

String matching is the problem of finding strings that refer to the same real-world entity. For example, the string David Smith in one database may refer to the same person as David R. Smith in another database. Similarly, the strings 1210 W. Dayton St, Madison WI and 1210 West Dayton, Madison WI 53706 refer to the same physical address.

String matching plays a critical role in many data integration tasks, including schema matching, data matching, and information extraction. Consequently, in this chapter we examine this problem in depth. Section 4.1 defines the string matching problem. Section 4.2 describes popular similarity measures that can be used to compute a similarity score between any two given strings. Finally, Section 4.3 discusses how to efficiently apply such a measure to match a large number of strings.

4.1 Problem Description

The problem we address in this chapter is the following. Given two sets of strings X and Y , we want to find all pairs of strings (x, y) , where $x \in X$ and $y \in Y$, such that x and y refer to the same real-world entity. We refer to such pairs as *matches*. Figure 4.1(a-b) shows two example databases representing sets of persons, and Figure 4.1(c) shows the matches between them. For example, the first match (x_1, y_1) states that strings Dave Smith and David D. Smith refer to the same real-world person.

Solving the matching problem raises two major challenges: accuracy and scalability. Matching strings accurately is difficult because strings that refer to the same real-world entity are often very different. The reasons for appearing different include typing and OCR errors (e.g., David Smith is misspelled as Davod Smith), different formatting conventions (10/8/2009 vs. Oct 8, 2009), custom abbreviation, shortening of strings or omission (Daniel Walker Herbert Smith vs. Dan W. Smith), different names or nicknames (William Smith vs. Bill Smith), and shuffling parts of the string (Dept. of Computer Science, UW-Madison vs. Computer Science Dept., UW-Madison). Additionally, in some cases the data source(s) do not contain enough information to determine whether two strings refer to the same entity (e.g., trying to decide whether two authors mentioned in two different publications are the same person).

To address the accuracy challenge, a common solution is to define a *similarity measure* s that takes two strings x and y and returns a score in the range $[0, 1]$. The intention is that the higher the score, the more likely that x and y match. We then say that x and y match if $s(x, y) \geq t$, where t is a prespecified threshold. Many similarity measures have been proposed, and we discuss the main ones in Section 4.2.

<u>Set X</u>	<u>Set Y</u>	<u>Matches</u>
x_1 =Dave Smith	y_1 =David D. Smith	(x_1, y_1)
x_2 =Joe Wilson	y_2 =Daniel W. Smith	(x_3, y_2)
x_3 =Dan Smith		
(a)	(b)	(c)

FIGURE 4.1 An example of matching person names. Column (c) shows the matches between the databases in (a) and (b).

The second challenge is to apply the similarity measure $s(x,y)$ to a large number of strings. Since string similarity is typically applied to data entries, applying $s(x,y)$ to all pairs of strings in the Cartesian product of sets X and Y would be quadratic in the size of the data and therefore impractical. To address this challenge, several solutions have been proposed to apply $s(x,y)$ only to the most promising pairs. We discuss the key ideas of these solutions in Section 4.3.

4.2 Similarity Measures

A broad range of measures have been proposed to match strings. A similarity measure maps a pair of strings (x,y) into a number in the range $[0,1]$ such that a higher value indicates greater similarity between x and y . The terms *distance* and *cost measures* have also been used to describe the same concept, except that smaller values indicate higher similarity.

Broadly speaking, current similarity measures fall into four groups: sequence-based, set-based, hybrid, and phonetic measures. We now describe each one in turn.

4.2.1 Sequence-Based Similarity Measures

The first set of similarity measures we consider views strings as sequences of characters and computes a cost of transforming one string into the other. We begin with a basic measure called *edit distance*, and then consider several more elaborate versions.

Edit Distance

The edit distance measure, also known as Levenshtein distance, $d(x,y)$, computes the *minimal* cost of transforming string x to string y . Transforming a string is carried out using a sequence of the following operators: delete a character, insert a character, and substitute one character for another. For example, the cost of transforming the string $x = \text{David Smiths}$ to the string $y = \text{David Smith}$ is 4. The required operations are: inserting a *d* (after David), substituting *m* by *i*, substituting *i* by *m*, and deleting the last character of x , which is *s*.

It is not hard to see that the minimal cost of transforming x to y is the same as the minimal cost of transforming y to x (using in effect the same transformation). Thus, $d(x,y)$ is well-defined and symmetric with respect to both x and y .

Intuitively, edit distance tries to capture the various ways people make editing mistakes, such as inserting an extra character (e.g., Davidd instead of David) or swapping two characters (e.g., Simth instead of Smith). Hence, the smaller the edit distance, the more similar the two strings are.

The edit distance function $d(x,y)$ is converted into a similarity function $s(x,y)$ as follows:

$$s(x,y) = 1 - \frac{d(x,y)}{\max(\text{length}(x), \text{length}(y))}$$

For example, the similarity score between David Smiths and Davidd Smith is

$$s(\text{David Smiths}, \text{Davidd Smith}) = 1 - \frac{4}{\max(12, 12)} = 0.67$$

The value of $d(x,y)$ can be computed using dynamic programming. Let $x = x_1x_2 \dots x_n$ and $y = y_1y_2 \dots y_m$, where the x_i and y_j are characters. Let $d(i,j)$ denote the edit distance between $x_1x_2 \dots x_i$ and $y_1y_2 \dots y_j$ (which are the i th and j th prefixes of x and y).

The key to designing a dynamic-programming algorithm is to establish a recurrence equation that enables computing $d(i,j)$ from previously computed values of d . Figure 4.2(a) shows the recurrence equation in our case. To understand the equation, observe that we can transform string $x_1x_2 \dots x_i$ into string $y_1y_2 \dots y_j$ in one of four ways: (a) transforming $x_1x_2 \dots x_{i-1}$ into $y_1y_2 \dots y_{j-1}$, then copying x_i into y_j if $x_i = y_j$, (b) transforming $x_1x_2 \dots x_{i-1}$ into $y_1y_2 \dots y_{j-1}$, then substituting x_i with y_j if $x_i \neq y_j$, (c) deleting x_i , then transforming $x_1x_2 \dots x_{i-1}$ into $y_1y_2 \dots y_j$, and (d) transforming $x_1x_2 \dots x_i$ into $y_1y_2 \dots y_{j-1}$, then inserting y_j . The value $d(i,j)$ is the minimum of the costs of the above transformations. Figure 4.2(b) simplifies the equation in Figure 4.2(a) by merging the first two lines.

Figure 4.3 shows an example computation using the simplified equation. The figure illustrates computing the distance between the strings $x = \text{dva}$ and $y = \text{dave}$. We start with the matrix in Figure 4.3(a), where the x_i are listed on the left and the y_j at the top. Note that we have added x_0 and y_0 , two null characters at the start of x and y , to simplify the implementation. Specifically, this allows us to quickly fill in the first row and column, by setting $d(i, 0) = i$ and $d(0, j) = j$.

$$d(i, j) = \min \begin{cases} d(i-1, j-1) & \text{if } x_i = y_j \text{ // copy} \\ d(i-1, j-1) + 1 & \text{if } x_i \neq y_j \text{ // substitute} \\ d(i-1, j) + 1 & \text{// delete } x_i \\ d(i, j-1) + 1 & \text{// insert } y_j \end{cases} \quad d(i, j) = \min \begin{cases} d(i-1, j-1) + c(x_i, y_j) & \text{// copy or substitute} \\ d(i-1, j) + 1 & \text{// delete } x_i \\ d(i, j-1) + 1 & \text{// insert } y_j \end{cases}$$

$$c(x_i, y_j) = \begin{cases} 0 & \text{if } x_i = y_j \\ 1 & \text{otherwise} \end{cases}$$

FIGURE 4.2 (a) The recurrence equation for computing edit distance between strings x and y using dynamic programming and (b) a simplified form of the equation.

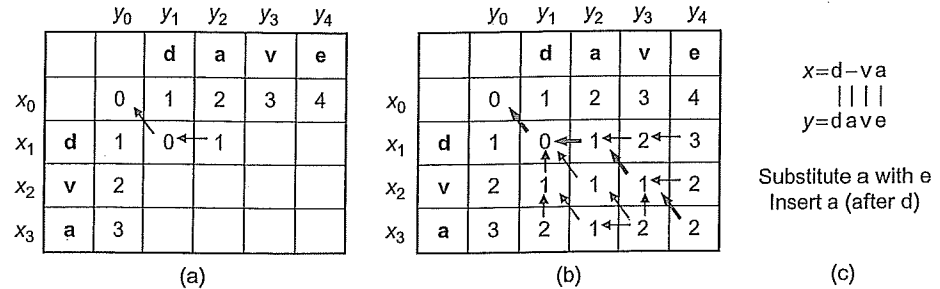


FIGURE 4.3 Computing the edit distance between dva and dave using the dynamic programming equation in Figure 4.2(b). (a) The dynamic programming matrix after filling in several cells, (b) the filled-in matrix, and (c) the sequence of edit operations that transforms dva into dave, found by following the bold arrows in part (b).

