Letters

# Large-margin nearest neighbor classifiers via sample weight learning

Qinghua Hu *, Pengfei Zhu, Yongbin Yang, Daren Yu

*Harbin Institute of Technology, PO 458, Harbin 150001, PR China*

## ARTICLE INFO

## ABSTRACT

The nearest neighbor classification is a simple and yet effective technique for pattern recognition. Performance of this technique depends significantly on the distance function used to compute similarity between examples. Some techniques were developed to learn weights of features for changing the distance structure of samples in nearest neighbor classification. In this paper, we propose an approach to learning sample weights for enlarging margin by using a gradient descent algorithm to minimize margin based classification loss. Experimental analysis shows that the distances trained in this way reduce the loss of the margin and enlarge the hypothesis margin on several datasets. Moreover, the proposed approach consistently outperforms nearest neighbor classification and some other state-of-the-art methods.

© 2010 Elsevier B.V. All rights reserved.

## 1. Introduction

The nearest neighbor (1-NN) rule, first proposed by Fix and Hodges, is one of the oldest and simplest pattern classification algorithms [1]. The performance of 1-NN classification depends crucially on the way that distances are computed. Distance functions are used to measure the similarity between two patterns for identifying the nearest neighbor of a query pattern.

In recent years, there has been growing interest in finding variants of the 1-NN rule via learning distance functions from samples [2–4] or prototype editing [5,6]. Goldberger et al. proposed a method via a stochastic variant of the leave-one-out 1-NN score on the training data to learn a Mahalanobis distance measure [10]. Weinberger et al. [12] introduced a method for learning a Mahalanobis distance metric from training samples by semidefinite programming. In addition, Roberto et al. developed a prototype weighting algorithm derived by approximately minimizing the leaving-one-out classification error of the given training set [20]. Wang et al. developed an extremely simple adaptive distance measure that assigns different weights to each sample [11]. Experimental analysis shows that these metrics learned from training samples perform well on some real-world tasks.

In the last few years, the large margin principle has been used to design classification algorithms [7,12]. We can show that 1-NN generates the classification rule with the largest margin once the proper definition of margin is used [8]. Motivated by this issue, several algorithms are developed to estimate the weights of features. The weights are then used in feature selection and distance learning [8,15,21]. The classification margin in the derived feature space is large.

In fact, margin based loss functions, instead of margin itself, are used in characterizing the complexity of classification [17,21]. Classification loss is a function of margin, such as hinge loss used in SVM [9,22]. Breiman showed that AdaBoost minimizes the exponential loss function of margin [13]. Mason et al. proposed AnyBoost classification procedures that perform gradient descent in the function space with general loss functions of the margin [14]. Weinberger et al. proposed LMNN for $k$ nearest neighbor classification from labeled examples using hinge loss [16].

In this paper, we introduce a technique to enlarge classification margin by learning weights of samples. The optimization objective is to minimize classification loss with a gradient decent algorithm. Experimental analysis shows that all the metrics trained with SWL-MLM greatly reduce the loss of the hypothesis margin and enlarge the hypothesis margin on several benchmark datasets, and the proposed approach consistently outperforms the nearest neighbor rule and A-NN [11].

The rest of the paper is organized as follows. Section 2 introduces some basic concepts about hypothesis margin and margin based loss functions. In Section 3, a sample weighted distance learning algorithm is proposed via minimizing margin based loss. Section 4 presents experimental results on artificial and real-world datasets. Finally, conclusions are given in Section 5.

## 2. Preliminaries

Some related definitions and theorems about margin, generalization and classification loss are reviewed in this section.

---

* Corresponding author.
   *E-mail address:* huqinghua@hit.edu.cn (Q. Hu).

## 2.1. Hypothesis margin

According to the statistic learning theory, the confidence of a classifier with respect to its predictions can be measured by class margin. There are two natural ways to define the margin of an instance with respect to a classification rule. One is to define margin as the distance between the instance and the decision boundary, such as support vector machines [23].

