

Nearest Neighbor Search using Additive Binary Tree

Sung-Hyuk Cha and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition

State University of New York at Buffalo, U. S. A.

E-mail: {scha, srihari}@cedar.buffalo.edu

Abstract

Classifying an unknown input is a fundamental problem in Pattern Recognition. One standard method is finding its nearest neighbors in a reference set. It would be very time consuming if computed feature by feature for all templates in the reference set; this naïve method is $O(nd)$ where n is the number of templates in the reference set and d is the number of features or dimension. For this reason, we present a technique for quickly eliminating most templates from consideration as possible neighbors. The remaining candidate templates are then evaluated feature by feature against the query vector. We utilize frequencies of features as a pre-processing to reduce query processing time burden. The most notable advantage of the new method over other existing techniques occurs where the number of features is large and the type of each feature is binary although it works for other type features. We improved our OCR system at least twice (without a threshold) or faster (with higher threshold value).

Keywords – Nearest Neighbor searching, Filtration, Additive Binary Tree, OCR, GSC Classifier

1. Introduction

One common method for classifying an input vector of unknown class involves finding the most or top k similar templates in the reference set. Not surprisingly, this problem, so-called k -nearest neighbor or simply KNN problem has received a great deal of attention because of its significant and practical value in pattern recognition (see [2, 3] for extensive surveys).

One straight-forward method is computing feature by feature for all templates in the reference set and it takes $O(nd)$ where n is the number of templates in the reference set and d is the number of features. This is very time consuming for users to wait for the output. Hence, there is a wealth of literature regarding computational expenses of the KNN problem dating as far back as 1970. Papadimitiou and Bentley showed $O(n^{1/d})$ worst-case algorithm [13] and Friedman et al suggested possible $O(\log n)$ expected time algorithm [6]. There are three general algorithmic techniques for reducing the computational burden: computing partial distance, pre-structuring, and editing the stored prototypes [4].

First, the partial distance technique is often called a *sequential decision technique*; decision for match between two vectors can be made before all features in the vector are examined. It requires a predetermined threshold value to reduce computation time. Next, the pre-structuring method focuses on preprocessing the prototype set into certain well-organized

structures for the fast classification processing. Many approaches utilizing *multidimensional search trees* that partition the space appear in the literature [7, 11, 12, 1]. In these approaches, the range of each feature must be large. Otherwise, if features are binary, we achieve little speedup. Furthermore, the dimension of feature space must be low. Quite often in image pattern recognition, each feature is thresholded for binarization and the dimension is high.

Finally, the prototype editing or reduction method reduces the size of prototype set to improve speed with sacrifice of accuracy. The condensed nearest neighbor rule [9] and the reduced nearest neighbor rule [8] are used to select a subset of training samples to be the prototype set. In this approach, we must sacrifice accuracy for speed. Hong *et al* [10] successfully implemented a fast nearest neighbor classifier for the use of Japanese Character Recognition. They combined a non-iterative method for CNN and RNN and a hierarchical prototype organization method to achieve a great speed-up with a little accuracy drop.

The new algorithm utilizes both partial distance and pre-structuring techniques. We reduce computation time by using an *Additive Binary Tree (ABT)* data structure that contains additive information which is frequency information in binary features. The idea behind the *ABT* approach in finding the nearest neighbor is filtration by which the unnecessary computation can be eliminated. It makes this approach unique from the others such as redundancy reduction or metric method. First, take a quick glance at the reference set and select candidates for match. Next, take a harder look only at those candidates selected from previous filtration to select a fewer candidates, and so on. After several filtration, take a complete thorough look only at the final candidates to verify them. All matches whose distance is less than or equal to a threshold are guaranteed to be in all candidate sets.

In this paper, we will describe the *additive binary tree (ABT)* and the new nearest neighbor search algorithm based on several famous similarity measures. We will give the simulated experimental results. Finally, we report the experimental results on our *OCR* system using the gradient, structural and concavity or simply *GSC* classifier [5].

