

A Fast Algorithm for Finding k -Nearest Neighbors with Non-metric Dissimilarity

Bin Zhang and Sargur N. Srihari
CEDAR, Department of Computer Science and Engineering
State University of New York at Buffalo
Buffalo, NY 14228.
Email: {binzhang,srihari}@cedar.buffalo.edu.

Abstract

Fast nearest neighbor (NN) finding has been extensively studied. While some fast NN algorithms using metrics rely on the essential properties of metric spaces, the others using non-metric measures fail for large-size templates. However, in some applications with very large size templates, the best performance is achieved by NN methods based on the dissimilarity measures resulting in a special space where computations cannot be pruned by the algorithms based-on the triangular inequality. For such NN methods, the existing fast algorithms except condensing algorithms are not applicable. In this paper, a fast hierarchical search algorithm is proposed to find k - NN s using a non-metric measure in a binary feature space. Experiments with handwritten digit recognition show that the new algorithm reduces on average dissimilarity computations by more than 90% while losing the accuracy by less than 0.1%, with a 10% increase in memory.

1 Introduction

K -nearest neighbor (k - NN) search has been extensively used as a powerful non-parametric technique for pattern recognition[1]. However, k - NN search requires intensive dissimilarity computations, particularly for large training set, search of the whole set is unacceptable. Therefore, speeding-up k - NN search is a key to make k - NN classification useful in practice.

Any algorithm for speeding-up k - NN search falls into one of the two categories. The first category is to condense template set[2]. The idea is to remove redundant patterns in a template set, thus reduce the size of the template set and search complexity. The second is to reorganize the patterns of a training set during the training process so that the k - NN search in the classification process can be done efficiently [3] [4] [5].

The algorithms[3] [5], which belong to the second category, rely on the essential properties of metric spaces. However, in some applications, the best performance is achieved by NN methods based on the dissimilarity measures resulting in a special sub-space where addition of any two dissimilarity measures is always bigger than any single measure, so these computation pruning methods based on triangular inequality completely fails for such applications. Although the algorithms in [4] are applicable for non-metric dissimilarity measures, they are effective only with small-size (≤ 10) feature vectors.

In this paper, we describe properties of a non-metric dissimilarity measure $D(\cdot, \cdot)$. $D(\cdot, \cdot)$ that has been successfully used in character recognition at CEDAR [6]. The k - NN search is performed on 512-dimensional binary feature vectors. For such NN methods with $D(\cdot, \cdot)$, the existing fast algorithms of the second category are not applicable.

The dissimilarity measure $D(\cdot, \cdot)$ is actually a variation of Sokal-and-Michner similarity, one of the eight similarity measures for binary vector matching discussed in [7]. As a comprehensive exploration of the properties of each of the eight similarity measures in [7] is beyond the scope of this paper, here we focus on $D(\cdot, \cdot)$.

A fast binary vector matching algorithm [8] was proposed recently. The idea is to reorganize a binary feature vector into an Additive Binary Tree (ABT) and transform the dissimilarity computation of two binary vectors into comparison of the two corresponding ABTs. Given a threshold, the dissimilarity computation is saved by stopping the comparison at some parent level.

While ABT saves computation of the dissimilarity between two binary vectors (local speeding-up), we approach this problem in a different way, that is, all templates in a training set is organized as a hierarchical tree based on dissimilarity measures and class labels for efficient k - NN search using $D(\cdot, \cdot)$ (global speeding-up). Although the algorithm aims at pruning k - NN computations with $D(\cdot, \cdot)$, it is suitable for NN methods with any non-metric dissimilarity measures.

The rest of the paper is organized as follows. In Section 2, we define $D(\cdot, \cdot)$ and explore its non-metric properties. In Section 3, we discuss k-nearest neighbor classification using the measure. In Section 4, we present the hierarchical search algorithm. In Section 5, the experimental results are given and analyzed. We draw conclusions in Section 6.

2 Non-metric Dissimilarity

2.1 Metric Space

A *metric* defined in a set S is a *distance function* d which assigns to each pair of points of S a distance between them, and which satisfies the following four properties: non-negativity, commutativity, triangular inequality, and reflexivity.

2.2 Dissimilarity for Binary Feature Vectors

A binary vector Z with N dimensions is defined as:

$$Z = (z_1, z_2, \dots, z_N). \quad (1)$$

where

$$z_i \in \{0, 1\}, \forall i \in \{1, 2, \dots, N\}.$$

We define Ω as the set of all N -dimensional binary vectors.

