

# Adaptive Nearest Neighbor Classifier Based on Supervised Ellipsoid Clustering

Guo-Jun Zhang<sup>1,2</sup>, Ji-Xiang Du<sup>1,2</sup>, De-Shuang Huang<sup>1,2</sup>, Tat-Ming Lok<sup>3</sup>,  
and Michael R. Lyu<sup>4</sup>

<sup>1</sup> Intelligent Computing Lab, Hefei Institute of Intelligent Machines, Chinese Academy of Sciences, P.O. Box 1130, HeFei Anhui 230031, China  
{zhanggj, dshuang}@iim.ac.cn

<sup>2</sup> Department of Automation, University of Science and Technology of China, Hefei, China

<sup>3</sup> Information Engineering Dept., The Chinese University of Hong Kong, Hong Kong

<sup>4</sup> Computer Science & Engineering Dept., The Chinese University of Hong Kong, Hong Kong

**Abstract.** Nearest neighbor classifier is a widely-used effective method for multi-class problems. However, it suffers from the problem of the curse of dimensionality in high dimensional space. To solve this problem, many adaptive nearest neighbor classifiers were proposed. In this paper, a locally adaptive nearest neighbor classification method based on supervised learning style which works well for the multi-classification problems is proposed. In this method, the ellipsoid clustering learning is applied to estimate an effective metric. This metric is then used in the  $K$ -NN classification. Finally, the experimental results show that it is an efficient and robust approach for multi-classification.

## 1 Introduction

One of most popular classification approaches is the Nearest Neighbor (NN) method. However, the NN rule becomes less appealing in the case of limited training samples and high dimensional feature space due to the curse of dimensionality. Severe bias can be caused in such situations. Recently, several methods [1-4] characterized by query-based local distance functions were proposed to reduce this bias. However, the computation seems not to be efficient and there are too many parameters despite of the relatively high accuracy they achieve. Are there any better alternatives which could estimate the local relevance efficiently?

Our method uses the supervised ellipsoid clustering (SEC) to produce boundary and utilizes its boundary to estimate the local features relevance of the query with a scheme provided in LAMANN[1]. The relevance is then used in the weighted Euclidean distance during the  $K$ -NN classifications. The algorithm is referred to as Ellipsoids-boundary weighting adaptive Nearest Neighbor Algorithm (EWANN Algorithm) because the class boundary is constructed according to the surface of ellipsoids.

This paper is organized as follows. Section 2 presents the feature relevance theory. Some empirical evaluation of our method is given in Section 3. Finally, a concluding remark is included in Section 4 .

## 2 Ellipsoids-Boundary Weighting Adaptive Nearest Neighbor Algorithm

It is often noted that different classes may have different discriminant features. This means that distance computation does not always have equal strength on features in the feature space. Considering this, many locally adaptive nearest neighbor classifiers are designed to find these relevant dimensions.

Supervised ellipsoid clustering (SEC) is to cover each class region with ellipsoids. Each ellipsoid will represent a set of points [5]. Ellipsoids can be generated by an incremental learning procedure. The ellipsoids are created, constricted, or enlarged gradually at the present of each training sample. After SEC, each class is stuffed with a set of ellipsoids. In other views, the outer surface of the ellipsoids of a certain class composes the boundary of that class, i.e. the boundary of the class against the other classes is formed by the surface of its ellipsoids.

The gradient vector of the point on the boundary identifies a direction along which nearby data points are well separated. We use this gradient vector to measure local feature relevance and weighting features accordingly.

The feature relevance  $R(x)$  can be given as  $R_i(x) = |N_d \cdot u_i| = |N_{di}|$ , where  $d$  is the nearest point on the boundary to the query,  $N_d$  the gradient vector at point  $d$ ,  $u_i$  the vector unit. After  $r$  is transformed from  $R(x)$  by a scheme given by Jing Peng [1], it could be applied in the weighted distance computation during the NN classification. The resulting algorithm is summarized in Fig. 1.