2. Preliminary

The performance of pattern classification depends significantly on its definition of similarity or dissimilarity measure between pattern vectors. Several definitions have been proposed and encountered in various fields such as information retrieval and biological taxonomy [4]. Among many

definitions, *Euclidean distance*, *absolute differences*, and *Minkowski distance* are the most famous ones. Section 3 discusses the algorithm where the distance measure function is the absolute difference. We chose this definition as it is the most simple measure to explain the new algorithm. *Euclidean* and *Minkowski's distances* are almost identical when features are binary. The absolute difference, well known as the *city block* or *Manhattan metric*, is defined as follows:

$$D[x, y] = \sum_{i=1}^d |x_i - y_i| \quad (1)$$

where d is the dimension of the vector. A normalized inner product appears often as a non-metric similarity function:

$$S[x, y] = \frac{x^t y}{\|x\| \|y\|} \text{ where } \|x\| = \sqrt{x^t x} \quad (2)$$

It is the cosine of the angle between two vectors. Another similarity definition used in *OCR* system using *GSC* classifier developed in *CEDAR* is:

$$S[x, y] = x^t y + \frac{\bar{x}^t \bar{y}}{\gamma} \quad (3)$$

where \bar{x} is the negation of the vector x and γ is the contribution factor ($\gamma > 1$). This definition has the highest accuracy in our *OCR* system among known definitions and it performs best when $\gamma = 1.9$. The difference among these definitions can be explained by its weight where features are binary. The *Minkowski* family definition including *Manhattan* distance gives the same weight for the case where both patterns have the feature and the case where both patterns do not have the feature. The normalized inner product computes the credit for the case where both patterns have the feature but also reflects the other case as it is normalized. The definition 3 gives a full credit for the case where both patterns have features and also allows to control the credit for the other case. Section 4 discusses the third definition of similarity and *OCR*.

We consider the case where a threshold value, t , is known. The *KNN with threshold* problem is that we reject a query vector if the nearest neighbor is not close enough, meaning that the distance exceeds t . Let M be a set of matches whose dissimilarity value is less than or equal to t or whose similarity value is greater than or equal to t . Let $C(x)$ denote the class of the vector x and n is the number of template vectors in the reference set, R .

Definition 1 k-nearest neighbor with threshold

Input: A query vector q , a reference set, $R = \{r_1, r_2, \dots, r_n\}$ and a threshold value, t .

Output: $C(q) = \begin{cases} \text{on the votes of } k\text{-NN} & : |M| \geq k \\ \text{on the votes of } |M|\text{-NN} & : 0 < |M| < k \\ \text{Rejected} & : M = \emptyset \end{cases}$

In this model, the classification decision is made by up to k nearest neighbors because not enough matches may be in the set M . This model is frequently used in many applications considering the fact that the higher reject rate usually lessens the error rate.

We introduce an additive binary tree data structure, *ABT* in short. It is a binary tree whose nodes are sums of their direct children. Consider the example in Fig. 1. The root is $\sum_{i=1}^d x_i$. Nodes at the second level are $\sum_{i=1}^{d/2} x_i$ and $\sum_{i=d/2+1}^d x_i$ from left to right. Leaves are x_1, x_2, \dots, x_d . Let $B(x)$ be a list of nodes, $\{B_1(x), B_2(x), \dots, B_{2d-1}(x)\}$ in breadth first order.

Property 1 Additive Binary Tree

1. There are 2^{l-1} number of nodes at level, l .
2. Total number of nodes = $2d - 1$.
3. The depth of *ABT* = $\log d$.
4. $B_i = B_{2i} + B_{2i+1}$

The structure is named “additive” because additive informations produced by property 4. are appended in every vector. Each element in a vector lies in the leaf level from left to right correspondingly. If vectors are binary, the root is the frequency of 1’s and each node is the frequency of 1’s in sub-part of the vector.

3. k-NN using *ABT* in City block distance

There are three phases: the pre-processing, candidate selection and verification, and voting for classification phases. During the pre-processing phase, the *ABT*’s for all templates are built by adding a pair of elements and storing it in their parent node. Building an *ABT* for one template takes $O(d)$ and thus, total pre-processing takes $O(nd)$ to build *ABT*’s for all templates. The space required is also $O(nd)$ which is at most twice bigger than the size of the reference set.

In the second phase, we search for matches before choosing top k similar templates for votes. The objective of the candidate selection and verification procedure is quickly to find the smallest subset M of the reference set which contains top or up to k templates. Consider the following pseudo code for finding matches.

