric Data**

W

Ward's method, **Measures of Dissimilarity**

web stickiness example (permutation test), **Example: Web Stickiness**

web testing

bandit algorithms in, **Multi-Arm Bandit Algorithm**

deciding how long a test should run, **Power and Sample Size**

Weibull distribution, **Poisson and Related Distributions**

calculating, **Weibull Distribution**

weighted mean, **Estimates of Location**

expected value, **Expected Value**

weighted median, **Estimates of Location, Median and Robust Estimates**

formula for calculating, **Mean**

weighted regression, **Multiple Linear Regression, Weighted Regression**

weights, **Simple Linear Regression**

component loadings, **A Simple Example**

whiskers (in boxplots), **Percentiles and Boxplots**

wins, **Multi-Arm Bandit Algorithm**

within cluster sum of squares (SS), **K-Means Clustering**

X

XGBoost, **XGBoost-Hyperparameters and Cross-Validation**

hyperparameters, **Hyperparameters and Cross-Validation**

Z

z-distribution, **Standard Normal and QQ-Plots**

(see also normal distribution)

z-score, **Normal Distribution, Strategies for Imbalanced Data, K-Nearest Neighbors, Standardization (Normalization, Z-Scores)**

converting data to, **Standard Normal and QQ-Plots**

rescaling variables, **Standardization (Normalization, Z-Scores)**

About the Authors

Peter Bruce founded and grew the Institute for Statistics Education at Statistics.com, which now offers about 100 courses in statistics, roughly a third of which are aimed at the data scientist. In recruiting top authors as instructors and forging a marketing strategy to reach professional data scientists, Peter has developed both a broad view of the target market and his own expertise to reach it.

Andrew Bruce has over 30 years of experience in statistics and data science in academia, government, and business. He has a PhD in statistics from the University of Washington and has published numerous papers in refereed journals. He has developed statistical-based solutions to a wide range of problems faced by a variety of industries, from established financial firms to internet startups, and offers a deep understanding of the practice of data science.

Colophon

The animal on the cover of *Practical Statistics for Data Scientists* is a lined shore crab (*Pachygrapsus crassipes*), also known as a striped shore crab. It is found along the coasts and beaches of the Pacific Ocean in North America, Central America, Korea, and Japan. These crustaceans live under rocks, in tidepools, and within crevices. They spend about half their time on land, and periodically return to the water to wet their gills.

The lined shore crab is named for the green stripes on its brown-black carapace. It has red claws and purple legs, which also have a striped or mottled pattern. The crab generally grows to be 3–5 centimeters in size; females are slightly smaller. Their eyes are on flexible stalks that can rotate to give them a full field of vision as they walk.

Crabs are omnivores, feeding primarily on algae, but also mollusks, worms, fungi, dead animals, and other crustaceans (depending on what is available). They moult many times as they grow to adulthood, taking in water to expand and crack open their old shell. Once this is achieved, they spend several difficult hours getting free, and then must hide until the new shell hardens.

Many of the animals on O'Reilly covers are endangered; all of them are important to the world. To learn more about how you can help, go to animals.oreilly.com.

The cover image is from *Pictorial Museum of Animated Nature*. The cover fonts are URW Typewriter and Guardian Sans. The text font is Adobe Minion Pro; the heading font is Adobe Myriad Condensed; and the code font is Dalton Maag's Ubuntu Mono.

Preface

What to Expect

Conventions Used in This Book

Using Code Examples

Safari® Books Online

How to Contact Us

Acknowledgments

1. Exploratory Data Analysis

Elements of Structured Data

Further Reading

Rectangular Data

Data Frames and Indexes

Nonrectangular Data Structures

Further Reading

Estimates of Location

Mean

Median and Robust Estimates

Example: Location Estimates of Population and Murder Rates

Further Reading

Estimates of Variability

Standard Deviation and Related Estimates

Estimates Based on Percentiles

Example: Variability Estimates of State Population

Further Reading

Exploring the Data Distribution

Percentiles and Boxplots

Frequency Table and Histograms

Density Estimates

Further Reading

Exploring Binary and Categorical Data

Mode

Expected Value

Further Reading

Correlation

Scatterplots

Further Reading

Exploring Two or More Variables

Hexagonal Binning and Contours (Plotting Numeric versus Numeric Data)

Two Categorical Variables

Categorical and Numeric Data

Visualizing Multiple Variables

Further Reading

Summary

2. Data and Sampling Distributions

Random Sampling and Sample Bias

Bias

Random Selection

Size versus Quality: When Does Size Matter?