<p><b>INPUT:</b> Ellipsoids of all classes produced by SEC, query <math>q</math> and parameter <math>K</math> for <math>K</math>-NN</p> <p><b>Step 1.</b> Find ellipsoid <math>E_n</math> which is nearest to query <math>q</math></p> <p><b>Step 2.</b> Find the point <math>d</math> which is the nearest point to <math>q</math> on the <math>E_n</math> .</p> <p><b>Step 3.</b> Compute the gradient vector <math>N_d</math></p> <p><b>Step 4.</b> Compute <math>R(x)</math> and transform it to <math>r</math> by the scheme given by Jing Peng</p> <p><b>Step 5.</b> Use <math>r</math> in weighted distance computation and apply <math>K</math>-NN rule</p>
---

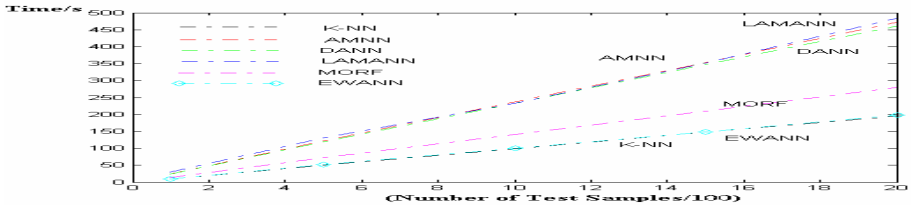
**Fig. 1.** Ellipsoids-boundary Weighting Adaptive Nearest Neighbor Algorithm

### 3 Empirical Evaluation

In the following, we compare several competing classification methods using a number of data sets. The classifiers are 5NN which uses the simple five NN rule, DANN using the discriminant adaptive NN rule.[2], Morph which uses he morphing rule [3],ADAMENN using the adaptive Metric Nearest Neighbor Algorithm [4] and LAMANN using the large margin nearest neighbor.[1]. The data sets were taken from the UCI Machine Learning Database Repository. We randomly select 60% of samples of each class as training samples and other 40% for testing.

**Table 1.** Average classification error rate

	Iris	Heart	Diabetes	Cancer
EWANN	<b>5.8</b>	24.4	<b>24.3</b>	23.6
5NN	<b>5.8</b>	24.8	25.4	24.1
DANN	<b>5.8</b>	23.7	24.8	22.7
Morph	<b>5.8</b>	<b>22.7</b>	25.7	<b>22.8</b>
ADAMENN	<b>5.8</b>	22.9	25.0	25.0
LAMNN	<b>5.8</b>	24.0	24.8	23.1



**Fig. 2.** On-line computing time performance

From Table 1, it can be found that EWANN achieved the best performance of the 2/4, followed closely by LAMANN. The result shows that our method is as competing as other adaptive classifiers. What’s more, as it is shown in Fig. 2, the proposed method is superior to other adaptive nearest neighbor methods in terms of online computing.

### 4 Conclusions

This paper presents a new flexible metric method for effective and efficient pattern classification. It employs SEC to generate class boundary, and use its information for the following nearest neighbor classification. The experimental results show clearly that the proposed algorithm can potentially improve the performance of K-NN in several classification problems.

## References

1. Domeniconi, C., Gunopulos, D., and Peng, J.: Large Margin Nearest Neighbor Classifiers, IEEE Transaction on Neural Networks, VOL 16, No. 4, (2005), 899-909
2. Hastie, T., Tibshirani, R.: "Discriminant Adaptive Nearest Neighbor Classification", IEEE Trans. on PAMI, Vol. 18, No. 6, (1996), 607-615
3. Peng, J., Heisterkamp, D. R., Dai, H.K.: LDA/SVM Driven Nearest Neighbor Classifiers, IEEE Transaction on Neural Networks, VOL 14, No. 4, (2003), 940-942
4. Domeniconi, C., Peng, J., Gunopulos, D.: Locally Adaptive Metric Nearest Neighbor Classification, IEEE Transaction on PAMI, (2002), 1281-1285
5. Kositsky, M., Ullman, S.: Learning Class Regions by the Union of Ellipsoids, Proceedings of the 13th ICPR, IEEE Computer Society Press, (1996), 750-757