Algorithm 1 Candidate Selection and Verification

```

1 Build ABT(q)
2 for every template  $x_i$ 
3   for every level  $l = 1$  to  $\log d$ 
4     if  $\sum_{j=2^{l-1}}^{2^l-1} |B_j(x_i) - B_j(q)| > t$ 
5       break
6     else if  $l = \text{leaf}$ 
7       Verify the Match
8       if verified, update  $t$ 

```

The inner loop computes the sum of absolute differences of every node in level, l . The sequential decision technique may be applied to this step; we may not have to compare all nodes in the level because the value may exceed the threshold before all nodes in a specific node are examined. We use the parent level operations as filtration functions. Those candidate templates survived from the parent level are only considered in the next level. After finding all matches, search for top or k -nearest neighbors only in the set, M .

In the line 8, we would like to update the threshold if k number of matches are found and maximum distance is less

than t . This is critical because the lower the t is, the more filtration occurs. Also, if a threshold value is not given at the start, then it is assigned at the point right after k templates are examined and updated if a less value is found later on.

3.1. Example

Consider a following set with $n = 3, d = 8$ and $t = 2$,

$$q = 00000100, \quad R = \begin{cases} r_1 = \underline{01100111} \\ r_2 = \underline{00011100} \\ r_3 = \underline{01110000} \end{cases}$$

When a brute force method is performed, 24 number of comparisons take place. The underlined parts of each reference vectors are those contributing the comparisons when the sequential decision technique is used; there are 19 number of comparisons. Fig. 1 shows *ABTs* for the query vector and all templates. First, the root level is compared. r_2 and r_3 are

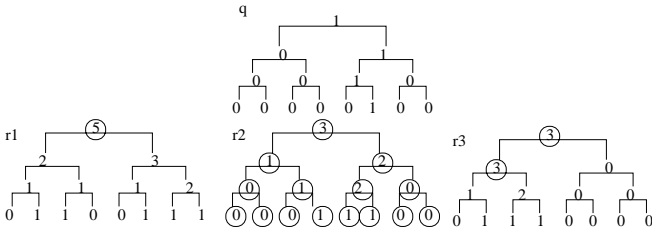


Figure 1. Sample *ABTs*

considered as candidates. r_1 is discarded because $(|B_1(r_1) - B_1(q)| = 4) > (t = 2)$. Next, the second level of trees are compared. For r_2 , $|B_2(r_2) - B_2(q)| + |B_3(r_2) - B_3(q)| = 2$ while that of r_3 is 4. But only the $B_2(r_2)$ is computed because $|B_2(r_2) - B_2(q)|$ already exceeds the threshold. Hence, only r_2 is considered as a candidate. At the level 3, r_2 still is a candidate and at the leaf level, it is verified as the nearest neighbor to q . $|R| = 3, |L_1| = 2, |L_2| = 1, |L_3| = 1, |M| = 1$ The number of comparisons occurred is 18, the number of circled nodes as shown in the example in Fig. 1.

3.2. Correctness

Let f_M and f_l be functions for matching and for selecting candidates at the level, l : $f_M(x, q) = \sum_{i=1}^d |x_i - q_i|$ and $f_l(x, q) = \sum_{i=2^{l-1}}^{2^l-1} |B_i(x) - B_i(q)|$. Let L_l be a set of templates, x 's chosen by the function $f_l(x, q) \leq t$.

Theorem 1 A set of matches is always a subset of all candidate sets. $M = L_{1 \log d} \subset \dots \subset L_l \subset L_{l-1} \subset \dots \subset L_1 \subset R$.

Proof: $f_{L_1} \leq \dots \leq f_{L_{1 \log d}} = f_M$. Proved by the well-known fact $|a + b + c| \leq |a| + |b| + |c|$. For example, $x = (x_1, x_2, x_3, x_4)$ and $q = (q_1, q_2, q_3, q_4)$. Then, $f_M(x, q) = |x_1 - q_1| + |x_2 - q_2| + |x_3 - q_3| + |x_4 - q_4|$ and $f_{l=1}(x, q) = |x_1 + x_2 + x_3 + x_4 - q_1 - q_2 - q_3 - q_4|$. And $f_{l=2}(x, q) = |x_1 + x_2 - q_1 - q_2| + |x_3 + x_4 - q_3 - q_4|$. Clearly, $f_1 \leq f_2 \leq f_M$. Now, if there exists $x \in M$ meaning that $f_M \leq t$, then $f_l \leq t$ and $x \in L_l$. Thus, $M \subset L_{l+1} \subset L_l \subset L_{l-1} \subset L_1 \subset R$. ■

The theorem means that all patterns whose distance is within

the threshold are guaranteed in all candidate sets. Hence, this method is a sort of dynamic prototype reduction method but there is no loss of accuracy at all. It is as accurate as the brute force search method but much faster.