Besides, given an instance $x$, an alternative definition, hypothesis margin, with respect to a set of samples $S$, is defined as

$$\theta(x) = \tfrac{1}{2}(\|x - NM(x)\| - \|x - NH(x)) \tag{1}$$

where $NM(x)$ and $NH(x)$ are the nearest points from the same and different labels, called nearest miss (NM) and nearest hit (NH), respectively, and $\|x - NM(x)\|$ and $\|x - NH(x)\|$ mean the distances between $x$ and $NM(x)$ and $NH(x)$ [8,23]. Here, the size of the margin depends on the distance function. Based on the above definition, Crammer et al. derived a training algorithm that selects a good set of prototypes using large margin principles [7]. Moreover, this type of margin has also been used in AdaBoost [18] with the L1-norm.

Throughout this paper, we mainly discuss the margin for the nearest neighbor (1-NN) classification. As to margins for 1-NN, Crammer et al. proved that the hypothesis margin lower bounds the sample margin and it is easy to compute the hypothesis margin of an instance with respect to a set of samples [7].

Generally, a weighted distance can be written as

$$\Delta(x_1, x_2) = \sqrt{\sum_{i=1}^{N} w_i^2 |f(x_1, a_i) - f(x_2, a_i)|^2} \tag{2}$$

where $w_i$ is the weight of feature $a_i$. If we learn weights for different features from labeled training data, we call this technique feature weight learning [15,22].

On the other hand, one can also assign different weights to samples for changing the distance structure of samples. In this case, the sample weighted distance from $x_2$ to $x_1$ is computed with $w(x_1)\|x_2 - x_1\|$, where $w(x_1)$ is the weight of sample $x_1$. Correspondingly, the sample weighted distance from $x_1$ to $x_2$ is $w(x_2)\|x_1 - x_2\|$. Generally speaking, $w(x_1)\|x_2 - x_1\| \neq w(x_2)\|x_1 - x_2\|$ for $w(x_1) \neq w(x_2)$. Thus the sample weighted distance does not satisfy the conditions of a general distance function.

Given a sample weighted distance function, the sample weighted hypothesis margin of an instance $x$ with respect to a set of samples $S$ is

$$\theta^w(x) = \tfrac{1}{2}[w(NM(x))\|x - NM(x)\| - w(NH(x))\|x - NH(x)\|] \tag{3}$$

This hypothesis margin reflects the distance that instance would be misclassified if it walks close to the samples with different decisions.

## 2.2. Margin based generalization bound and classification loss

It has been shown [1] that the nearest neighbor rule has asymptotic error rate that is at most twice the Bayes error rate, independent of the distance metric used. However the training error is thus too rough to provide information on the generalization performance of 1-NN. A more detailed measure is needed so as to provide meaningful generalization bounds. Gilad-Bachrach gave us the generalization bound for 1-NN [8] with respect to the classification margin.

**Definition 1.** Let $D$ be a distribution over $X\{\pm 1\}$ and $h : \xi \rightarrow X\{\pm 1\}$ is a classification function. We denote $error_D(h)$ by the generalization error of $h$ with respect to $D$:

$$error_D(h) = \Pr_{x,y \sim D}(h(x) \neq y) \tag{4}$$

For a sample $S = \{(x_k, y_k)\}_{k=1}^{m} \in (\xi \times \{\pm 1\})^m$ and a constant $\gamma > 0$, we denote the $\gamma$-sensitive training error to be

$$\hat{error}_s(h) = \frac{1}{m}|\{(k : h(x_k \neq y_k)) \text{ or } x_k \text{ has sample margin} \prec \gamma\}| \tag{5}$$

$\gamma$ is a given constant. $\hat{error}_s(h)$ denotes the average number of the samples that are misclassified or have sample margin less than $\gamma$.

