



**Figure 2.14.** Discretizing  $x$  and  $y$  attributes for four groups (classes) of points.

a *department name* attribute might have dozens of different values. In this situation, we could use our knowledge of the relationships among different departments to combine departments into larger groups, such as *engineering*, *social sciences*, or *biological sciences*. If domain knowledge does not serve as a useful guide or such an approach results in poor classification performance, then it is necessary to use a more empirical approach, such as grouping values together only if such a grouping results in improved classification accuracy or achieves some other data mining objective.

### 2.3.7 Variable Transformation

A **variable transformation** refers to a transformation that is applied to all the values of a variable. (We use the term variable instead of attribute to adhere to common usage, although we will also refer to attribute transformation on occasion.) In other words, for each object, the transformation is applied to the value of the variable for that object. For example, if only the magnitude of a variable is important, then the values of the variable can be transformed by taking the absolute value. In the following section, we discuss two important types of variable transformations: simple functional transformations and normalization.

### Simple Functions

For this type of variable transformation, a simple mathematical function is applied to each value individually. If  $x$  is a variable, then examples of such transformations include  $x^k$ ,  $\log x$ ,  $e^x$ ,  $\sqrt{x}$ ,  $1/x$ ,  $\sin x$ , or  $|x|$ . In statistics, variable transformations, especially  $\sqrt{x}$ ,  $\log$ , and  $1/x$ , are often used to transform data that does not have a Gaussian (normal) distribution into data that does. While this can be important, other reasons often take precedence in data mining. Suppose the variable of interest is the number of data bytes in a session, and the number of bytes ranges from 1 to 1 billion. This is a huge range, and it can be advantageous to compress it by using a  $\log_{10}$  transformation. In this case, sessions that transferred  $10^8$  and  $10^9$  bytes would be more similar to each other than sessions that transferred 10 and 1000 bytes ( $9 - 8 = 1$  versus  $3 - 1 = 2$ ). For some applications, such as network intrusion detection, this may be what is desired, since the first two sessions most likely represent transfers of large files, while the latter two sessions could be two quite distinct types of sessions.

Variable transformations should be applied with caution because they change the nature of the data. While this is what is desired, there can be problems if the nature of the transformation is not fully appreciated. For instance, the transformation  $1/x$  reduces the magnitude of values that are 1 or larger, but increases the magnitude of values between 0 and 1. To illustrate, the values  $\{1, 2, 3\}$  go to  $\{1, \frac{1}{2}, \frac{1}{3}\}$ , but the values  $\{1, \frac{1}{2}, \frac{1}{3}\}$  go to  $\{1, 2, 3\}$ . Thus, for all sets of values, the transformation  $1/x$  reverses the order. To help clarify the effect of a transformation, it is important to ask questions such as the following: What is the desired property of the transformed attribute? Does the order need to be maintained? Does the transformation apply to all values, especially negative values and 0? What is the effect of the transformation on the values between 0 and 1? Exercise 21 on page 129 explores other aspects of variable transformation.

### Normalization or Standardization

The goal of standardization or normalization is to make an entire set of values have a particular property. A traditional example is that of “standardizing a variable” in statistics. If  $\bar{x}$  is the mean (average) of the attribute values and  $s_x$  is their standard deviation, then the transformation  $x' = (x - \bar{x})/s_x$  creates a new variable that has a mean of 0 and a standard deviation of 1. If different variables are to be used together, e.g., for clustering, then such a transformation is often necessary to avoid having a variable with large values