The extra space required for *ABT* is $d - 1$ per vector which is the number of inner nodes by definition. Therefore, the total space is $O(nd)$. Selecting Candidate process is made in the fewer computational time than that of the verification process. Number of comparisons needed at the level, l , is the number of nodes at the level which is 2^{l-1} . That of the top root level, for example, is 1 while that of the leaf level is m . The number of operations is multiplied by 2 in every filtration process. The number of candidates is, on the other hand, reduced in every filtration process.

3.3. Simulated experiment

We present simulation results which demonstrate the performance of the new algorithm on 4 different reference sets with different feature numbers, $d = \{8, 16, 32, 64\}$. The size of the reference set, $|R| = 10,000$ for each reference set. Each feature in vectors is generated by a random function which returns a number between 0 and 99. Possible $\min(D(x, q)) = 0$ which is an exact match and possible $\max(D(x, q)) = 6336$ when $d = 64$. Corresponding test sets, Q 's with $d = \{8, 16, 32, 64\}$, are prepared as well. The size of each query set is $|Q| = 10,000$. The experiment is performed on machine with a 300MHz UltraSparc2 CPU, SunOS 5.6, 588 MIPS, and 256 Memory. Fig. 2 shows the performance of a naïve method and methods using *ABT* structure with several different threshold values.

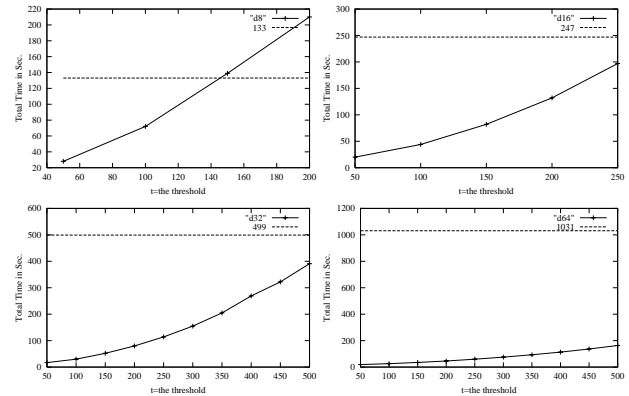


Figure 2. Cumulated elapsed time for 10,000 queries over 10,000 templates

Observation 1 The smaller threshold value, t , the faster *ABT* algorithm runs but the reject rate increases.

A smaller threshold means fewer candidates and the algorithm runs faster. As the threshold increases, on the contrary, candidates abound. Too small threshold, on the other hand, might reject most of inputs.

Observation 2 The larger dimension d is, the better performance the *ABT* algorithm achieves the better execution time performance.

It is necessary to compare the results of d 's with the same ratio of t . Consider the case that $d = 64, t = 300$, then the *ABT* runs 13.8 times faster than the naive method. On the other hands, when $d = 32, t = 150$, it is only 9.6 times faster.

Filtering from top to all the way down to the level, $l - 1$ may not give always the most expeditious running time. It would be wise to stop filtering and jump to the verification stage if the number of candidates is small. Table 1 indicates this behavior.

Table 1. Comparisons for methods.

Method	time in millisec.
Naïve	103.3
Sequential Decision	24.1
Filtration $l = 1$	16.7
Filtration $l = 2$	12.4
Filtration $l = 3$	7.2
Filtration $l = 4$	9.5

Observation 3 *There exists a level, \hat{l} which results the minimum elapsed time.*

The deeper the level is, the more computations which is the number of nodes at the level, are required. \hat{l} is the starting point that $|L_i| \times 2^i \leq |L_{i+1}| \times 2^{i+1}$ which means that we gain no advantage of filtration from this level. Therefore, finding \hat{l} is important because it gives not only the minimum elapsed time, but also less required space for *ABT* data structure.

Observation 4 *The ABT technique is superior to the sequential decision technique experimentally.*

In the worst case when every reference is a match, both sequential decision and *ABT* would not give any speed up. Technique using *ABT* would be even twice slower than the naive method in case of a full *ABT*.