**Theorem 1.** *Let $D$ be a distribution over $R^N \in \{\pm 1\}$ which is supported on a ball of radius $R$ in $R^N$. Let $\delta > 0$ and $S$ be a sample of size $m$. With probability $1 - \delta$ over the random choice of $S$, for any $\gamma \in (0,1]$:*

$$error(h) \leq \hat{error}_s^r(h) + \sqrt{\frac{2}{m}\left[d\ln\left(\frac{34em}{d}\right)\log_2(578m) + \ln\left(\frac{8}{\gamma\delta}\right) + (N+1)\ln N\right]}$$

*where $h$ is the nearest neighbor classification rule and $d = (64R/r)^N$.*

From Theorem 1 [8], we can conclude that margin has great impact on the generalization of 1-NN. Large margin can reduce the generalization error of this rule. We can improve the performance of the nearest neighbor rule by enlarging the margin.

In practice, loss functions are used in machine learning for finding the optimal classification functions [13]. In these techniques, a margin based loss is associated for each hypothesis with respect to a sample. Hence, given a set of samples $S = \{x_1, \ldots, x_n\}$, we can define the hypothesis margin based loss function as

$$L(S) = \frac{1}{n}\sum_{x \in S} l(\theta(x)) \tag{6}$$

where $\theta(x)$ is the hypothesis margin of the sample $x$, $l(\theta(x))$ is the loss of the sample $x$ and $L(S)$ is the loss of the whole set $S$.

Some loss functions were proposed for different applications [23]. For instance, AdaBoost uses the exponential loss function, while SVM uses the hinge loss. In this paper, four loss functions are used to learn the sample weights, including the linear loss function $l(\theta) = 1 - \theta$, the logistic loss function $l(\theta) = \log(1 + \exp(-\theta))$, the exponential loss function $l(\theta) = \exp(-\theta^2/u)$ and the surrogate loss function $l(\theta) = \exp(-\theta) - \theta$. Usually users employ non-negative loss functions. Some of the loss functions we use here would produce negative values to guarantee that objective functions are consistent.

## 3. Distance learning for minimizing loss of hypothesis margin

In the above section, four margin based loss functions are introduced. We can minimize the loss of hypothesis margin by sample weighted metric learning. The sample weighted loss functions are defined as follows:

$$L^w(S) = \frac{1}{n}\sum_{x \in S} l(\theta^w(x)) = \frac{1}{n}\sum_{x \in S} l\left\{\frac{1}{2}[w(NM(x))\|x - NM(x)\| - w(NH(x))\|x - NH(x)\|]\right\} \tag{7}$$

We obtain the four loss functions after sample is weighted. The four sample weighted loss functions are

$$L_1^w(S) = \frac{1}{n}\sum_{x \in S} 1 - \frac{1}{2}[w(NM(x))\|x - NM(x)\| - w(NH(x))\|x - NH(x)\|]$$

$$L_2^w(S) = \frac{1}{n}\sum_{x \in S} \exp\left\{-\frac{1}{2}[w(NM(x))\|x - NM(x)\| - w(NH(x))\|x - NH(x)\|]\right\}$$

$$L_3^w(S) = \frac{1}{n}\sum_{x \in S} \log\left\{1 + \exp\left(-\frac{1}{2}[w(NM(x))\|x - NM(x)\| - w(NH(x))\|x - NH(x)\|]\right)\right\}$$

$$L_4^w(S) = \frac{1}{n}\sum_{x \in S} \exp\left\{-\frac{1}{2}[w(NM(x))\|x - NM(x)\|, -w(NH(x))\|x - NH(x)\|]\right\}$$
$$\quad - \frac{1}{2}[w(NM(x))\|x - NM(x)\| - w(NH(x))\|x - NH(x)\|]$$

In the previous section, we have concluded that we can reduce the generalization error of the 1-NN via minimizing loss of hypothesis Margin. Now, we minimize loss of hypothesis margin via sample weighted distance learning.

The weights of samples are usually set as the same value. In this case, we consider that $w_j=1$. However, it is known that the importance of samples is different. Distances are measured differently depending on the location of the samples. Thus we develop an algorithm to optimize the weight vector $W = \langle w_1, w_2, \ldots, w_j, \ldots, w_n \rangle$ of samples. The optimization objective function is the loss function of hypothesis margin.