Now we can use the equation in Figure 4.2(b) to fill in the rest of the matrix. For example, we have $d(1, 1) = \min\{0 + 0, 1 + 1, 1 + 1\} = 0$. Since this value is obtained by adding 0 to $d(0, 0)$, we add an arrow pointing from cell (1, 1) to cell (0, 0). Similarly, $d(1, 2) = 1$ (see Figure 4.3(a)). Figure 4.3(b) shows the entire filled-in matrix. The edit distance between x and y can then be found to be 2, in (3, 4), the bottom rightmost cell.

In addition to computing the edit distance, the algorithm also shows the sequence of edit operations, by following the arrows. In our example, since the arrow goes "diagonal," from (3, 4) to (2, 3), we know that x_3 (character a) has been copied into or substituted with y_4 (character e). The arrow then goes "diagonal" again, from (2, 3) to (1, 2). So again x_2 (character v) has been copied into or substituted with y_3 (character v). Next, the arrow goes "horizontal," from (1, 2) to (1, 1). This means a gap - has been inserted into x and aligned with a in y (which denotes an insertion operation). This process stops when we have reached cell (0, 0). The transformation is depicted in Figure 4.3(c).

The cost of computing the edit distance is $O(|x||y|)$. In practice, the lengths of x and y often are roughly the same, and hence we often refer to the above cost as quadratic.

The Needleman-Wunch Measure

The Needleman-Wunch measure generalizes the Levenshtein distance. Specifically, it is computed by assigning a score to each *alignment* between the two input strings and choosing the score of the best alignment, that is, the maximal score. An alignment between two strings x and y is a set of *correspondences* between the characters of x and y , allowing for *gaps*. For example, Figure 4.4(a) shows an alignment between the strings $x = dva$ and $y = deeve$, where d corresponds to d, v to v, and a to e. Note that a gap of length 2 has been inserted into x , and so the two characters ee of y do not have any corresponding characters in x .

The score of an alignment is computed using a *score matrix* and a *gap penalty*. The matrix assigns a score for a correspondence between every pair of characters and therefore allows for penalizing transformations on a case-by-case basis. Figure 4.4(b) shows a sample score matrix where a correspondence between two identical characters scores 2,

d--va
| | | |
dave

Match a with e
a (after d)

(c)

Recurrence equation in matrix, and (c) the alignment in part (b).

Matrix. For example, substituted by adding 0 $d(1,2) = 1$ (see since between x

the sequence of the "diagonal," substituted with 2). So again x_2 text, the arrowed into x and stops when we

lengths of x and y quadratic.

Specifically, it input strings An alignment characters of x and seen the strings a gap of length corresponding

penalty. The letters and there- 4.4(b) shows a characters scores 2,

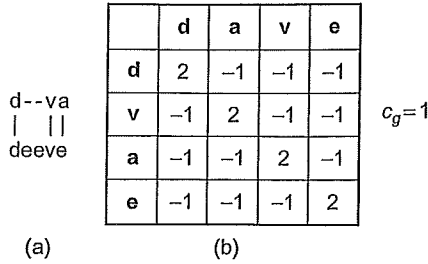


FIGURE 4.4 An example for the Needleman-Wunch function: (a) an alignment between two strings $x = dva$ and $y = deeve$, with a gap in x , and (b) a score matrix that assigns to each pair of characters a score and a gap penalty c_g .

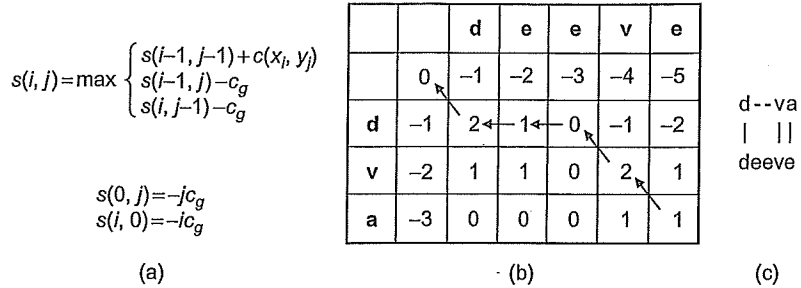


FIGURE 4.5 An example for the Needleman-Wunch function: (a) the recurrence equation for computing the similarity score using dynamic programming, (b) the dynamic-programming matrix between dva and $deeve$ using the equation in part (a), and (c) the optimal alignment between the above two strings, found by following the arrows in part (b).

and -1 otherwise. A gap of length 1 has a penalty c_g (set to 1 in Figure 4.4). A gap of length k has a linear penalty of kc_g .

The score of an alignment is then the sum of the scores of all correspondences in the alignment, minus the penalties for the gaps. For example, the score of the alignment in Figure 4.4(a) is 2 (for correspondence d-d) + 2 (for v-v) - 1 (for a-e) - 2 (penalty for the gap of length 2) = 1. This is the best alignment between x and y (that is, the one with the highest score) and therefore is the Needleman-Wunch score between the two.

As described, the Needleman-Wunch measure generalizes the Levenshtein distance in three ways. First, it computes similarity scores instead of distance values. Second, it generalizes edit costs into a score matrix, thus allowing for more fine-grained score modeling. For example, the letter o and the number 0 are often confused in practice (e.g., by OCR systems). Hence, a correspondence between these two should have a higher score than one between a and 0, say. As another example, when matching bioinformatic sequences, different pairs of amino acids may have different semantic distances. Finally, the Needleman-Wunch measure generalizes insertion and deletion into gaps and generalizes the cost of such operations from 1 to an arbitrary penalty c_g .

The Needleman-Wunch score $s(x, y)$ is computed with dynamic programming, using the recurrence equation shown in Figure 4.5(a). We note three differences between the

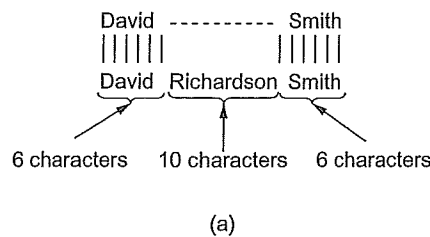
algorithm used here and the one used for computing the edit distance (Figure 4.2(b)). First, here we compute the *max* instead of *min*. Second, we use the score matrix $c(x_i, y_j)$ in the recurrence equation instead of the unit costs of edit operations. Third, we use the gap cost c_g instead of the unit gap cost. Note that when initializing this matrix, we must set $s(i, 0) = -ic_g$ and $s(0, j) = -jc_g$ (instead of i and j as in the edit distance case).

Figure 4.5(b) shows the fully filled-in matrix for $x = \text{dva}$ and $y = \text{deeve}$, using the score matrix and gap penalty in Figure 4.4(b). Figure 4.5(c) shows the best alignment, found by following the arrows in the matrix of Figure 4.5(b).

The Affine Gap Measure

The affine gap measure is an extension of the Needleman-Wunch measure that handles longer gaps more gracefully. Consider matching $x = \text{David Smith}$ with $y = \text{David R. Smith}$. The Needleman-Wunch measure can do this match very well, by opening a gap of length 2 in x , right after David, then aligning the gap with R. However, consider matching the same $x = \text{David Smith}$ with $y' = \text{David Richardson Smith}$, as shown in Figure 4.6(a). Here the gap between the two strings is 10 characters long. Needleman-Wunch does not match very well because the cost of the gap is too high. For example, assume that each character correspondence has a score of 2 and that c_g is 1; then the score of the above alignment under Needleman-Wunch is $6 \cdot 2 - 10 = 2$.

In practice, gaps tend to be longer than one character. Hence, assigning a uniform penalty to each character in the gap in a sense will unfairly punish long gaps. The affine gap measure addresses this problem by distinguishing between the cost of *opening* a gap and the cost of *continuing* the gap. Formally, the measure assigns to each gap of length k a cost $c_o + (k - 1)c_r$, where c_o is the cost of opening the gap (i.e., the cost of the very first character of the gap), and c_r is the cost of continuing the gap (i.e., the cost of the remaining $k - 1$ characters of the gap). Cost c_r is less than c_o , thereby lessening the penalty of a long gap. Continuing with the example in Figure 4.6(a), if $c_o = 1$ and $c_r = 0.5$, the score of the alignment is $6 \cdot 2 - 1 - 9 \cdot 0.5 = 6.5$, much higher than the score 2 obtained under Needleman-Wunch.



$$s(i, j) = \max \{ M(i, j), I_x(i, j), I_y(i, j) \}$$

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + c(x_i, y_j) \\ I_x(i-1, j-1) + c(x_i, y_j) \\ I_y(i-1, j-1) + c(x_i, y_j) \end{cases}$$

$$I_x(i, j) = \max \begin{cases} M(i-1, j) - c_o \\ I_x(i-1, j) - c_r \end{cases}$$

$$I_y(i, j) = \max \begin{cases} M(i, j-1) - c_o \\ I_y(i, j-1) - c_r \end{cases}$$

FIGURE 4.6 (a) An example of two strings where there is a long gap, (b) the recurrence equations for the affine gap measure.