3.4. Auxiliary

Lookup Table: Suppose $d = 512$ and all features are binary, then a vector is stored as an array of 16 unsigned integer values. We build a lookup-table for counts of 1's for all 16 bit combination. The size of lookup-table is 262K bytes. This look-up table for unsigned integers is built at preprocessing stage. To find $D[x, y]$, we perform table look-ups for high and low 16 bits of each unsigned integer ($= x \otimes y$). There are total 32 table look-ups to determine the distance. Where vectors are stored as unsigned integers, we consider counts of 1's of each unsigned integer value as leaves of *ABT*. Lookup table enables significant speedup over the method counting the number of bits without the table.

Ordered List: The simulated experiment tells that the smaller threshold is, the faster algorithm runs (see Observation 1). The threshold value is dynamically changed as more templates are examined. It would be desirable if the smaller threshold value is assigned early. Therefore, we would like to order the templates so that the less threshold value is set early in the search.

Fact 1 *The closer in frequency two vectors are, the smaller in absolute difference distance they tend to be.*

If we order templates by its frequency and search from a template whose difference of frequency is smallest, then the low threshold value might be set in the earlier search stage resulting in a great speed up. First, order the reference set by B_1 . For each bin, the average number of templates is n/m . Next, order each bin by B_2 . During the search stage, start searching from the template, x whose $\min(|B_1(x) - B_1(q)|)$ and then $\min(|B_2(x) - B_2(q)|)$. As a result of ordering, not only we find the low threshold value quickly, but also, our search space is reduced down to the size of L_2 .

Selection: We have seen that *ABT* facilitates filtering some reference templates from consideration by using the lower bounds at each level. In addition to the lower bounds, *ABT* also provides upper bounds. If an upper bound is less than or equal to the threshold, the match is verified without further calculation. This selection technique is useful in a query such that we would like to find all matches to the query vector within a threshold. Consider w -ary vectors x and y where each element can have a value between 0 and $w - 1$.

Lemma 1 $D[x, y]$ has upper and lower bounds at the root level: $|B_1(x) - B_1(y)| \leq D[x, y] \leq \min(B_1(x) + B_1(y), (wd - B_1(x)) + (wd - B_1(y)))$.

Proof: The lower bound is given in Theorem 1. Similarly, $\sum_{i=1}^d |x_i - y_i| \leq \sum_{i=1}^d x_i + \sum_{i=1}^d y_i$ because $|a + b| \leq |a| + |b|$. Now, $D[x, y] = \sum_{i=1}^d |x_i - y_i| = \sum_{i=1}^d |(w - x_i) - (w - y_i)|$. This is the absolute difference of two inverse vectors and it is the same as the absolute difference of the original vectors. Clearly, $\sum_{i=1}^d |(w - x_i) - (w - y_i)| \leq \sum_{i=1}^d (w - x_i) + \sum_{i=1}^d (w - y_i)$. Thus, $D[x, y] \leq \min(\sum_{i=1}^d x_i + \sum_{i=1}^d y_i, \sum_{i=1}^d (w - x_i) + \sum_{i=1}^d (w - y_i))$

Consider two vectors, x, y where they are binary, $w = 1$. $\begin{cases} x = 00011111 & : B_1(x) = 5 \\ y = 01001011 & : B_1(y) = 4 \end{cases}$ According to the lemma 1, we achieve $1 \leq D[x, y] \leq 8$ immediately.

Theorem 2 $\sum_{i=2^{l-1}}^{2^{2l}-1} (B_i(x) - B_i(y)) \leq D[x, y] \leq \min(\sum_{i=2^{l-1}}^{2^{2l}-1} (B_i(x) + B_i(y)), \sum_{i=2^{l-1}}^{2^{2l}-1} ((\frac{wd}{t} - B_i(x)) + (\frac{wd}{t} - B_i(y))))$

Proof: Again, the lower bounds are proved in Theorem 1. For the upper bounds, divide vectors into 2^{l-1} number of sub-vectors. For each sub-vectors, Lemma 1 holds.

For the above example, we have the upper and lower bounds at the second level of *ABT*: $1 \leq D[x, y] \leq 3$.

4. Using *ABT* for GSC classifier

Among many classifiers in *OCR*, the *Gradient, Structural, and Concavity classifier*, simply known as *GSC classifier*, has 98% accuracy on 24,000 handwritten digits taken from various databases [5]. It is based on the philosophy that feature sets can be designed to extract certain types of information from the image. These types are gradient, structural, and concavity informations. Gradient features use the stroke shapes on a small scale. Next, structural features are based on the

stroke trajectories on the intermediate scale. Finally, concavity features use stroke relationship at long distances. Gradient, Structural, and Concavity has 192, 192, and 128 feature vectors correspondingly. In all, there are 512 number of features in a vector.