Then, the **unit binary vector** I is a binary vector with every element equal to 1.

The **complement** of a binary vector Z is

$$\bar{Z} = I - Z. \quad (2)$$

The L_1 - *norm* of a binary vector Z is

$$|Z| = \sum_{i=1}^N z_i, \quad (3)$$

which is the number of 1's in Z .

Given two binary vectors, X and $Y \in \Omega$, the Sokal-and-Michner (S-M) similarity[7] is given by

$$S_{sm}(X, Y) = (X \odot Y + \bar{X} \odot \bar{Y})/N. \quad (4)$$

where \odot is the inner product operator defined as $X \odot Y = \sum_{i=1}^N x_i y_i$.

A generalization of S-M similarity measure is

$$S(X, Y) = X \odot Y + \beta \bar{X} \odot \bar{Y}, 0 \leq \beta \leq 1. \quad (5)$$

The complement of $S(X, Y)$ is the **dissimilarity** given by

$$\begin{aligned} D(X, Y) &= N - X \odot Y - \beta \bar{X} \odot \bar{Y} \\ &= N - D1(X, Y) - \beta D0(X, Y), 0 \leq \beta \leq 1, \end{aligned} \quad (6)$$

where $D1(X, Y) = X \odot Y$ is the number of positions on which X and Y both have a 1, and $D0(X, Y) = \bar{X} \odot \bar{Y}$ is the number of positions on which X and Y both have a 0.

When $\beta = 1$, the bits that are 1 or 0 have the same weight, and the distance $D(X, Y)$ is the same as Hamming distance. When $\beta = 0.5$, distance between two 1s is 0, but two zeros have a distance of 0.5. Intuitively, when $\beta \neq 1$, two bits equal to 0 are more far away from each other than two 1s.

With the formula (2), we can rewrite (6) as:

$$D(X, Y) = (1 - \beta)N + \beta(|X| + |Y|) - (1 + \beta)X \odot Y. \quad (7)$$

2.3 Non-metric Property of $D(X, Y)$

When $\beta = 1$, the dissimilarity $D(X, Y)$ is Manhattan distance which is a metric. But if $\beta \neq 1$, $D(X, Y)$ is not a metric because it only satisfies three properties of a metric, i.e., non-negativity, commutativity, and triangular inequality, and violates reflexivity in general. These properties are stated here with proof omitted. We conclude that $D(X, Y)$ is not a metric unless $\beta = 1$.

3 k-Nearest Neighbor Classification Using $D(X, Y)$

The k-nearest neighbor rule of classification assigns a class label for unknown classification by examining its k-nearest neighbors of known classification [1]. It has been used in binary classification using $D(X, Y)$ as the distance measure at CEDAR [6]. For example, the GSC classifier extracts Gradient, Structural, and Concavity (GSC) features from a binary image, and quantizes these features into a binary vector of 512 bits, then conducts 6-nearest neighbor search on a training set for majority voting. A 6-NN character classifier GSC recognizer achieves high accuracy with $\beta = 0.5$.

3.1 Performance in Handwritten Digit Recognition

Theoretically, $D(X, Y) \in [0, N]$, however, in GSC application, the domain of $D(X, Y)$ spans a very narrow range. Traditionally, $N = 512$ and $\beta = 0.5$. Given a testing set \mathcal{F} of 13,565 handwritten digits and a training set \mathcal{G} of 14,678 numeral patterns, there are $13565 * 14678 = 1.991070e + 08$ dissimilarities, resulting in a set ϕ^a . The distribution of dissimilarities in ϕ^a is plotted in Figure 1. From this experiment we also find:

$$207 \leq D(X, Y) \leq 373, \forall X \in \mathcal{F}, \forall Y \in \mathcal{G}. \quad (8)$$

Whatever, $\beta = 0.5$ is selected intuitively, it is not an optimal value for GSC application. Figure 2(a) shows the relationship between recognition rate and the coefficient β .