In fact, we should note that the margin based loss $L^w(S)$ not only depends on the weights $w$, but also is up to the $NM(x)$ and $NH(x)$ in that for each $x \in S$, the $NM(x)$ and $NH(x)$ may change if the weights are varied. Accordingly, we can define the loss function as follows:

$$L^w(S) = f(w, T(w)) \tag{8}$$

where $T(w)$ is a function that can find the nearest miss and nearest hit of $x$. Now we compute the partial derivatives of $L^w(S)$ as

$$\frac{\partial L^w(S)}{\partial w} = \frac{\partial f}{\partial w} + \frac{\partial f}{\partial T(w)} \frac{\partial T(w)}{\partial w} \tag{9}$$

It is easy to get $\partial f/\partial w$. However, the second term $T(w)$ is not a continuous function of $w$. While this formulation can still be followed to some extent, the development becomes rather cumbersome and it does not really lead to useful approximations. Therefore, we assume that the samples neighborhoods remain unchanged for sufficiently small variations of the weights, like the technique in [20].

Based on the above assumption, we know $L^w(S)$ is smooth, we can use gradient descent to minimize it. The minimization of $L^w(S)$ by gradient descent consists in an iterative procedure which updates the weights $w(i)$ by a small amount, in the negative direction of the gradient of $L^w(S)$:

$$w(i) = w(i) - \eta \frac{\partial L^w(S)}{\partial w(i)} \tag{10}$$

With regard to the linear loss function, the update equations are given as

$$w(i) = w(i) - \left\{ \eta_1 \sum_{\substack{\forall x \in S \\ index(NH(x))=i}} \Delta(w(NH(x))) - \eta_2 \sum_{\substack{\forall x \in S \\ index(NM(x))=i}} \Delta(w(NM(x))) \right\}$$

where $\partial L^w(S)/\partial w(NH(x)) = \sum_{\forall x \in S} \Delta(w(NH(x)))$, $\partial L^w(S)/\partial w(NM(x)) = \sum_{\forall x \in S} \Delta(w(NM(x)))$.

Given a set of training samples, we can iteratively search the weight with the following procedure.

**Algorithm 1.** Sample weight learning via minimizing loss of margin (SWL-MLM)

1: **procedure** INITIALIZE($w = <1,1,\ldots,1>$, loss=0, loss1=1, $\varepsilon > 0.001$)
2: $\forall x \in U$, compute $NM(x)$ and $NH(x)$, $\mu = \frac{1}{m} \sum_{x \in U} (\|x - NM(x)\| - \|x - NH(x)\|)$
3: **while** $|loss1 - loss| > \varepsilon$ **do**
4: loss1 = loss
5: **for** i = 1, 2 n **do**
6: $w(NH(x)) = w(NH(x)) - \eta_1 \Delta(w(NH(x)))$
7: $w(NM(x)) = w(NM(x)) + \eta_2 \Delta(w(NM(x)))$
8: **end for**
9: Compute the loss after samples are weighted, $loss = L^w(S)$
10: **end while**
11: **end procedure**

The values of $\eta$ are referred to as learning rates or learning step factors. It can take just a fixed value for all samples or may vary on different samples. For instance, they may be inversely proportional to the variance of each sample. Two different learning rates $\eta_1$ and $\eta_2$ are set for the gradient $\Delta(w(NM(x)))$ and $\Delta(w(NH(x)))$, because we cannot guarantee the distance between the nearest hit of the sample $x$ and the class center, that is, the nearest hit of the sample $x$ is not necessarily closer to the class center and it may be anywhere in the sample space. However, the nearest miss of is on or close to the classification boundary to a certainty. Hence, we set $\eta_1$ as zero and give $\eta_2$ an appropriate positive value.