The K -NN approach is often used in the *GSC classification*. It selects the top k similar templates in the training set and returns its class based on the vote of these k template vectors. The definition of this similarity, $S[x, y]$, currently used in the GSC classifier follows:

$$S[x, y] = \sum_{i=0}^n S[x_i, y_i] \text{ where } S[x_i, y_i] = \begin{cases} 1 & : x_i = y_i = 1 \\ 1/\gamma & : x_i = y_i = 0 \\ 0 & : \text{otherwise} \end{cases}$$

When both x_i and y_i are 1's, it is the case that the desired feature exists. It is denoted as $S_{11}[x, y] = x^t y$. When both x_i and y_i are 0's, it is the case that the feature does not exist. It is denoted as $S_{00}[x, y] = \bar{x}^t \bar{y}$. It would be reasonable to rely more upon S_{11} than S_{00} . The definition of the similarity measure becomes:

$$S[x, y] = S_{11}[x, y] + \frac{S_{00}[x, y]}{\gamma}$$

γ is the contribution factor and usually $\gamma \geq 1$. The recent experiment shows that $\gamma = 1.9$ gives the best performance. Note that when $\gamma = 1$, it is the inverse of the city block distance between two vectors.

4.1. k -NN for GSC classifier

There is a slight difference in line 4 of the previous algorithm 1 as its definition of similarity differs. Let $B_i(x)$ and $W_i(x)$ be the number of 1's and 0's at the i 's node of vector x respectively. The line 4 is replaced by:

$$\text{if } \sum_{i=2^{l-1}}^{2^{2l}-1} (\min(B_i(x), B_i(q)) + \frac{\min(W_i(x), W_i(q))}{\gamma}) \geq t.$$

Fact 2 $B_i(x) = \frac{d}{2^{l-1}} - W_i(x)$ where l is the level of i 's node.

Lemma 2 If $\min(W_i(x), W_i(y)) = W_i(x)$, then $\min(B_i(x), B_i(y)) = B_i(y)$ and vice versa.

Proof: $B_i(x) = \frac{d}{2^{l-1}} - W_i(x)$ and $B_i(y) = \frac{d}{2^{l-1}} - W_i(y)$. Now if $B_i(x) \leq B_i(y)$, then $\frac{d}{2^{l-1}} - W_i(x) \leq \frac{d}{2^{l-1}} - W_i(y)$. By rearranging the formula, we get $W_i(x) \geq W_i(y)$. ■

Lemma 3 $S[x, y] \leq \min(B_1(x), B_1(y)) + \frac{\min(W_1(x), W_1(y))}{\gamma}$.

Proof: When only frequency values are available, the maximum value $S_{11}(x, y)$ is when 1's in $B_1(x)$ and $B_1(y)$ are aligned together. $S_{11}(x, y)$ cannot exceed $\min(B_1(x), B_1(y))$. The same for $S_{00}(x, y)$. ■ The maximum similarity value that two vectors x and y can have at the root level of *ABT* is $\min(B_1(x), B_1(y)) + \frac{\min(W_1(x), W_1(y))}{\gamma}$. Suppose $\gamma = 2$. In the examples in Fig. 1, $f_1(q, r_1) = \min(1, 5) + \frac{\min(73)}{2}$ which is 2.5. Similarly, $f_1(q, r_2) = 3.5$, $f_1(q, r_3) = 3.5$. These values are always larger than or equal to $S[x, y]$'s defined in GSC classifier. When there is an exact match meaning that $x = y$, $S[x, y] = B_1(x) + \frac{W_1(x)}{\gamma}$.

Therefore, the maximum similarity value varies from vector to vector by their frequencies.

Let f_l be a function for selecting candidates at the level, l : $f_l(x, q) = \sum_{i=2^{l-1}}^{2^{2l}-1} (\min(B_i(x), B_i(q)) + \frac{\min(W_i(x), W_i(q))}{\gamma})$. Let L_l be a set of templates, x 's chosen by the function $f_l(x, q) \geq t$.