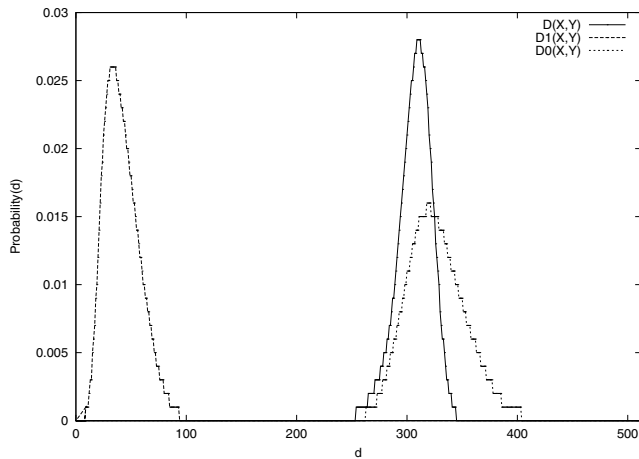


Figure 1. Distributions of $D(X, Y)$ ($\in \phi^a$), $D1 = X \odot Y$, and $D0 = \bar{X} \odot \bar{Y}$ with $\beta = 0.5$, $\forall X \in \mathcal{F}, \forall Y \in \mathcal{G}$. Here, d represents dissimilarity value.

From Figure 2(a), we may find that the recognition rate gets to the maxima 99.425% when $\beta = 0.42$, whereas, the rate drops down to 98.85% when $\beta = 1.0$. With $\beta = 0.5$, the rate is 99.38%.

The experiments above show that in practice, β should not be set to 1.0 for better performance, thus, the dissimilarity measure in GSC is not a metric, as discussed in Section 2.

By the result (8), we can conclude that:

$$\forall D_i(X, Y), D_j(X, Y) \text{ and } D_k(X, Y) \in \phi^a,$$

$$D_i(X, Y) \leq D_j(X, Y) + D_k(X, Y). \quad (9)$$

A more general result can be seen in Figure 2(b). It's easy to find that the inequality (9) holds if $\beta \leq 0.58$.

3.2 Traditional Fast k-NN Search Method

The branch and bound algorithm [5] is a representative of the traditional fast nearest neighbor search methods. This algorithm applies two rules to reduce the computations, here we just reiterate the first rule and show that this rule is not applicable to GSC dissimilarity measure.

Let n be a node of the search tree representing a set of patterns T_n , H_n be the hyper center of the set T_n , R_n be the radius of the set T_n , and Y be the current nearest pattern to Z with a distance L .

Then, the first rule shown as Figure 3 claims that no $X_i \in T_n$ can be the nearest neighbor to Z , if $L + R_n < d(Z, H_n)$.

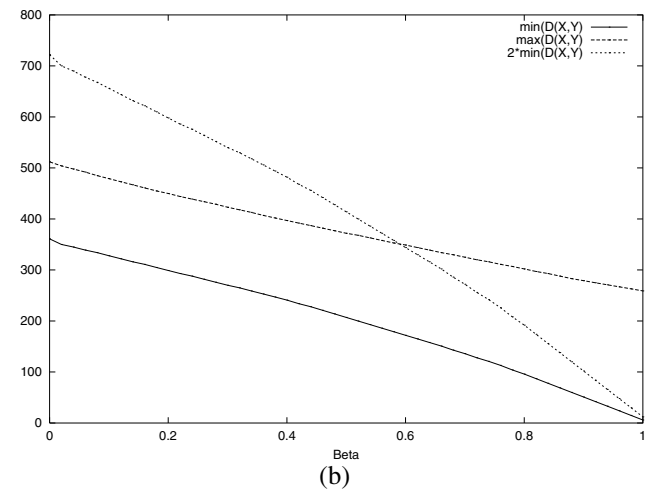
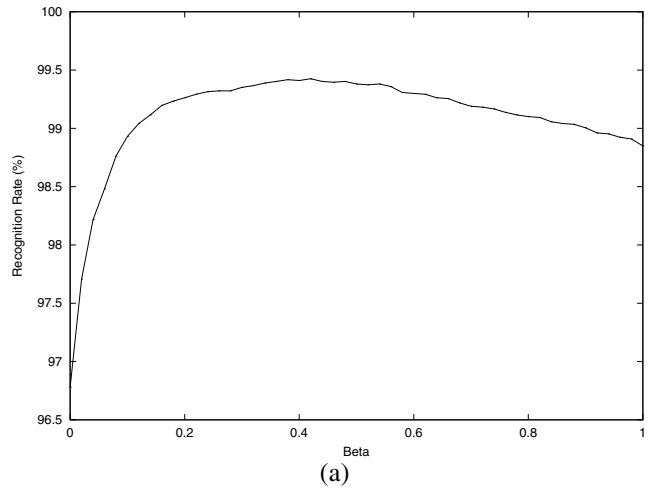