In essence, SWL-MLM mainly assigns greater weights to the boundary samples whose hypothesis margin is very small or negative and does not change weights of non-boundary samples. In order to simplify the learning process, the learning rate $\eta_1$ for the gradient $\Delta(w(NH(x)))$ can be set as zero. Thus, the weights of non-boundary samples keep unchanged. The closer a sample point is to the classification boundary, the positive gradients $\Delta(w(NM(x)))$ of more heterogeneous samples are added to the weight of this sample.

The computational complexity of SWL-MLM is $O(Tm)$, where $T$ is the number of iterations and $m$ is the size of the samples $S$.

## 4. Experimental analysis

In this section, we test the sample-weighted technique with some real-world datasets. We gathered ten datasets from UCI machine learning repository [19]. The difference of hypothesis margin before and after samples are weighted is given in Table 1. We can see that the average hypothesis margin of the datasets greatly increases after samples are weighted.

We also compute the hypothesis margin based loss before and after samples are weighted. From Table 2, we can see that the loss has been reduced on all the datasets regardless of loss functions. In fact, we can easily get the variation trend of loss through the variance of the hypothesis margin because the hypothesis margin is in inverse proportion to the margin based loss. From Tables 2 and 3 we can find that different loss functions give different margins and loss and even if the loss function gives the largest increase of margin, it is not reflected as the minimum decrease of loss using that function for the given dataset. Actually, as to each loss function, their gradient is different. Therefore, the weights of the samples we learned are different with respect to each loss function. That is why different loss functions give different margins and loss. In addition, hypothesis margin or loss is the average margin or loss of all the samples. As to different loss functions, the margin of each sample is different after samples are weighted. Hence, even if the loss function gives the largest increase of margin, we cannot guarantee that it gives the minimum decrease of loss.

**Table 1**
Hypothesis margin before and after sample-weighted.

| Data | Before weighted | After weighted | | | |
|---|---|---|---|---|---|
| | | L-SWL | LL-SWL | EL-SWL | SL-SWL |
| Heart | 0.4478 | 2.0365 | 2.7339 | 1.3832 | 2.3494 |
| Hepatitis | 0.5326 | 2.4738 | 2.6575 | 1.7667 | 2.6117 |
| Horse | 1.2092 | 12.9244 | 12.8163 | 8.4353 | 13.2502 |
| Iono | 0.4248 | 0.9491 | 1.1195 | 4.1819 | 1.0204 |
| WDBC | 0.2594 | 0.4389 | 0.4748 | 0.4345 | 0.4649 |
| WPBC | 0.0599 | 0.2164 | 0.4661 | 0.2021 | 0.3200 |
| Wine | 0.2744 | 0.4046 | 0.4356 | 0.3991 | 0.4261 |
| German | 0.3131 | 5.3582 | 5.4221 | 5.7092 | 0.8086 |
| Crx | 4.3239 | 9.8169 | 8.1156 | 5.6307 | 7.2531 |
| Derm | 2.3138 | 2.6544 | 2.6778 | 2.6482 | 2.6802 |

**Table 2**
Comparison of margin based loss before and after sample-weighted.

| Data | L-SWL | | L-SWL | | LL-SWL | | SL-SWL | |
|---|---|---|---|---|---|---|---|---|
| | Before | After | Before | After | Before | After | Before | After |
| Heart | 0.5522 | −1.0365 | 0.5466 | 0.5316 | 0.4724 | 0.2021 | −0.9120 | −2.7850 |
| Hepatitis | 0.4674 | −1.4738 | 0.5062 | 0.3219 | 0.4838 | 0.1207 | 0.1885 | −1.4173 |
| Horse | −0.2092 | −11.9244 | 0.3682 | 0.0070 | 0.3770 | 0.0309 | −3.0640 | −17.925 |
| Iono | 0.5752 | 0.0509 | 0.5241 | 0.3642 | 0.6352 | 0.2450 | 0.2852 | −0.4800 |
| WDBC | 0.7406 | 0.5611 | 0.5770 | 0.5083 | 0.7498 | 0.6098 | 0.5272 | 0.2185 |
| WPBC | 0.9401 | 0.7836 | 0.6658 | 0.5269 | 0.7617 | 0.3847 | 0.8904 | 0.4751 |
| Wine | 0.7256 | 0.5954 | 0.5683 | 0.5055 | 0.7307 | 0.5690 | 0.4948 | 0.2446 |
| German | 0.6869 | −4.3582 | 0.6137 | 0.0307 | 0.4359 | 0.3718 | 0.6612 | −0.0180 |
| Crx | −3.3239 | −8.8169 | 0.8491 | 0.8098 | 0.4253 | 0.1039 | −3.2846 | −6.2395 |
| Derm | −1.3138 | −1.6757 | 0.1916 | 0.1506 | 0.5233 | 0.4554 | −2.0700 | −2.5036 |