Corollary 1 All members in M are guaranteed to be in all candidate sets. $M = L_{\log d} \subset \dots \subset L_l \subset L_{l-1} \dots \subset L_1 \subset R$

Proof: Let $B(x)$ and $W(x)$ be the frequency of 1's and 0's of one internal node. Let $B_l(x), W_l(x)$ and $B_r(x), W_r(x)$ be those of the left and right children of the node. By the definition, $B(x) = B_l(x) + B_r(x)$ and $W(x) = W_l(x) + W_r(x)$. Consider two vectors x and y . Suppose $B(x) = \min(B(x), B(y))$, then the maximum possible value at the internal node can have becomes $B(x) + \frac{W(y)}{\gamma} = B_l(x) + \frac{W_l(y)}{\gamma} + B_r(x) + \frac{W_r(y)}{\gamma}$. Now at its children's nodes, there are three possible cases and one impossible case. First, if $B_l(x) = \min(B_l(x), B_l(y))$ and $B_r(x) = \min(B_r(x), B_r(y))$, then the maximum value at the both children nodes can have is the same as that of their parent's node by definition: $B_l(x) + B_r(x) \leq B_l(y) + B_r(y)$ and $W_l(x) + W_r(x) \geq W_l(y) + W_r(y)$. Next, if $B_l(x) = \min(B_l(x), B_l(y))$ but $B_r(y) = \min(B_r(x), B_r(y))$, then it is smaller than or equal to that of their parent's node.

$B_l(x) + \frac{W_l(y)}{\gamma} + B_r(y) + \frac{W_r(x)}{\gamma} \leq B_l(x) + \frac{W_l(y)}{\gamma} + B_r(x) + \frac{W_r(y)}{\gamma}$. When $B_l(y) = \min(B_l(x), B_l(y))$ but $B_r(x) = \min(B_r(x), B_r(y))$, it is also smaller than or equal to that of their parent's node. It is impossible by definition if $B_l(y) = \min(B_l(x), B_l(y))$ but $B_r(y) = \min(B_r(x), B_r(y))$. Therefore, the upper bound of a parent level is always greater than or equal to that of its children's nodes. This is true for all internal nodes. The upper bound at a certain level is the sum of all upper bounds at all nodes in that level. Let f_l be the upper bound function at a level, l , then $f_{L_1} \geq \dots \geq f_{L_{\log d}} = f_M$. Thus, $M = L_{\log d} \subset \dots \subset L_1 \subset R$. ■

This corollary is the sine qua non which guarantees the correctness and speedup of Algorithm for GSC classifier.

The following lemma further accelerates the search process, which we state it without a proof.

Lemma 4 $S_{00}(x, y) = d - B_1(x) - B_1(y) + S_{11}(x, y)$.

This means that once we compute the similarity for $S_{11}(x, y)$, $S_{00}(x, y)$ is computed in constant time.

4.2. Experiment

Before embarking on the results, it is important to discuss the relationship between error and reject rates of a recognizer. As shown in Fig. 3, the error versus reject percentages of recognition graph is a great way to evaluate the performance of OCR systems. A good recognizer must be as close to the axes as possible. Typically, the higher reject rate, the lower error rate. The threshold value must be made depending on the costs of rejects and errors. As the threshold value increases, the reject rate also increases and the error rate decreases.

We use two sets of isolated mixed hand-printed/cursive character images. The first set is the reference set whose

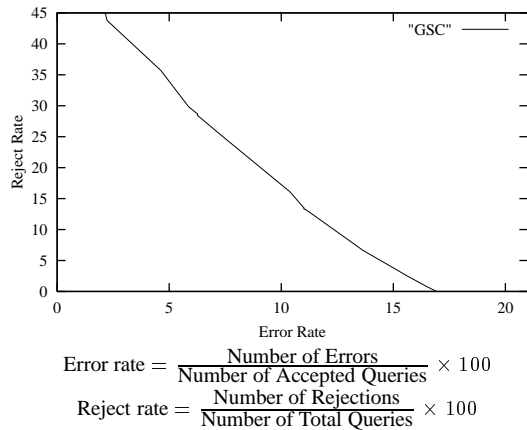


Figure 3. Error vs. Reject graph for GSC classifier

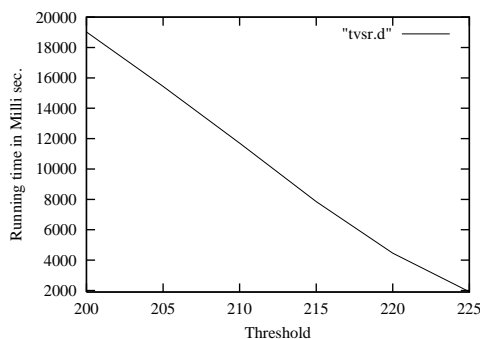


Figure 4. Threshold vs. running time

size is 21 800 and approximately 800 per character (A-Z). The other set is a test set whose size is 1681. While the average running time of the brute force method is 39.470 millisecond, one using ABT is 19.714 millisecond where the depth of tree is 2 and $t = 0$.

The higher threshold value the classifier has, the faster it runs. as shown in Fig. 4 The higher threshold value means the more rejects which save running time. Hence, the importance of a threshold is twofold; it not only enables to control the error and reject rates, but also gives speed up.

5. Finale

Nearest Neighbor searching is an extremely well studied area and a couple of techniques to speed up the OCR system were previously developed. They are prototype thinning and clustering techniques which cause a little degradation in performance. Because even a slight degradation in performance is often too costly in real OCR applications, we studied it and a new fast nearest neighbor search algorithm with no degradation is proposed. Filtration with a threshold is the key idea for expedition. To perform this, we introduced ABT structure. The additive information of pattern feature vectors ensures that the query processing time reduces significantly. Furthermore, the idea is effective even in combination with other techniques like prototype thinning or clustering.

There are two parameters that can affect the speed of search : i) depth of ABT and ii) the number of branches.

As stated in Observation 3, the depth of ABT influences the speed significantly. We introduced the additive binary tree, however, the tree can be an arbitrary branch tree. Additive N-ary Tree might perform better than the binary tree.

Moreover, the idea of filtration using ABT can extend to other definitions with a little embellishment. Consider the normalized inner product defined in equation 2. The maximum value for this definition is 1 when there is an exact match. We exclude the exceptional case where all features are 0 as it makes the denominator be 0. A filtration occurs in case that $\sum_{i=2^{2t}-1} \frac{\min(B_i(x), B_i(q))}{\|x\| \|q\|} \geq t$ Again, consider the example in Fig. 1 with a threshold = 0.5. At the root level filtration, r_1 is filtrated because the upper bound is 0.447.

References

- [1] A. J. Broder. Strategies for efficient incremental nearest neighbor search. *Pattern Recognition*, 23(12):171–178, 1990.
- [2] B. V. Dasarathy. Visiting nearest neighbors- a survey of nearest neighbor pattern classification techniques. In *Proceedings of the International Conference on Cybernetics and Society*, p 630–636. IEEE, 1977.
- [3] B. V. Dasarathy. *Nearest Neighbor Pattern Classification Techniques*. IEEE Computer Society Press, 1991.
- [4] R. O. Duda, D. G. Stork, and P. E. Hart. *Pattern Classification and Scene Analysis*. John Wiley & Sons, Inc., 2nd edition, 2000.
- [5] J. T. Favata, G. Srikantan, and S. N. Srihari. Hand-printed character/digit recognition using a multiple feature/resolution philosophy. In *IWFHR-IV*, p 57–66, 1994.
- [6] J. Friedman, J. Bentley, and R. Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Tran. on Mathematical Software*, 3:209–226, 1977.
- [7] K. Fukunaga and P. Narendra. A branch and bound algorithm for computing k-nearest neighbors. *IEEE Tran. on Computers*, 24:750–743, 1975.
- [8] G.W. Gates. The reduced nearest neighbor rule. *IEEE Tran. on Information Theory*, IT-18(3):431–433, 1972.
- [9] P.E. Hart. The condensed nearest neighbor rule. *IEEE Tran. on Information Theory*, IT-14(3):515–516, 1972.
- [10] T. Hong, S. W. Lam, J. J. Hull, and S. N. Srihari. The design of a nearest-neighbor classifier and its use for japanese character recognition. In *Proceedings of IC-DAR*, p 370–377. IEEE, 1995.
- [11] B. S. Kim and S. B. Park. A fast k-nearest neighbor finding algorithm based on the ordered partition. *IEEE Tran. on PAMI*, 8(6):761–766, 1986.
- [12] H. Niemann and R. Goppert. An efficient branch-and-bound nearest neighbour classifier. *Pattern Recognition Letters*, 7:67–72, 1988.
- [13] C. Papadimitriou and J. Bentley. A worst-case analysis of nearest neighbor searching by projection. In *In Automata, Languages and Programming LNCS*, v. 85, p 470–482. Springer-Verlag, 1980.