Figure 2. Accuracy and dissimilarity measures vs β : (a) recognition rate vs coefficient (β). (b) $\min(D(X, Y, \beta)), \max(D(X, Y, \beta))$, and $2 * \min(D(X, Y, \beta))$

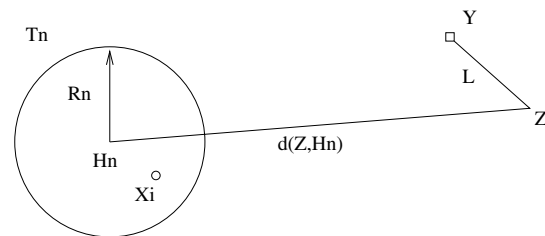


Figure 3. Illustration of the reduction rule

As for GSC setup, according to the formula (9), the condition, $L + R_n < d(Z, H_n)$, can never be true, therefore, no computations can be pruned by this rule. It's the same case with the 2nd rule. Finally, we can conclude that *the triangular-inequality based rules in the traditional fast k-NN algorithms for pruning distance computations will not be applicable for k-NN search with $D(X, Y)$.*

In the following section, we describe a fast algorithm to search k -NNs using the dissimilarity measure $D(X, Y)$ without considering geometric relations among a set of feature vectors.

4 The Hierarchical Search Algorithm

Let $\mathcal{T} = \{Z_1, Z_2, \dots, Z_m | Z_i \in \Omega, 1 \leq i \leq m\}$ be a set of training patterns, and each pattern Z_i has a class label $C(Z_i)$. The target is to efficiently find the k -nearest neighbors of a testing pattern X ($X \in \Omega$) in \mathcal{T} with the dissimilarity measure $D(X, Y)$ defined as the formula (6).

Because of the special property of $D(X, Y)$, we cannot take advantage of geometric relations of feature vectors by applying reduction rules based on triangular inequality. The algorithm proposed in the following uses only dissimilarities between templates and template class labels (in training set). This makes the algorithm actually effective for any dissimilarity measures, metric or non-metric. Of course, for metric measures, we have better algorithms to prune k -NN computations [3] [5].

The algorithm includes two stages, training stage and classification stage. In the training stage, a training set is restructured into a multi-level tree, and in classification stage hierarchical searching is conducted by choosing for comparison a subset of templates at each level of the tree based on the information from the searching result of the exact upper level till any *decision level*. Since only a small sub-set of templates is used for comparison with a given template with unknown class, a great number of dissimilarity computations is able to be avoided, thus speeding k -NN search.

4.1 Building-up of The Search Tree

Given a full template set \mathcal{T} , the bottom level of the tree \mathcal{B} consists of all patterns in \mathcal{T} , i.e., a decision level. In this algorithm, we'll establish another decision level, **hyper level**, where each pattern (**hyper pattern**) is a cluster center of some patterns in \mathcal{T} with the same class label, hence we can make decision in this level without searching the next (bottom) level.

For a pattern $Z \in \mathcal{T}$, we explore its local properties of, $\gamma(Z)$ as dissimilarity between Z and its nearest neighbor with a different class label, and $\Psi(Z)$ as the set of nearest neighbors with class label $C(Z)$ and dissimilarity to Z less

than $\gamma(Z)$, and $\ell(Z)$ as the size of the set $\Psi(Z)$. Mathematically, they are defined as:

$$\gamma(Z) = \{D(Z, Y) | C(Y) \neq C(Z) \text{ and } C(V) = C(Z), \forall D(Z, V) < D(Z, Y), Y, V \in \mathcal{T}\} \quad (10)$$

$$\Psi(Z) = \{V | D(Z, V) < \gamma(Z), \forall V \in \mathcal{T}\} \quad (11)$$

$$\ell(Z) = \{\sum \aleph(Y) | \forall Y \in \Psi(Z)\} \quad (12)$$

where, $\aleph(Z) = 1$.