**Table 3**
Classification accuracy of SWL-MLM and 1-NN.

| Data | 1-NN | L-SWL | LL-SWL | EL-SWL | SL-SWL |
|---|---|---|---|---|---|
| Heart | 76.6 ± 9.4 | 79.3 ± 6.6 | 80.4 ± 5.8 | 79.3 ± 5.8 | 81.9 ± 5.5 |
| Hepatitis | 82.5 ± 7.6 | 85.7 ± 8.6 | 84.8 ± 8.2 | 85.8 ± 8.8 | 85.3 ± 8.0 |
| Horse | 87.2 ± 4.2 | 90.2 ± 2.6 | 90.5 ± 3.4 | 89.8 ± 3.4 | 90.0 ± 3.4 |
| Iono | 86.4 ± 4.9 | 89.6 ± 6.0 | 90.1 ± 5.5 | 90.4 ± 5.6 | 90.1 ± 4.9 |
| WDBC | 95.4 ± 3.3 | 97.2 ± 2.2 | 97.2 ± 2.2 | 97.4 ± 2.2 | 97.4 ± 2.2 |
| WPBC | 70.6 ± 6.8 | 77.2 ± 6.5 | 74.7 ± 6.3 | 75.7 ± 5.7 | 74.3 ± 5.4 |
| Wine | 94.8 ± 5.1 | 96.6 ± 2.9 | 96.6 ± 2.9 | 96.6 ± 2.9 | 96.6 ± 2.9 |
| Crx | 79.1 ± 11.6 | 84.2 ± 16.8 | 81.4 ± 13.0 | 83.3 ± 15.9 | 82.9 ± 12.9 |
| Derm | 96.1 ± 5.7 | 97.9 ± 3.0 | 98.2 ± 2.3 | 97.9 ± 2.9 | 97.9 ± 3.0 |
| IRIS | 96.0 ± 5.6 | 97.3 ± 4.7 | 97.3 ± 4.7 | 97.3 ± 4.7 | 97.3 ± 4.7 |

**Table 4**
Classification accuracy of A-1-NN Relief LSVM and LVQ.

| Data | A-1-NN | Relief | LSVM | LVQ |
|---|---|---|---|---|
| Heart | 77.0 ± 5.5 | 82.9 ± 6.8 | 83.3 ± 5.3 | 78.9 ± 10.2 |
| Hepatitis | 83.2 ± 9.4 | 79.7 ± 9.7 | 86.2 ± 7.7 | 83.8 ± 10.3 |
| Horse | 88.1 ± 3.1 | 90.7 ± 4.3 | 93.0 ± 4.4 | 88.3 ± 6.3 |
| Iono | 90.1 ± 4.1 | 77.6 ± 5.7 | 87.6 ± 6.5 | 86.9 ± 5.3 |
| WDBC | 96.5 ± 2.5 | 94.7 ± 2.6 | 97.7 ± 2.5 | 96.0 ± 1.9 |
| WPBC | 72.7 ± 10.3 | 77.2 ± 8.0 | 7.4 ± 7.7 | 69.7 ± 7.3 |
| Wine | 96.6 ± 2.9 | 97.2 ± 3.0 | 98.9 ± 2.4 | 94.9 ± 6.3 |
| German | 71.3 ± 2.8 | 69.3 ± 4.0 | 73.7 ± 4.7 | 70.9 ± 4.9 |
| Crx | 80.7 ± 11.7 | 79.8 ± 15.1 | 85.5 ± 18.5 | 84.1 ± 3.5 |
| Derm | 96.3 ± 0.5 | 96.5 ± 2.4 | 96.5 ± 2.8 | 97.3 ± 3.2 |
| IRIS | 97.3 ± 4.7 | 94.7 ± 4.2 | 97.3 ± 4.7 | 95.7 ± 3.7 |