Sample Mean versus Population Mean

Further Reading

Selection Bias

Regression to the Mean

Further Reading

Sampling Distribution of a Statistic

Central Limit Theorem

Standard Error

Further Reading

The Bootstrap

Resampling versus Bootstrapping

Further Reading

Confidence Intervals

Further Reading

Normal Distribution

Standard Normal and QQ-Plots

Long-Tailed Distributions

Further Reading

Student's t-Distribution

Further Reading

Binomial Distribution

Further Reading

Poisson and Related Distributions

Poisson Distributions

Exponential Distribution

Estimating the Failure Rate

Weibull Distribution

Further Reading

Summary

3. Statistical Experiments and Significance Testing

A/B Testing

Why Have a Control Group?

Why Just A/B? Why Not C, D...?

For Further Reading

Hypothesis Tests

The Null Hypothesis

Alternative Hypothesis

One-Way, Two-Way Hypothesis Test

Further Reading

Resampling

Permutation Test

Example: Web Stickiness

Exhaustive and Bootstrap Permutation Test

Permutation Tests: The Bottom Line for Data Science

For Further Reading

Statistical Significance and P-Values

P-Value

Alpha

Type 1 and Type 2 Errors

Data Science and P-Values

Further Reading

t-Tests

Further Reading

Multiple Testing

Further Reading

Degrees of Freedom

Further Reading

ANOVA

F-Statistic

Two-Way ANOVA

Further Reading

Chi-Square Test

Chi-Square Test: A Resampling Approach

Chi-Squared Test: Statistical Theory

Fisher's Exact Test

Relevance for Data Science

Further Reading

Multi-Arm Bandit Algorithm

Further Reading

Power and Sample Size

Sample Size

Further Reading

Summary

4. Regression and Prediction

Simple Linear Regression

The Regression Equation

Fitted Values and Residuals

Least Squares

Prediction versus Explanation (Profiling)

Further Reading

Multiple Linear Regression

Example: King County Housing Data

Assessing the Model

Cross-Validation

Model Selection and Stepwise Regression

Weighted Regression

Prediction Using Regression

The Dangers of Extrapolation

Confidence and Prediction Intervals

Factor Variables in Regression

Dummy Variables Representation

Factor Variables with Many Levels

Ordered Factor Variables

Interpreting the Regression Equation

Correlated Predictors

Multicollinearity

Confounding Variables

Interactions and Main Effects

Testing the Assumptions: Regression Diagnostics

Outliers

Influential Values

Heteroskedasticity, Non-Normality and Correlated Errors

Partial Residual Plots and Nonlinearity

Polynomial and Spline Regression

Polynomial

Splines

Generalized Additive Models

Further Reading

Summary

5. Classification

Naive Bayes

Why Exact Bayesian Classification Is Impractical

The Naive Solution

Numeric Predictor Variables

Further Reading

Discriminant Analysis

Covariance Matrix

Fisher's Linear Discriminant

A Simple Example

Further Reading

Logistic Regression

Logistic Response Function and Logit

Logistic Regression and the GLM

Generalized Linear Models

Predicted Values from Logistic Regression

Interpreting the Coefficients and Odds Ratios

Linear and Logistic Regression: Similarities and Differences

Assessing the Model

Further Reading

Evaluating Classification Models

Confusion Matrix

The Rare Class Problem

Precision, Recall, and Specificity

ROC Curve

AUC

Lift

Further Reading

Strategies for Imbalanced Data

Undersampling

Oversampling and Up/Down Weighting

Data Generation

Cost-Based Classification

Exploring the Predictions

Further Reading

Summary

6. Statistical Machine Learning

K-Nearest Neighbors

A Small Example: Predicting Loan Default

Distance Metrics

One Hot Encoder

Standardization (Normalization, Z-Scores)

Choosing K

KNN as a Feature Engine

Tree Models

A Simple Example

The Recursive Partitioning Algorithm

Measuring Homogeneity or Impurity

Stopping the Tree from Growing

Predicting a Continuous Value

How Trees Are Used

Further Reading

Bagging and the Random Forest

Bagging

Random Forest

Variable Importance

Hyperparameters

Boosting

The Boosting Algorithm

XGBoost

Regularization: Avoiding Overfitting

Hyperparameters and Cross-Validation

Summary

7. Unsupervised Learning

Principal Components Analysis

A Simple Example

Computing the Principal Components

Interpreting Principal Components

Further Reading

K-Means Clustering

A Simple Example

K-Means Algorithm

Interpreting the Clusters

Selecting the Number of Clusters

Hierarchical Clustering

A Simple Example

The Dendrogram

The Agglomerative Algorithm

Measures of Dissimilarity

Model-Based Clustering

Multivariate Normal Distribution

Mixtures of Normals

Selecting the Number of Clusters

Further Reading

Scaling and Categorical Variables

Scaling the Variables

Dominant Variables

Categorical Data and Gower's Distance

Problems with Clustering Mixed Data

Summary

Bibliography

Index