4.2(b)). First, $c(x_i, y_j)$ in the gap cost set $s(i, 0) =$

the score found by

that handles R. Smith. The of length 2 in the same $x =$ gap between well because correspondence Needleman-

ng a uniform os. The affine opening a gap p of length k the very first f the remain- e penalty of a).5, the score tained under

.j))
 y_j
 y_j
 y_j

or the affine gap

Figure 4.6(b) shows the recurrence equations for the affine gap measure. Deriving these equations is rather involved. Since we now penalize opening a gap and continuing a gap differently, in every stage, for each cell (i, j) of the dynamic-programming matrix we keep track of three values:

- $M(i, j)$: the best score between $x_1 \dots x_i$ and $y_1 \dots y_j$ given that x_i is aligned to y_j
- $I_x(i, j)$: the best score given that x_i is aligned to a gap
- $I_y(i, j)$: the best score given that y_j is aligned to a gap

The best score for the cell, $s(i, j)$, is then the maximum of these three scores.

In deriving the above recurrence equations, we make the following assumption about the cost function. We assume that in an optimal alignment, we will never have an insertion followed directly by a deletion, or vice versa. This means we will never have the situation depicted in Figure 4.7(a) or Figure 4.7(b). We can guarantee this property by setting the cost $-(c_o + c_r)$ to be lower than the lowest mismatch score in the score matrix. Under such conditions, the alignment in Figure 4.7(c) will have a higher score than those to its left.

We now explain how to derive the equations for $M(i, j)$, $I_x(i, j)$, and $I_y(i, j)$ in Figure 4.6(b). Figure 4.8 explains how to derive the equation for $M(i, j)$. This equation considers the case where x_i is aligned with y_j (Figure 4.8(a)). Thus, $x_1 \dots x_{i-1}$ is aligned with $y_1 \dots y_{j-1}$. This can only happen in one of three ways, as shown in Figures 4.8(b)–(d): x_{i-1} is

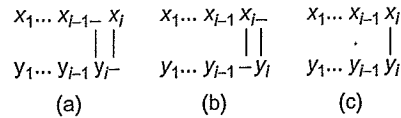


FIGURE 4.7 By setting the gap penalties and the score matrix appropriately, the alignment in (c) will always score higher than those in (a) and (b).

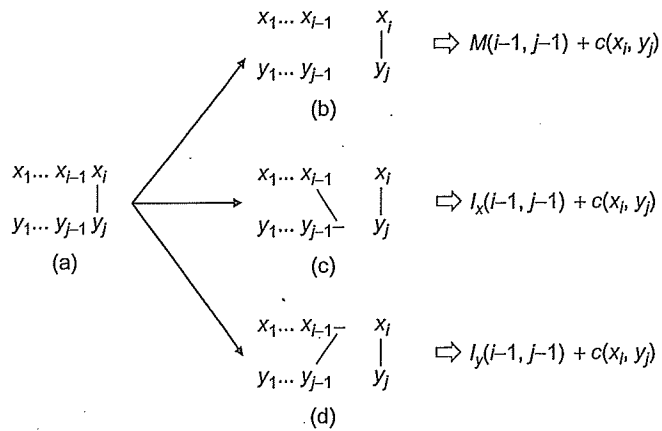


FIGURE 4.8 The derivation of the equation for $M(i, j)$ in Figure 4.6(b).

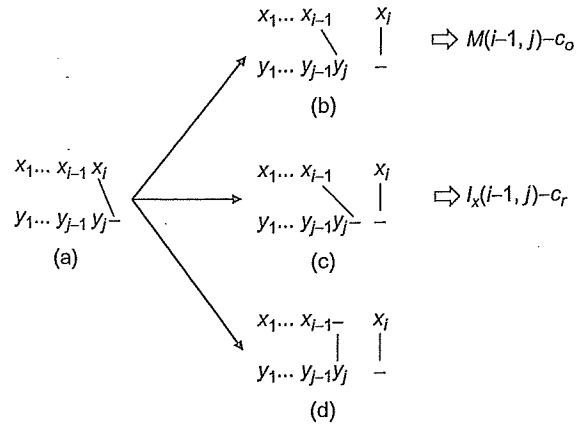


FIGURE 4.9 The derivation of the equation for $I_x(i, j)$ in Figure 4.6(b).

aligned with y_{j-1} , x_{i-1} is aligned into a gap, or y_{j-1} is aligned into a gap. These three ways give rise to the three cases in the equation for $M(i, j)$ in Figure 4.6(b).

Figure 4.9 shows how to derive the equation for $I_x(i, j)$ in Figure 4.6(b). This equation considers the case where x_i is aligned into a gap (Figure 4.9(a)). This can only happen in one of three ways, as shown in Figure 4.9(b)–(d): x_{i-1} is aligned with y_j , x_{i-1} is aligned into a gap, or y_j is aligned into a gap. The first two ways give rise to the two cases shown in Figure 4.6(b). The third case cannot happen, because it means an insertion followed directly by a deletion. This violates the assumption we described earlier. The equation for $I_y(i, j)$ in Figure 4.6(b) is derived in a similar fashion. The complexity of computing the affine gap measure remains $O(|x||y|)$.

The Smith-Waterman Measure

The previous measures consider *global* alignments between the input strings. That is, they attempt to match all characters of x with all characters of y .

Global alignments may not be well suited for some cases. For example, consider the two strings Prof. John R. Smith, Univ of Wisconsin and John R. Smith, Professor. The similarity score based on global alignments will be relatively low. In such a case, what we really want is to find two substrings of x and y that are most similar, and then return the score between the substrings as the score for x and y . For example, here we would want to identify John R. Smith to be the most similar substrings of the above two strings. This means matching the above two strings by ignoring certain prefixes (e.g., Prof.) and suffixes (e.g., Univ of Wisconsin in x and Professor in y). We call this *local alignment*.

The Smith-Waterman measure is designed to find matching substrings by introducing two key changes to the Needleman-Wunch measure. First, the measure allows the match to restart at any position in the strings (no longer limited to just the first position). The restart is captured by the first line of the recurrence equation in Figure 4.10(a). Intuitively, if the global match dips below zero, then this line has the effect of *ignoring the prefix*

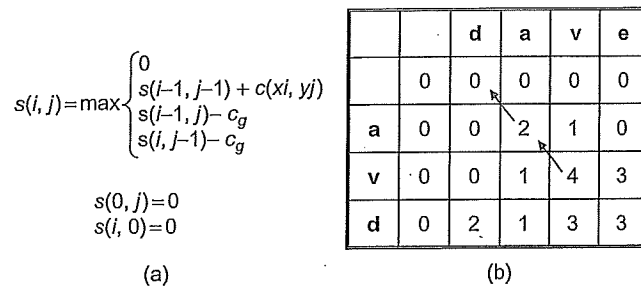


FIGURE 4.10 An example for the Smith-Waterman measure: (a) the recurrence equation for computing the similarity score using dynamic programming and (b) the dynamic-programming matrix between avd and dave using the equation in part (a).

and restarting the match. Similarly, note that all values in the first row and column of the dynamic-programming matrix are zeros, instead of $-ic_g$ and $-jc_g$ as in the Needleman-Wunsch case. Applying this recurrence equation to the strings avd and dave produces the matrix in Figure 4.10(b).

The second key change is that after computing the matrix using the recurrence equation, the algorithm starts retracing the arrows from the largest value in the matrix (4 in our example) rather than starting from the lower-right corner (3 in the matrix). This change effectively ignores suffixes if the match they produce is not optimal. Retracing ends when we meet a cell with value 0, which corresponds to the start of the alignment. In our example we can read out the best local alignment between avd and dave, which is av.

The Jaro Measure

The Jaro measure was developed mainly to compare short strings, such as first and last names. Given two strings x and y , we compute their Jaro score as follows.

- Find the common characters x_i and y_j such that $x_i = y_j$ and $|i - j| \leq \min \{|x|, |y|\} / 2$.
Intuitively, common characters are those that are identical and are positionally “close to one another.” It is not hard to see that the number of common characters x_i in x is the same as the number of common characters y_j in y . Let this number be c .
- Compare the i th common character of x with the i th common character of y . If they do not match, then we have a transposition. Let the number of transpositions be t .
- Compute the Jaro score as

$$jaro(x, y) = 1/3[c/|x| + c/|y| + (c - t/2)/c]$$

As an example, consider $x = jon$ and $y = john$. We have $c = 3$. The common character sequence in x is jon , and so is the sequence in y . Hence, there is no transposition, and $t = 0$. Thus, $jaro(x, y) = 1/3(3/3 + 3/4 + 3/3) = 0.917$. In contrast, the similarity score according to edit distance would be 0.75.