In order to show the influence of the new technique on the classification, we compare it to 1-NN, A-NN [11], Relief, LSVM and LVQ. In Table 4, the experiment results indicate that 1-NN using the sample-weighted technique is much better than the original NN classifier and A-NN. Regardless of loss functions, the improvements of the sample-weighed method by minimizing the hypothesis margin based loss functions are statistically significant. These results confirm that the first nearest neighbor identified according to the sample-weighted technique is more likely to have the same class label as the query pattern than the first nearest neighbor identified with all the sample in the dataset treated as equally important.

It is shown that Relief is based on the large-margin principle. We can calculate weights for features by Relief, and the weights can also be used for feature weighting. From Table 4, we can see that SWL-MLM outperforms Relief when the weights learned by Relief is used to weight features. Besides, the proposed method is very similar to SVM in the sense that SVM is a linear combination of

some prototypes (support vectors) and in the paper the sample weight can be seen as the weight that the SVM set to each prototype as well. In this sense, we compare the proposed method to LSVM (one of the best classifier in pattern recognition). In Table 4, we can see that although SWL-MLM is better than LSVM on some datasets, the average classification accuracy of LSVM is higher than the proposed one. Besides, the average classification accuracy of LVQ which maximizes the margin by prototype learning, is lower than SWL-MLM.

It is worth noting that no significant difference is found among different loss functions based algorithms. As to some tasks, line loss gets the best performance, however, sometimes, exponential loss produces the highest classification accuracy. As a whole, no loss function consistently outperforms other functions. The difference in performance is not significant enough for discerning. In this sense, we can select any of them in real-world applications or tries some of them to find the best solution.

## 5. Conclusions

The performance of NN classification depends crucially on the way that distances are computed. In this paper, with margin based loss functions, we obtain sample weighted metrics for samples. We test the sample-weighted technique SWL-MLM with some real-world datasets and find that the classification loss is greatly reduced. Besides, we show how the sample-weighted technique impacts on the classification boundary and the hypothesis margin. In these experiments, the proposed approach outperforms nearest neighbor classification and several state-of-the-art methods in terms of classification accuracy.
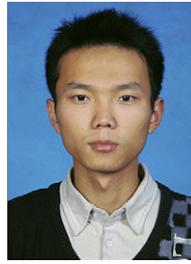
## References

[1] P. Hart, T. Cover, Nearest neighbor pattern classification, IEEE Transactions on Information Theory 13 (1967) 21–27.
[2] C. Stanfill, D. Waltz, Toward memory-based reasoning, Communications of the ACM 29 (1986) 1213–1228.
[3] I. Kononenko, Estimating attributes: analysis and extensions of RELIEF, Technical Report, Faculty of Electrical Engineering and Computer Science, University of Ljubljana, 1993.
[4] R. Kohavi, P. Langley, Y. Yung, The utility of feature weighting in nearest-neighbor algorithms, in: Proceedings of the Ninth European Conference on Machine Learning, 1997.

[5] R. Paredes, E. Vidal, D. Keysers, An evaluation of the WPE algorithm using tangent distance, in: Proceedings of the International Conference on Pattern Recognition, 2002, pp. 48–51.

[6] R. Paredes, E. Vidal, Weighting prototypes a new editing approach, Proceedings of the 15th International Conference on Pattern Recognition, vol. 2, 2000, pp. 25–28.

[7] K. Crammer, R. Gilad-Bachrach, A. Navot, A. Tishby, Margin analysis of the LVQ algorithm, Advances in Neural Information Processing Systems 15 (2003) 462–469.

[8] R. