

examples of *eager learners*. **Eager learners**, when given a set of training tuples, will construct a generalization (i.e., classification) model before receiving new (e.g., test) tuples to classify. We can think of the learned model as being ready and eager to classify previously unseen tuples.

Imagine a contrasting lazy approach, in which the learner instead waits until the last minute before doing any model construction to classify a given test tuple. That is, when given a training tuple, a **lazy learner** simply stores it (or does only a little minor processing) and waits until it is given a test tuple. Only when it sees the test tuple does it perform generalization to classify the tuple based on its similarity to the stored training tuples. Unlike eager learning methods, lazy learners do less work when a training tuple is presented and more work when making a classification or numeric prediction. Because lazy learners store the training tuples or “instances,” they are also referred to as **instance-based learners**, even though all learning is essentially based on instances.

When making a classification or numeric prediction, lazy learners can be computationally expensive. They require efficient storage techniques and are well suited to implementation on parallel hardware. They offer little explanation or insight into the data's structure. Lazy learners, however, naturally support incremental learning. They are able to model complex decision spaces having hyperpolygonal shapes that may not be as easily describable by other learning algorithms (such as hyperrectangular shapes modeled by decision trees). In this section, we look at two examples of lazy learners: *k-nearest-neighbor classifiers* (Section 9.5.1) and *case-based reasoning classifiers* (Section 9.5.2).

### 9.5.1 **k-Nearest-Neighbor Classifiers**

The *k*-nearest-neighbor method was first described in the early 1950s. The method is labor intensive when given large training sets, and did not gain popularity until the 1960s when increased computing power became available. It has since been widely used in the area of pattern recognition.

Nearest-neighbor classifiers are based on learning by analogy, that is, by comparing a given test tuple with training tuples that are similar to it. The training tuples are described by  $n$  attributes. Each tuple represents a point in an  $n$ -dimensional space. In this way, all the training tuples are stored in an  $n$ -dimensional pattern space. When given an unknown tuple, a **k-nearest-neighbor classifier** searches the pattern space for the  $k$  training tuples that are closest to the unknown tuple. These  $k$  training tuples are the  $k$  “nearest neighbors” of the unknown tuple.

“Closeness” is defined in terms of a distance metric, such as Euclidean distance. The Euclidean distance between two points or tuples, say,  $X_1 = (x_{11}, x_{12}, \dots, x_{1n})$  and  $X_2 = (x_{21}, x_{22}, \dots, x_{2n})$ , is

$$\text{dist}(X_1, X_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}. \quad (9.22)$$

In other words, for each numeric attribute, we take the difference between the corresponding values of that attribute in tuple  $X_1$  and in tuple  $X_2$ , square this difference, and accumulate it. The square root is taken of the total accumulated distance count. Typically, we normalize the values of each attribute before using Eq. (9.22). This helps prevent attributes with initially large ranges (e.g., *income*) from outweighing attributes with initially smaller ranges (e.g., binary attributes). Min-max normalization, for example, can be used to transform a value  $v$  of a numeric attribute  $A$  to  $v'$  in the range  $[0, 1]$  by computing

$$v' = \frac{v - \min_A}{\max_A - \min_A}, \quad (9.23)$$

where  $\min_A$  and  $\max_A$  are the minimum and maximum values of attribute  $A$ . Chapter 3 describes other methods for data normalization as a form of data transformation.

For  $k$ -nearest-neighbor classification, the unknown tuple is assigned the most common class among its  $k$ -nearest neighbors. When  $k = 1$ , the unknown tuple is assigned the class of the training tuple that is closest to it in pattern space. Nearest-neighbor classifiers can also be used for numeric prediction, that is, to return a real-valued prediction for a given unknown tuple. In this case, the classifier returns the average value of the real-valued labels associated with the  $k$ -nearest neighbors of the unknown tuple.

*“But how can distance be computed for attributes that are not numeric, but nominal (or categorical) such as color?”* The previous discussion assumes that the attributes used to describe the tuples are all numeric. For nominal attributes, a simple method is to compare the corresponding value of the attribute in tuple  $X_1$  with that in tuple  $X_2$ . If the two are identical (e.g., tuples  $X_1$  and  $X_2$  both have the color blue), then the difference between the two is taken as 0. If the two are different (e.g., tuple  $X_1$  is blue but tuple  $X_2$  is red), then the difference is considered to be 1. Other methods may incorporate more sophisticated schemes for differential grading (e.g., where a larger difference score is assigned, say, for blue and white than for blue and black).