Building-up of the search tree is as follows:

Step 1: Compute the localities of each pattern $Z \in \mathcal{T}$, $\gamma(Z)$, $\Psi(Z)$, and $\ell(Z)$. And sort Z in \mathcal{T} by $\ell(Z)$ in descendant order. The bottom level of the tree, \mathcal{B} , is set to be empty.

Step 2: Take the top pattern Z_1 in \mathcal{T} as a hyper pattern, put $\Psi(Z_1)$ into \mathcal{B} , and remove all patterns in $\Psi(Z_1)$ from \mathcal{T} . A link between Z_1 and each pattern of $\Psi(Z_1)$ in \mathcal{B} is established. For the patterns left in \mathcal{T} , we do the same thing on the top pattern.

Step 3: Repeat *Step 2* till \mathcal{T} becomes empty. In this way we build up the hyper level \mathcal{H} and the bottom level \mathcal{B} in the search tree.

Step 4: Select a threshold η , we group the patterns in \mathcal{H} so that the radius of each group is less or equal to η , all group centers form another higher level of the search tree, \mathcal{P} .

Step 5: Increase the threshold η and repeat *Step 4* on \mathcal{P} till only one pattern is left in the resulted higher level.

It's easy to see that each intermediate node in the tree may have different number of sub-nodes. Moreover, the hyper level \mathcal{H} and the bottom level \mathcal{B} are 'meaningful' so that we can make decision at these two levels, but all levels above the hyper level in the tree are 'meaningless', i.e., the class of a given pattern cannot be decided at those levels, instead, they act as routes to direct efficient search down to the lower levels.

A search tree contains at least a hyper level and a bottom level. The total number of levels in a search tree is decided by the parameter η , which is an ascendant function of recursive depth in *Step 4*.

4.2 Finding k -NN in the Search Tree

Given a testing pattern X , its k -NN in the full template set can be found by traversing the search tree. k -NN search begins from the 2nd level of the search tree, and selects the top ζ nodes ($\zeta \geq K$) closest to X at each level measured by the dissimilarity $D(X, Y)$ as (6), then compares X with only the descendant nodes of the ζ nodes selected to find another set of ζ nodes closest to X for the further searching. The searching process above is repeated till any decision level in the tree.

Here, the key is how to make decision at *the hyper level*. The rule is that the class of a given pattern X is decided only if all clusters into which X falls have the same class label. Because of this strict condition, the recognition rate in *the hyper level* is quite high, and the algorithm stops at *the hyper level* for most patterns. Once X fails to be classified at *the hyper level*, we search further by examining the descendants (in *the bottom level*) of the ζ closest patterns to X in *the hyper level*, where the final decision can be made.

The bigger ζ , the more dissimilarity computations, and the higher accuracy. In the next section, we'll show the affect of ζ on accuracy and computation complexity.

5 Performance Variations with ζ

A training set \mathcal{G} of 14678 feature vectors extracted from handwritten digits is used for building up the search tree. The hierarchical search tree has 5 levels. From the root to the bottom of the search tree, the size in number of nodes at each level is 1, 4, 76, 1360, and 14678. The tree contains 16119 nodes in total, an increase of 10% compared with the original full template set.

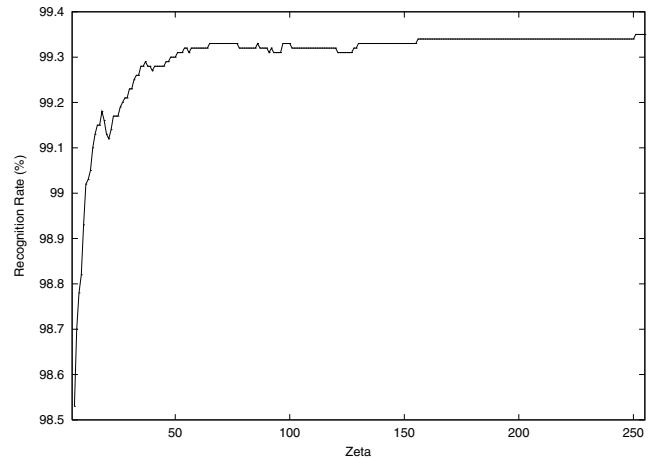
We conduct the hierarchical search on \mathcal{G} with a testing set \mathcal{F} of 13565 handwritten digits. Figure 4 shows the dynamic performance of the hierarchical search algorithm as a function of the parameter ζ . The bigger ζ , the more dissimilarity computations, but higher recognition rate. The comparison of performance of direct search and hierarchical search with different ζ is presented in Table 1.

	Direct	Hierarchical		
		$\zeta = 6$	$\zeta = 50$	$\zeta = 251$
Avg. No. of Nodes visited	14678	386	1441	2244
Recognition Rate (%)	99.38	98.53	99.30	99.35

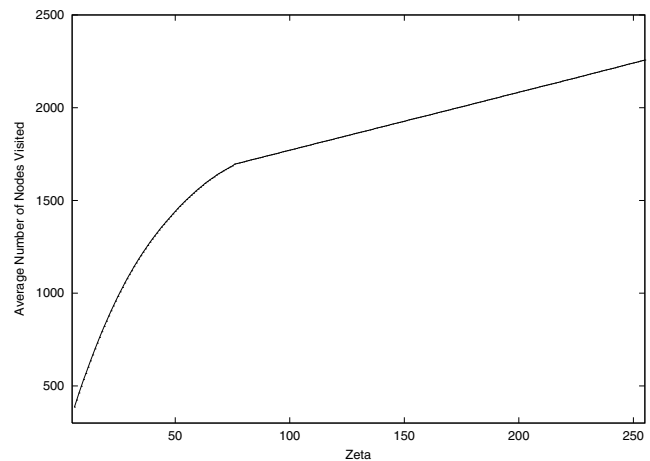
Table 1. Comparison of performance of direct k -NN search and hierarchical k -NN search

Compared with the direct full search of the training set, when $\zeta = 251$, the proposed hierarchical search algorithm reduces the dissimilarity computations by 84.7% with an accuracy drop by only 0.03%; when $\zeta = 50$, the proposed method reduces the dissimilarity computations by 90.2% with an accuracy loss 0.08%.

The experiments also show the effectiveness of *the hyper level*. With the same testing set and $\zeta = 100$, the search algorithm stops at *the hyper level* for 79% patterns, and the accuracy gets to 99.96%.



(a)



(b)

Figure 4. Dynamic performance of the hierarchical search algorithm: (a) relationship between recognition rate and number of top choices selected at each level. (b) relationship between number of nodes visited and number of top choices selected at each level.

6 Conclusions

A fast hierarchical search algorithm is proposed to find k - NN s with a non-metric dissimilarity measure $D(X, Y)$, to which the existing fast k - NN search algorithms are not applicable. The new algorithm takes advantage of only dissimilarities between templates and template class labels (in training set), thus it is applicable to NN methods using any non-metric dissimilarity measure.

Experiments with handwritten digits show that the proposed hierarchical search saves the dissimilarity computations by more than 90% with an accuracy drop only 0.08% when the parameter ζ in the algorithm are properly set.

The future research will focus on exploring mathematical fundamentals behind the algorithm.

References

- [1] T. M. Cover and P. E. Hart, "Nearest neighbor pattern classification," *IEEE Trans. Inform. Theory*, vol. IT-13, pp. 21–27, Jan 1968.
- [2] P. E. Hart, "Condensed nearest neighbor rule," *IEEE Trans. Inform. Theory*, vol. IT-14, pp. 515–516, May 1968.
- [3] A. J. Broder, "Strategies for efficient incremental nearest neighbor search," *Pattern Recognition*, vol. 23, no. 1/2, pp. 171–178, Nov 1986.
- [4] J. H. Friedman, F. Baskett, and L. J. Shustek, "An algorithm for finding nearest neighbors," *IEEE Trans. Computers*, vol. C-24, no. 10, pp. 1000–1006, Oct 1975.
- [5] K. Funkunaga and P.M. Navendra, "A branch and bound algorithm for computing k-nearest neighbors," *IEEE Trans. Computers*, vol. C-24, no. 7, pp. 750–753, July 1975.
- [6] J. T. Favata and G. Srikantan, "A multiple feature/resolution approach to handprinted character/digit recognition," *International Journal of Imaging Systems and Technology*, vol. 7, pp. 304–311, 1996.
- [7] J.D. Tubbs, "A note on binary template matching," *Pattern Recognition*, vol. 22, no. 4, pp. 359–365, 1989.
- [8] S. Cha and S. N. Srihari, "A fast nearest neighbor search algorithm by filtration," *Pattern Recognition*, vol. 35, no. 2, pp. 515–525, Feb 2002.