Now suppose $x = \text{jon}$ and $y = \text{ojhn}$. Here the common character sequence in x is jon and the common character sequence in y is ojn . Thus $t = 2$, and $\text{jaro}(x, y) = 1/3(3/3 + 3/4 + (3 - 2/2)/3) = 0.81$.

The cost of computing the Jaro distance is $O(|x||y|)$, due to the cost of finding common characters.

The Jaro-Winkler Measure

The Jaro-Winkler measure is designed to capture cases where two strings have a low Jaro score, but share a prefix and thus are still likely to match. Specifically, the measure introduces two parameters: PL , which is the length of the longest common prefix between the two strings, and PW , which is the weight to give the prefix. The measure is computed using the following formula:

$$\text{jaro-winkler}(x, y) = (1 - PL * PW) * \text{jaro}(x, y) + PL * PW$$

4.2.2 Set-Based Similarity Measures

The previous class of measures considers strings as sequences of characters. We now describe similarity measures that view strings as sets or multi-sets of *tokens*, and use set-related properties to compute similarity scores.

There are many ways to generate tokens from the input strings. A common method is to consider the *words* in the string (as delimited by the space character) and possibly stem the words. Common stop words (e.g., the, and, of) are typically excluded. For example, given the string *david smith*, we may generate the set of tokens $\{\text{david}, \text{smith}\}$.

Another common type of token is *q-grams*, which are substrings of length q that are present in the string. For example, the set of all 3-grams of *david smith* is $\{\#\#\text{d}, \#\text{da}, \text{dav}, \text{avi}, \dots, \text{ith}, \text{h}\#\#, \text{th}\#\}$. Note that we have appended the special character $\#$ to the start and the end of the string, to handle 3-grams in these positions.

We discuss several set-based similarity measures. The bibliographic notes contain pointers to others that have been proposed in the literature. In our discussion we often assume *sets* of tokens, but we note that these measures have also been considered for the multi-set case, and what we discuss below can be generalized to that case.

The Overlap Measure

Let B_x and B_y be the sets of tokens generated for strings x and y , respectively. The overlap measure returns the number of common tokens $O(x, y) = |B_x \cap B_y|$.

Consider $x = \text{dave}$ and $y = \text{dav}$; then the set of all 2-grams of x is $B_x = \{\#\text{d}, \text{da}, \text{av}, \text{ve}, \text{e}\#\}$ and the set of all 2-grams of y is $B_y = \{\#\text{d}, \text{da}, \text{av}, \text{v}\#\}$. So $O(x, y) = 3$.

The Jaccard Measure

Continuing with the above notation, the Jaccard similarity score between two strings x and y is $J(x, y) = |B_x \cap B_y| / |B_x \cup B_y|$.

in x is jon and
 $1/3(3/3 + 3/4 +$
 iding common

ave a low Jaro
 measure intro-
 ix between the
 e is computed

cters. We now
 kens, and use

non method is
 l possibly stem
 . For example,

ngth q that are
 $d, \#da, dav, avi,$
 e start and the

notes contain
 sion we often
 sidered for the

ly. The overlap

$d, da, av, ve, e\#$)

o strings x and

$$\begin{array}{l} x = aab \Rightarrow B_x = \{a, a, b\} \\ y = ac \Rightarrow B_y = \{a, c\} \\ z = a \Rightarrow B_z = \{a\} \end{array} \quad \begin{array}{l} tf(a, x) = 2 \\ tf(b, x) = 1 \\ \dots \\ tf(c, z) = 0 \end{array} \quad \begin{array}{l} idf(a) = 3/3 = 1 \\ idf(b) = 3/1 = 3 \\ idf(c) = 3/1 = 3 \end{array}$$

	a	b	c
v_x	2	3	0
v_y	3	0	3
v_z	3	0	0

(a)

(b)

(c)

FIGURE 4.11 In the TF/IDF measure (a) strings are converted into bags of terms, (b) TF and IDF scores of the terms are computed, then (c) these scores are used to compute feature vectors.

Again, consider $x = \text{dave}$ with $B_x = \{\#d, da, av, ve, e\#$, and $y = \text{dav}$ with $B_y = \{\#d, da, av, v\#$. Then $J(x, y) = 3/6$.

The TF/IDF Measure

This measure employs the notion of TF/IDF score commonly used in information retrieval (IR) to find documents that are relevant to keyword queries. The intuition underlying the TF/IDF measure is that two strings are similar if they share distinguishing terms. For example, consider the three strings $x = \text{Apple Corporation, CA}$, $y = \text{IBM Corporation, CA}$, and $z = \text{Apple Corp}$. The edit distance and Jaccard measure would match x with y as $s(x, y)$ is higher than $s(x, z)$. However, the TF/IDF measure is able to recognize that Apple is a distinguishing term, whereas Corporation and CA are far more common, and thus would correctly match x with z .

When discussing the TF/IDF measure, we assume that the pair of strings being matched is taken from a *collection of strings*. Figure 4.11(a) shows a tiny such collection of three strings, $x = aab$, $y = ac$, and $z = a$. We convert each string into a bag of terms. Using IR terminology, we refer to such a bag of terms as a *document*. For example, we convert string $x = aab$ into document $B_x = \{a, a, b\}$.

We now compute *term frequency (TF) scores* and *inverse document frequency (IDF) scores* as follows:

- For each term t and document d , we compute the term frequency $tf(t, d)$ to be the number of times t occurs in d . For example, since a occurs twice in B_x , we have $tf(a, x) = 2$.
- For each term t , we compute the inverse document frequency $idf(t)$ to be the total number of documents in the collection divided by the number of documents that contain t (variations of this definition are also commonly used, see below). For example, since a appears in all three documents in Figure 4.11(a), we have $idf(a) = 3/3 = 1$. A higher value of IDF means that the occurrence of the term is more distinguishing.

Figure 4.11(b) shows the TF/IDF scores for the tiny example in Figure 4.11(a).

Next, we convert each document d into a feature vector v_d . The intuition is that two documents will be similar if their corresponding vectors are close to each other. The vector of d has a feature $v_d(t)$ for each term t , and the value of $v_d(t)$ is a function of the TF and IDF scores. Vector v_d thus has as many features as the number of terms in the collection.

Figure 4.11(c) shows the three vectors $v_x, v_y,$ and v_z for the three documents $B_x, B_y,$ and $B_z,$ respectively. In the figure, we use a relatively simple score: $v_d(t) = tf(t, d) \cdot idf(t)$. Thus, the score for feature a of v_x is $v_x(a) = 2 \cdot 1 = 2$, and so on.

Now we are ready to compute the TF/IDF similarity score between any two strings p and q . Let T be the set of all terms in the collection. Then conceptually the vectors v_p and v_q (of the strings p and q) can be viewed as vectors in the $|T|$ -dimensional space where each dimension corresponds to a term. The TF/IDF score between p and q can be computed as the cosine of the angle between these two vectors:

$$s(p, q) = \frac{\sum_{t \in T} v_p(t) \cdot v_q(t)}{\sqrt{\sum_{t \in T} v_p(t)^2} \cdot \sqrt{\sum_{t \in T} v_q(t)^2}} \quad (4.1)$$

For example, the TF/IDF score between the two strings x and y in Figure 4.11(a) is $\frac{2 \cdot 3}{\sqrt{2^2 + 3^2} \cdot \sqrt{3^2 + 3^2}} = 0.39$, using the vectors v_x and v_y in Figure 4.11(c).

It is not difficult to see from Equation 4.1 that the TF/IDF similarity score between two strings p and q is high if they share many frequent terms (that is, terms with high TF scores), unless these terms also commonly appear in other strings in the collection (in which case the terms have low IDF scores). Using the IDF component, the TF/IDF similarity score can effectively discount the importance of such common terms.

In the above example, we assumed $v_d(t) = tf(t, d) \cdot idf(t)$. This means that if we “double” the number of occurrences of t in document d , then $v_d(t)$ will also double. In practice this is found to be excessive: doubling the number of occurrences of t should increase $v_d(t)$ but not double it. One way of addressing this is to “dampen” the TF and IDF components by a logarithmic factor. Specifically, we can take

$$v_d(t) = \log(tf(t, d) + 1) \cdot \log(idf(t))$$

(In fact, $\log(idf(t))$ itself is also commonly referred to as the inverse document frequency.) In addition, the vector v_d is often normalized to length 1, by setting

$$v_d(t) = v_d(t) / \sqrt{\sum_{t \in T} v_d(t)^2}$$

This way, computing the TF/IDF similarity score $s(p, q)$ as in Equation 4.1 reduces to computing the dot product between the two normalized vectors v_p and v_q .

4.2.3 Hybrid Similarity Measures

We now describe several similarity measures that combine the benefits of sequence-based and set-based methods.

The Generalized Jaccard Measure

Recall that the Jaccard measure considers the number of overlapping tokens in the input strings x and y . However, a token from x and a token from y have to be identical in order to be considered in the overlap set, which may be restrictive in some cases.

$B_x, B_y,$ and $B_z,$
(t). Thus, the

two strings p
ectors v_p and
e where each
computed as

(4.1)

re 4.11(a) is

between two
h TF scores),
n which case
ity score can

s that if we
so double. In
s of t should
" the TF and

at frequency.)

1 reduces to

uence-based

in the input
al in order to

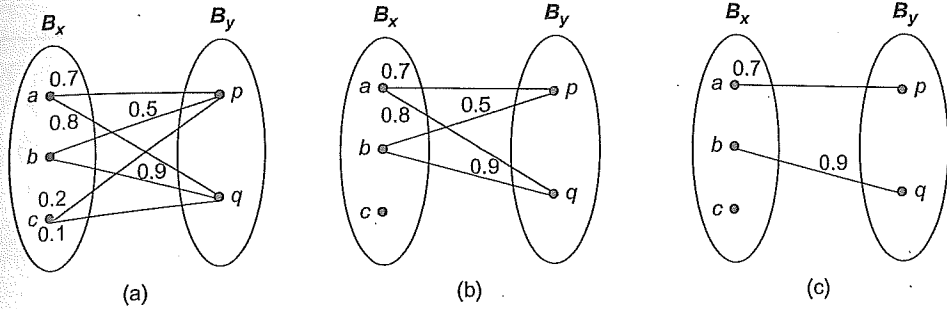


FIGURE 4.12 An example of computing the generalized Jaccard measure.

For example, consider matching the names of the nodes of two taxonomies describing divisions of companies. Each node is described by a string, such as Energy and Transportation and Transportation, Energy, and Gas. The Jaccard measure is a promising candidate for matching such strings because intuitively two nodes are similar if their names share many tokens (e.g., energy and transportation). However, in practice tokens are often misspelled, such as energy vs. eneryg. The generalized Jaccard measure will enable matching in such cases.

As with the Jaccard measure, we begin by converting the string x into a set of tokens $B_x = \{x_1, x_2, \dots, x_n\}$ and string y into a set of tokens $B_y = \{y_1, y_2, \dots, y_m\}$. Figure 4.12(a) shows two such strings x and y , with $B_x = \{a, b, c\}$ and $B_y = \{p, q\}$.

The next three steps will determine the set of pairs of tokens that are considered in the "softened" overlap set. First, let s be a similarity measure that returns values in the range $[0, 1]$. We apply s to compute a similarity score for each pair $(x_i \in B_x, y_j \in B_y)$. Continuing with the above example, Figure 4.12(a) shows all six such scores.

Second, we keep only those scores that equal or exceed a given threshold α . Figure 4.12(b) shows the remaining scores in our example, given $\alpha = 0.5$. The sets B_x and B_y , together with the edges that denote the remaining scores, form a bipartite graph G . In the third step we find the maximum-weight matching M in the bipartite graph G . Figure 4.12(c) shows this matching in our example. The total weight of this matching is $0.7 + 0.9 = 1.6$.

Finally, we return the normalized weight of M as the generalized Jaccard score between x and y . To normalize, we divide the weight by the sum of the number of edges in M and the number of "unmatched" elements in B_x and B_y . This sum is $|M| + (|B_x| - |M|) + (|B_y| - |M|) = |B_x| + |B_y| - |M|$. Formally,

$$GJ(x, y) = \frac{\sum_{(x_i, y_j) \in M} s(x_i, y_j)}{|B_x| + |B_y| - |M|}$$

Continuing with the above example, the generalized Jaccard score is $(0.7 + 0.9) / (3 + 2 - 2) = 0.53$.

The measure $GJ(x, y)$ is guaranteed to be between 0 and 1. It is a natural generalization of the Jaccard measure $J(x, y)$: if we constrain the elements of B_x and B_y to match only if they are identical, $GJ(x, y)$ reduces to $J(x, y)$. We discuss how to compute $GJ(x, y)$ efficiently later in Section 4.3.

The Soft TF/IDF Similarity Measure

This measure is similar in spirit to the generalized Jaccard measure, except that it uses the TF/IDF measure instead of the Jaccard measure as the “higher-level” similarity measure.

Consider an example with three strings $x = \text{Apple Corporation, CA}$, $y = \text{IBM Corporation, CA}$, and $z = \text{Aple Corp.}$ To match these strings, we would like to use the TF/IDF measure so that we can discount common terms such as Corporation and CA. Unfortunately, in this case the TF/IDF measure does not help us match x with z , because the term Apple in x does not match the misspelled term Aple in z . Thus, x and z do not share any term. As with the generalized Jaccard measure, we would like to “soften” the requirement that Apple and Aple match exactly, and instead require that they be similar to each other.

We compute the soft TF/IDF measure as follows. As with the TF/IDF measure, given two strings x and y , we create the two documents B_x and B_y . Figure 4.13(b) shows the documents created for the strings in Figure 4.13(a).

Next, we compute $close(x, y, k)$ to be the set of all terms in B_x that have at least one close term in B_y . Specifically, $close(x, y, k)$ is the set of terms $t \in B_x$ such that there exists a term $u \in B_y$ that satisfies $s'(t, u) \geq k$, where s' is a basic similarity measure (e.g., Jaro-Winkler) and k is a prespecified threshold. Continuing with our example in Figure 4.13(b), suppose $s'(a, a) = 1$, $s'(a, a') = 0.7$, \dots , as shown in the figure. Then $close(x, y, 0.75) = \{a, b, c\}$. Note that $d \in B_x$ is excluded because the closest term to d in B_y is d' , but d' is still too far from d (at a similarity score of 0.6).

In the final step, we compute $s(x, y)$ as in the traditional TF/IDF score, but giving a weight to each component of the TF/IDF formula according to the similarity score produced by s' . Specifically, let v_x and v_y be the feature vectors for x and y , as explained when we discussed the TF/IDF similarity measure. The vectors v_x and v_y are normalized to length 1, so that the traditional TF/IDF score can be computed as the dot product of the two vectors. Then we have

$$s(x, y) = \sum_{t \in close(x, y, k)} v_x(t) \cdot v_y(u_*) \cdot s'(t, u_*)$$

(a) $x = \text{abcd}$
 $y = \text{aa'b'cd'}$

(b) $B_x = \{a, b, c, d\}$

 $B_y = \{a, a', b', c, d'\}$
 $close(x, y, 0.75) = \{a, b, c\}$

(c) $s(x, y) = v_x(a) \cdot v_y(a) \cdot 1 +$
 $v_x(b) \cdot v_y(b') \cdot 0.8 +$
 $v_x(c) \cdot v_y(c) \cdot 1$

FIGURE 4.13 An example of computing soft TF/IDF similarity score for two strings x and y .

where $u_* \in B_y$ is the term that maximizes $s'(t, u)$ for all $u \in B_y$. Figure 4.13(c) shows how to compute $s(x, y)$ given the examples in Figures 4.13(a-b).

The Monge-Elkan Similarity Measure

The Monge-Elkan similarity measure can be effective for domains in which more control is needed over the similarity measure. To apply this measure to two strings x and y , first we break them into multiple substrings, say $x = A_1 \cdots A_n$ and $y = B_1 \cdots B_m$, where the A_i and B_j are substrings. Next, we compute

$$s(x, y) = \frac{1}{n} \sum_{i=1}^n \max_{j=1}^m s'(A_i, B_j)$$

where s' is a secondary similarity measure, such as Jaro-Winkler. To illustrate the above formula, suppose $x = A_1A_2$ and $y = B_1B_2B_3$. Then

$$s(x, y) = \frac{1}{2} [\max\{s'(A_1, B_1), s'(A_1, B_2), s'(A_1, B_3)\} + \max\{s'(A_2, B_1), s'(A_2, B_2), s'(A_2, B_3)\}]$$

Note that we ignore the order of the matching of the substrings and only consider the best match for the substrings of x in y . Furthermore, we can customize the secondary similarity measure s' to a particular application.

For example, consider matching strings $x = \text{Comput. Sci. and Eng. Dept., University of California, San Diego}$ and $y = \text{Department of Computer Science, Univ. Calif., San Diego}$. To employ the Monge-Elkan measure, first we break x and y into substrings such as *Comput.*, *Sci.*, and *Computer*. Next we must design a secondary similarity measure s' that works well for such substrings. In particular, it is clear that s' must handle matching abbreviations well. For example, s' may decide that if a substring A_i is a prefix of a substring B_j , such as *Comput.* and *Computer*, then they match, that is, their similarity score is 1.

4.2.4 Phonetic Similarity Measures

The similarity measures we have discussed so far match strings based on their *appearance*. In contrast, phonetic measures match strings based on their *sound*. These measures have been especially effective in matching names, since names are often spelled in different ways that sound the same. For example, Meyer, Meier, and Mire sound the same, as do Smith, Smithe, and Smythe. We describe the Soundex similarity measure, which is the most commonly used. We mention extensions of the basic Soundex measure in the bibliographic notes.

Soundex is used primarily to match surnames. It maps a surname x into a four-character code that captures the sound of the name. Two surnames are deemed similar if they share the same code. Mapping x to a code proceeds as follows. We use $x = \text{Ashcraft}$ as a running example in our description.

1. Keep the first letter of x as the first letter of the code. The first letter of the code for Ashcraft is A. The following steps are performed on the rest of the string x .

2. Remove all occurrences of W and H. Go over the remaining letters and replace them with digits as follows: replace B, F, P, V with 1; C, G, J, K, Q, S, X, Z with 2; D, T with 3; L with 4; M, N with 5; and R with 6. Note that we do not replace the vowels A, E, I, O, U, and Y. Continuing with our example, we convert Ashcraft into A226a13.
3. Replace each sequence of identical digits by the digit itself. So A226a13 becomes A26a13.
4. Drop all the nondigit letters (except the first one, of course). Then return the first four letters as the Soundex code. So A26a13 becomes A2613, and the corresponding Soundex code is A261.

Thus the Soundex code is always a letter followed by three digits, padded by 0 if there are not enough digits. For example, the soundex code for Sue is S000.

As described, the Soundex measure in effect “hashes” similar sounding consonants (such as B, F, P, and V) into the same digit, thereby mapping similar sounding names into the same soundex code. For example, it maps both Robert and Rupert into R163.

Soundex is not perfect. For example, it fails to map the similar sounding surnames Gough and Goff, or Jawornicki and Yavornitzky (an Americanized spelling of the former), into the same code. Nevertheless, it is a useful tool that has been used widely to match and index names in applications such as census records, vital records, ship passenger lists, and geneology databases. While Soundex was designed primarily for Caucasian surnames, it has been found to work well for names of many different origins (such as those appearing in the records of the U.S. Immigration and Naturalization Services). However, it does not work as well for names of East Asian origins, because much of the discriminating power of these names resides in the vowel sounds, which the code ignores.

4.3 Scaling Up String Matching

Once we have selected a similarity measure $s(x, y)$, the next challenge is to match strings efficiently. Let X and Y be two sets of strings to be matched, and t be a similarity threshold. A naive matching solution would be as follows:

```

for each string  $x \in X$  do
  for each string  $y \in Y$  do
    if  $s(x, y) \geq t$  then return  $(x, y)$  as a matched pair
  end for
end for

```

This $O(|X||Y|)$ solution is clearly impractical for large data sets. A more commonly employed solution is based on developing a method FindCands that can quickly find the strings that *may* match a given string x . Given such a method, we employ the following algorithm:

```

for each string  $x \in X$  do
  use method FindCands to find a candidate set  $Z \subseteq Y$ 

```



```

for each string  $y \in Z$  do
    if  $s(x, y) \geq t$  then return  $(x, y)$  as a matched pair
end for
end for
    
```

This solution, often called a *blocking solution*, takes $O(|X||Z|)$ time, which is much faster than $O(|X||Y|)$ because FindCands is designed so that finding Z is inexpensive and $|Z|$ is much smaller than $|Y|$. The set Z is often called the *umbrella set* of x . It should contain all true positives (i.e., all strings in Y that can possibly match x) and as few negative positives (i.e., those strings in Y that do not match x) as possible.

Clearly, the method FindCands lies at the heart of the above solution, and many techniques have been proposed for it. These techniques are typically based on indexing or filtering heuristics. We now discuss the basic ideas that underlie several common techniques for FindCands. In the following, we explain the techniques using the Jaccard and overlap measures. Later we discuss how to extend these techniques to other similarity measures.

4.3.1 Inverted Index Over Strings

This technique first converts each string $y \in Y$ into a document and then builds an inverted index over these documents. Given a term t , we can use the index to quickly find the list of documents created from Y that contain t , and hence the strings in Y that contain t .

Figure 4.14(a) shows an example of matching the two sets X and Y . We scan the set Y to build the inverted index shown in Figure 4.14(b). For instance, the index shows that the term *area* appears in just document 5, but the term *lake* appears in two documents, 4 and 6.

Given a string $x \in X$, the method FindCands uses the inverted index to quickly locate the set of strings in Y that share at least one term with x . Continuing with the above example,

Set X
1: {lake, mendota}
2: {lake, monona, area}
3: {lake, mendota, monona, dane}
Set Y
4: {lake, monona, university}
5: {monona, research, area}
6: {lake, mendota, monona, area}

(a)

Terms in Y	ID Lists
area	5
lake	4, 6
mendota	6
monona	4, 5, 6
research	5
university	4

(b)

FIGURE 4.14 An example of using an inverted index to speed up string matching.

given $x = \{\text{lake, mendota}\}$, we use the index in Figure 4.14(b) to find and merge the ID lists for lake and mendota, to obtain the umbrella set $Z = \{4, 6\}$.

This method is clearly much better than naively matching x with all strings in Y . Nevertheless, it still suffers from several limitations. First, the inverted list of some terms (e.g., stop words) can be very long, so building and manipulating such lists are quite costly. Second, this method requires enumerating all pairs of strings that share at least one term. The set of such pairs can still be very large in practice. The techniques described below address these issues.

4.3.2 Size Filtering

This technique retrieves only the strings in Y whose size make them match candidates. Specifically, given a string $x \in X$, we infer a constraint on the size of strings in Y that can possibly match x . We use a B-tree index to retrieve only the strings that satisfy the size constraints.

To derive the constraint on the size of strings in Y , recall that the Jaccard measure is defined as follows (where $|x|$ refers to the number of tokens in x):

$$J(x, y) = |x \cap y| / |x \cup y|$$

First, we can show that

$$1/J(x, y) \geq |y|/|x| \geq J(x, y) \quad (4.2)$$

To see why, consider the case where $|y| \geq |x|$. In this case, clearly $|y|/|x| \geq 1 \geq J(x, y)$. So we only have to prove $1/J(x, y) \geq |y|/|x|$, or equivalently that $|x \cup y|/|x \cap y| \geq |y|/|x|$. This inequality is true because $|x \cup y| \geq \max\{|x|, |y|\} = |y|$ and $|x \cap y| \leq \min\{|x|, |y|\} = |x|$. The case where $|y| < |x|$ can be proven similarly.

Now let t be the prespecified similarity threshold. If x and y match, then it must be that $J(x, y) \geq t$. Together with Equation 4.2, this implies that $1/t \geq |y|/|x| \geq t$ or, equivalently,

$$|x|/t \geq |y| \geq |x| \cdot t \quad (4.3)$$

Thus, given a string $x \in X$, we know that only strings that satisfy Equation 4.3 can possibly match x .

To illustrate, consider again the string $x = \{\text{lake, mendota}\}$ (the first string in set X in Figure 4.14(a)). Suppose $t = 0.8$. Using the above equation, if $y \in Y$ matches x , we must have $2/0.8 = 2.5 \geq |y| \geq 2 \cdot 0.8 = 1.6$. We can immediately see that no string in the set Y in Figure 4.14(a) satisfies this constraint.

Exploiting the above idea, procedure FindCands builds a B-tree index over the sizes of strings in Y . Given a string $x \in X$, it uses the index to find strings in Y that satisfy Equation 4.3 and returns that set of strings as the umbrella set Z . This technique is effective when there is significant variability in the number of tokens in the strings of X and Y .

the ID lists
 are in Y .
 are terms
 are costly.
 are term.
 are below

candidates.
 that can
 the size

measure is

(4.2)
 $\geq J(x, y)$.
 $|x|$. This
 $|x|$. The

it be that
 entirely,

(4.3)
 possibly

set X in
 we must
 set Y in

the sizes of
 by Equa-
 effective
 in Y .

$x: \{\text{lake, monona, area}\}$
 x'

$y: \{\text{lake, mendota, monona, area}\}$

(a)

Set X

1: {lake, mendota}
 2: {lake, monona, area}
 3: {lake, mendota, monona, dane}

Set Y

4: {lake, monona, university}
 5: {monona, research, area}
 6: {lake, mendota, monona, area}
 7: {dane, area, mendota}

(b)

Terms in Y	ID Lists
area	5, 6, 7
lake	4, 6
mendota	6, 7
monona	4, 5, 6
research	5
university	4
dane	7

(c)

FIGURE 4.15 An example of using prefix filtering to speed up string matching.

4.3.3 Prefix Filtering

The basic idea underlying this technique is that if two sets share many terms, then large subsets of them also must share terms. Using this principle, we can reduce the number of candidate strings that may match a string x .