*“What about missing values?”* In general, if the value of a given attribute  $A$  is missing in tuple  $X_1$  and/or in tuple  $X_2$ , we assume the maximum possible difference. Suppose that each of the attributes has been mapped to the range  $[0, 1]$ . For nominal attributes, we take the difference value to be 1 if either one or both of the corresponding values of  $A$  are missing. If  $A$  is numeric and missing from both tuples  $X_1$  and  $X_2$ , then the difference is also taken to be 1. If only one value is missing and the other (which we will call  $v'$ ) is present and normalized, then we can take the difference to be either  $|1 - v'|$  or  $|0 - v'|$  (i.e.,  $1 - v'$  or  $v'$ ), whichever is greater.

*“How can I determine a good value for  $k$ , the number of neighbors?”* This can be determined experimentally. Starting with  $k = 1$ , we use a test set to estimate the error rate of the classifier. This process can be repeated each time by incrementing  $k$  to allow for one more neighbor. The  $k$  value that gives the minimum error rate may be selected. In general, the larger the number of training tuples, the larger the value of  $k$  will be (so that classification and numeric prediction decisions can be based on a larger portion of the stored tuples). As the number of training tuples approaches infinity and  $k = 1$ , the

error rate can be no worse than twice the Bayes error rate (the latter being the theoretical minimum). If  $k$  also approaches infinity, the error rate approaches the Bayes error rate.

Nearest-neighbor classifiers use distance-based comparisons that intrinsically assign equal weight to each attribute. They therefore can suffer from poor accuracy when given noisy or irrelevant attributes. The method, however, has been modified to incorporate attribute weighting and the pruning of noisy data tuples. The choice of a distance metric can be critical. The Manhattan (city block) distance (Section 2.4.4), or other distance measurements, may also be used.

Nearest-neighbor classifiers can be extremely slow when classifying test tuples. If  $D$  is a training database of  $|D|$  tuples and  $k = 1$ , then  $O(|D|)$  comparisons are required to classify a given test tuple. By presorting and arranging the stored tuples into search trees, the number of comparisons can be reduced to  $O(\log(|D|))$ . Parallel implementation can reduce the running time to a constant, that is,  $O(1)$ , which is independent of  $|D|$ .

Other techniques to speed up classification time include the use of *partial distance* calculations and *editing* the stored tuples. In the **partial distance** method, we compute the distance based on a subset of the  $n$  attributes. If this distance exceeds a threshold, then further computation for the given stored tuple is halted, and the process moves on to the next stored tuple. The **editing** method removes training tuples that prove useless. This method is also referred to as **pruning** or **condensing** because it reduces the total number of tuples stored.

## 9.5.2 Case-Based Reasoning

**Case-based reasoning** (CBR) classifiers use a database of problem solutions to solve new problems. Unlike nearest-neighbor classifiers, which store training tuples as points in Euclidean space, CBR stores the tuples or “cases” for problem solving as complex symbolic descriptions. Business applications of CBR include problem resolution for customer service help desks, where cases describe product-related diagnostic problems. CBR has also been applied to areas such as engineering and law, where cases are either technical designs or legal rulings, respectively. Medical education is another area for CBR, where patient case histories and treatments are used to help diagnose and treat new patients.

When given a new case to classify, a case-based reasoner will first check if an identical training case exists. If one is found, then the accompanying solution to that case is returned. If no identical case is found, then the case-based reasoner will search for training cases having components that are similar to those of the new case. Conceptually, these training cases may be considered as neighbors of the new case. If cases are represented as graphs, this involves searching for subgraphs that are similar to subgraphs within the new case. The case-based reasoner tries to combine the solutions of the neighboring training cases to propose a solution for the new case. If incompatibilities arise with the individual solutions, then backtracking to search for other solutions may be necessary. The case-based reasoner may employ background knowledge and problem-solving strategies to propose a feasible combined solution.