We first explain this technique using the overlap similarity measure and then extend it to the Jaccard measure. Suppose that x and y are strings that have an overlap of tokens $|x \cap y| \geq k$. Then it is easy to see that any subset $x' \subseteq x$ of size at least $|x| - (k - 1)$ must overlap y . For example, consider the sets $x = \{\text{lake, monona, area}\}$ and $y = \{\text{lake, mendota, monona, area}\}$ in Figure 4.15(a). We have $|x \cap y| = 3 > 2$. Thus, the subset $x' = \{\text{lake, monona}\}$ in Figure 4.15(a) overlaps y (as does any other subset of size 2 of x).

We can exploit this idea in procedure FindCands as follows. Suppose we want to find all pairs (x, y) with overlap $O(x, y) \geq k$ (recall that $O(x, y) = |x \cap y|$ is the overlap similarity measure). Given a particular set x , we construct a subset x' of size $|x| - (k - 1)$, and use an inverted index to find all sets y that overlap x' . Figures 4.15(b-c) illustrate this idea. Suppose we want to match strings in the sets X and Y of Figure 4.15(b), using $O(x, y) \geq 2$. We begin by building an inverted index over the strings in set Y , as shown in Figure 4.15(c). Next, given a string such as $x_1 = \{\text{lake, mendota}\}$, we take the “prefix” of size $|x_1| - 1$, which is {lake} in this case, and let that be the set x'_1 . We use the inverted index to find all strings in Y that contain at least one token in x'_1 . This produces the set $\{y_4, y_6\}$. Note that if we use the inverted index to find all strings in Y that contain at least one token in x , we would end up with $\{y_4, y_6, y_7\}$, a larger candidate set. Thus, restricting index lookup to just a subset of x can significantly reduce the resulting set size.

Selecting the Subset Intelligently

So far we arbitrarily selected a subset x' of x and checked its overlap with the entire set y . We can do even better by selecting a *particular subset* x' of x and checking its overlap with only a *particular subset* y' of y . Specifically, suppose we have imposed an ordering \mathcal{O} over

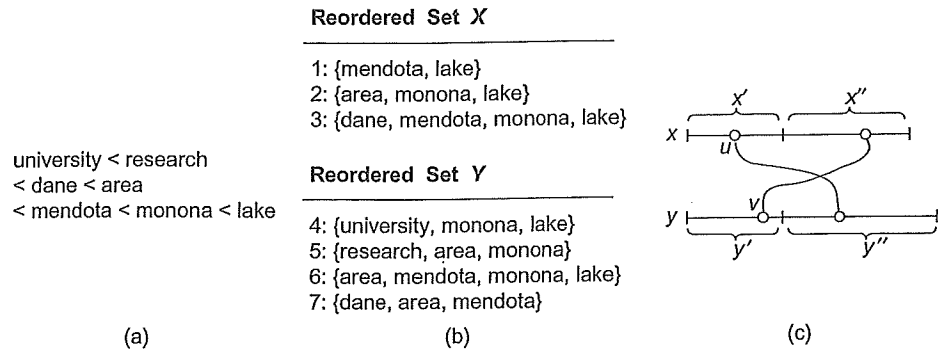


FIGURE 4.16 We can build an inverted index over only the prefixes of strings in Y , then use this index to perform string matching.

the universe of all possible terms. For example, we can order terms in increasing frequency, as computed over the union of sets X and Y of Figure 4.15(b). Figure 4.16(a) shows all terms found in $X \cup Y$ in this order.

We reorder the terms in each set $x \in X$ and $y \in Y$ according to the order \mathcal{O} . Figure 4.16(b) shows the reordered sets X and Y . Given a reordered set x , we refer to the subset x' that contains the first n terms of x as the prefix of size n (or the n th prefix) of x . For example, the 2nd prefix of $x_3 = \{\text{dane, mendota, monona, lake}\}$ is $\{\text{dane, mendota}\}$. Given the above, we can establish the following property:

Proposition 4.1. *Let x and y be two sets such that $|x \cap y| \geq k$. Let x' be the prefix of size $|x| - (k - 1)$ of x , and let y' be the prefix of size $|y| - (k - 1)$ of y . Then x' and y' overlap. \square*

Proof: Let x'' be the suffix of size $(k - 1)$ of x (see Figure 4.16(c)). Clearly, $x' \cup x'' = x$. Similarly, let y'' be the suffix of size $(k - 1)$ of y . Now suppose $x' \cap y' = \emptyset$. Then we must have $x' \cap y'' \neq \emptyset$ (otherwise, $|x \cap y| \geq k$ does not hold). So there exists an element u such that $u \in x'$ and $u \in y''$. Similarly, there exists an element v such that $v \in y'$ and $v \in x''$. Since $u \in x'$ and $v \in x''$, we have $u < v$ in the ordering \mathcal{O} . But since $v \in y'$ and $u \in y''$, we also have $v < u$ in the ordering \mathcal{O} , a clear contradiction. Thus, x' overlaps y' . \square

Given the above property, we can revise procedure FindCands as follows. Suppose again that we consider overlap of at least k (that is, $O(x, y) \geq k$).

- We reorder the terms in each set $x \in X$ and $y \in Y$ in increasing order of their frequency (as shown in Figure 4.16(b)).
- For each $y \in Y$, we create y' , the prefix of size $|y| - (k - 1)$ of y .
- We build an inverted index over all the prefixes y' . Figure 4.17(a) shows this index for our example, assuming $O(x, y) \geq 2$.
- For each $x \in X$, we create x' , the prefix of size $|x| - (k - 1)$ of x , then use the above inverted index to find all sets $y \in Y$ such that x' overlaps with y' .

Terms in Y	ID Lists
area	5, 6, 7
mendota	6
monona	4, 6
research	5
university	4
dane	7

Terms in Y	ID Lists
area	5, 6, 7
lake	4, 6
mendota	6, 7
monona	4, 5, 6
research	5
university	4
dane	7

(a)
(b)

FIGURE 4.17 The inverted indexes over (a) the prefixes of size $|y| - (k - 1)$ of all $y \in Y$ and (b) all $y \in Y$. The former is often significantly smaller than the latter in practice.

Consider for example $x = \{\text{mendota}, \text{lake}\}$, and therefore $x' = \{\text{mendota}\}$. Using mendota to look up the inverted index in Figure 4.17(a) yields y_6 . Thus, FindCands returns y_6 as the sole candidate that may match x . Note that if we check the overlap between x' and the entire y , then y_7 is also returned. Thus, checking the overlap between prefixes can reduce the size of the resulting set. In practice, this reduction can be quite significant.

It is also important to note that the size of the inverted index is much smaller. For comparison purposes, Figure 4.17(b) shows the inverted index for the entire sets $y \in Y$ (reproduced from Figure 4.15(c)). The index we create here does not contain an entry for the term lake, and its index list for mendota is also smaller than the same index list in the entire-sets index.

Applying Prefix Filtering to the Jaccard Measure

The following equation enables us to extend the prefix filtering method to the Jaccard measure.

$$J(x, y) \geq t \Leftrightarrow O(x, y) \geq \alpha = \frac{t}{1+t} \cdot (|x| + |y|) \quad (4.4)$$

The equation shows how to convert the Jaccard measure to the overlap measure, except for one detail. The threshold α as defined above is not a constant and depends on $|x|$ and $|y|$. Thus, we cannot build the inverted index over the prefixes of $y \in Y$ using α . To address this, we index the “longest” prefixes. In particular, it can be shown that we only have to index the prefixes of length $|y| - \lceil t \cdot |y| \rceil + 1$ of the $y \in Y$ to ensure that we do not miss any correct matching pairs.

4.3.4 Position Filtering

Position filtering further limits the set of candidate matches by deriving an upper bound on the size of the overlap between a pair of strings. As an example, consider the two strings $x = \{\text{dane}, \text{area}, \text{mendota}, \text{monona}, \text{lake}\}$ and $y = \{\text{research}, \text{dane}, \text{mendota}, \text{monona}, \text{lake}\}$.

Suppose we are considering $J(x, y) \geq 0.8$. In prefix filtering we will index the prefix of length $|y| - \lceil t \cdot |y| \rceil + 1$ of y , which is $y' = \{\text{research, dane}\}$ in this case (because $5 - \lceil 5 \cdot 0.8 \rceil + 1 = 2$). Similarly, the prefix of length $|x| - \lceil t \cdot |x| \rceil + 1$ of x is $x' = \{\text{dane, area}\}$. Since x' overlaps y' , in prefix filtering we will return the above pair (x, y) as a candidate pair.

However, we can do better than this. Let x'' be the rest of x , after x' , and similarly let y'' be the rest of y , after y' . Then it is easy to see that

$$O(x, y) \leq |x' \cap y'| + \min\{|x''|, |y''|\} \quad (4.5)$$

Applying this inequality to the above example, we have $O(x, y) \leq 1 + \min\{3, 3\} = 4$. However, using Equation 4.4 we have $O(x, y) \geq \frac{t}{1+t} \cdot (|x| + |y|) = \frac{0.8}{1+0.8} \cdot (5 + 5) = 4.44$. Hence, we can immediately discard the pair (x, y) from the set of candidate matches. More generally, position filtering combines the constraints from Equation 4.5 and Equation 4.4 to further reduce the set of candidate matches.

4.3.5 Bound Filtering

Bound filtering is an optimization for computing the generalized Jaccard similarity measure. Recall from Section 4.2.3 that the generalized Jaccard measure computes the normalized weight of the maximum-weight matching M in the bipartite graph connecting x and y :

$$GJ(x, y) = \frac{\sum_{(x_i, y_j) \in M} s(x_i, y_j)}{|B_x| + |B_y| - |M|}$$

In the equation, s is a secondary similarity measure, $B_x = \{x_1, x_2, \dots, x_n\}$ is the set of tokens that corresponds to x , and $B_y = \{y_1, y_2, \dots, y_m\}$ is the set that corresponds to y .

Computing $GJ(x, y)$ in a straightforward fashion would require computing the maximum-weight matching M in the bipartite graph, which can be very expensive. To address this problem, given a pair (x, y) we compute an upper bound $UB(x, y)$ and a lower bound $LB(x, y)$ on $GJ(x, y)$. FindCands uses these bounds as follows: if $UB(x, y) \leq t$, then we can ignore (x, y) as it cannot be a match; if $LB(x, y) \geq t$, then we return (x, y) as a match. Otherwise, we compute $GJ(x, y)$.

The upper and lower bounds are computed as follows. First, for each element $x_i \in B_x$, we find an element $y_j \in Y$ with the highest element-level similarity, such that $s(x_i, y_j) \geq \alpha$ (recall from the description of $GJ(x, y)$ that we consider only matches between $x_i \in B_x$ and $y_j \in B_y$ such that $s(x_i, y_j) \geq \alpha$). Let S_1 be the set of all such pairs.

For example, consider the two strings x and y together with the similarity scores between their elements in Figure 4.18(a) (reproduced from Figure 4.12(a)). Figure 4.18(b) shows the set $S_1 = \{(a, q), (b, q)\}$. Note that for element $c \in B_x$, there is no element in B_y such that the similarity score between them equals or exceeds α , which is 0.5 in this case.

of length
+ 1 = 2).
laps y' , in
arly let y''

(4.5)

= 4. How-
ence, we
generally,
o further

ity mea-
the nor-
necting x

of tokens

ting the
nsive. To
d a lower
; then we
a match.

it $x_i \in B_x$,
 $x_i, y_j \geq \alpha$
 $\in B_x$ and

ty scores
re 4.18(b)
n B_y such
use.

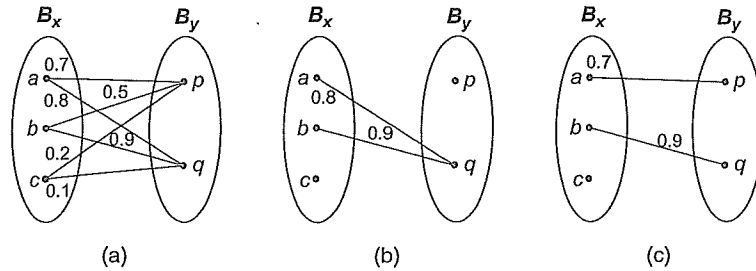


FIGURE 4.18 An example of computing an upper and lower bound for the generalized Jaccard measure.

Similarly, for each element $y_j \in B_y$, we find an element $x_i \in X$ with the highest element-level similarity, such that $s(x_i, y_j) \geq \alpha$. Let S_2 be the set of all such pairs. Continuing with our example, Figure 4.18(c) shows $S_2 = \{(a, p), (b, q)\}$.

The upper bound for $GJ(x, y)$ is given by the following formula:

$$UB(x, y) = \frac{\sum_{(x_i, y_j) \in S_1 \cup S_2} s(x_i, y_j)}{|B_x| + |B_y| - |S_1 \cup S_2|}$$

Note that the numerator of $UB(x, y)$ is at least as large as that of $GJ(x, y)$, and that the denominator of $UB(x, y)$ is no larger than that of $GJ(x, y)$. The lower bound is given by the following formula:

$$LB(x, y) = \frac{\sum_{(x_i, y_j) \in S_1 \cap S_2} s(x_i, y_j)}{|B_x| + |B_y| - |S_1 \cap S_2|}$$

Continuing with our example, $UB(x, y) = \frac{0.8+0.9+0.7+0.9}{3+2-3} = 1.65$ and $LB(x, y) = \frac{0.9}{3+2-1} = 0.225$.

4.3.6 Extending Scaling Techniques to Other Similarity Measures

So far we have discussed scaling techniques for the Jaccard measure or overlap measure. We now describe how to extend these techniques to multiple similarity measures. First, as noted earlier, we can easily prove that

$$J(x, y) \geq t \Leftrightarrow O(x, y) \geq \alpha = \frac{t}{1+t} \cdot (|x| + |y|)$$

by replacing $|x \cup y|$ in $J(x, y)$ with $|x| + |y| - |x \cap y|$. Thus, if a technique works for the overlap measure $O(x, y)$, there is a good chance that we can also extend it to work for the Jaccard measure $J(x, y)$, and vice versa. For example, earlier we described how to extend the prefix filtering technique originally developed for the overlap measure to the Jaccard measure.

In general, a promising way to extend a technique T to work for a similarity measure $s(x, y)$ is to translate $s(x, y)$ into constraints on a similarity measure that already works well

with T . For example, consider edit distance. Let $d(x, y)$ be the edit distance between x and y , and let B_x and B_y be the corresponding q -gram sets of x and y , respectively. Then we can show that

$$d(x, y) \leq \epsilon \Rightarrow O(x, y) \geq \alpha = (\max\{|B_x|, |B_y|\} + q - 1) - q\epsilon$$

Given the above constraint, we can extend prefix filtering to work with edit distance by indexing the prefixes of size $q\epsilon + 1$.

As yet another example, consider the TF/IDF cosine similarity $C(x, y)$. We can show that

$$C(x, y) \geq t \Leftrightarrow O(x, y) \geq \lceil t \cdot \sqrt{|x||y|} \rceil$$

Given this, we can extend prefix filtering to work with $C(x, y)$ by indexing the prefixes of size $|x| - \lceil t^2|x| \rceil + 1$ (this can be further optimized to just indexing the prefixes of size $|x| - \lceil t|x| \rceil + 1$). Finally, the above constraints can also help us extend position filtering to work with edit distance and cosine similarity measures.

Bibliographic Notes

Durbin et al. [196] provide an excellent description of the various edit distance algorithms, together with HMM-based probabilistic interpretations of these algorithms. Further discussion of string similarity measures and string matching can be found in [146, 204, 280, 370, 455]. Tutorials on string matching include [355, 563]. The Web site [118] describes numerous string similarity measures and provides open-source implementations.

Edit distance was introduced in [376]. The basic dynamic programming algorithm for computing edit distance is described in [455]. Variations of edit distance include Needleman-Wunsch [456], affine gap [566], Smith-Waterman [526], Jaro [331], and Jaro-Winkler [575]. Learning the parameters of edit distance and related similarity measures was discussed in [84, 86, 497].

The Jaccard measure was introduced in [329]. The notion of TF/IDF originated from the information retrieval community [414], and TF/IDF-based string similarity measures are discussed in [31, 120, 145, 264, 352]. Soft TF/IDF was introduced in [86, 148]. Generalized Jaccard was introduced in [469]. The Monge-Elkan hybrid similarity measure is introduced in [442]. Cohen et al. [86, 148] empirically compare the effectiveness of string similarity measures over a range of matching tasks.

The Soundex measure was introduced in [500, 501]. Other phonetic similarity measures include New York State Identification and Intelligence System (NYSIIS) [537], Oxford Name Compression Algorithm (ONCA) [250], Metaphone [477], and Double Metaphone [478].

The material on scaling up string matching in Section 4.3 was adapted from [370]. Building inverted indexes to scale up string matching was discussed in [508]. The technique of size filtering was discussed in [34]. Prefix indexes were introduced in [124]. Bayardo et al. [59] discuss how to combine these indexes with inverted indexes to further

een x and
en we can

scale up string matching. On et al. [469] discuss bound filtering, and Xiao et al. [582] discuss position indexes.

Gravano et al. [265] discuss q -gram-based string matching in RDBMSs. Koudas, Marathe, and Srivastava [354] discuss the trade-offs between accuracy and performance. Vernica, Carey, and Li [559] discuss string matching in the map reduce framework. Further techniques to scale up string matching were discussed in [388, 423, 548].

stance by

show that

refixes of
f size $|x| -$
g to work

gorithms,
rther dis-
204, 280,
describes
is.
algorithm
e include
and Jaro-
measures

l from the
asures are
neralized
itroduced
similarity

measures
ord Name
e [478].
om [370].
The tech-
in [124].
to further

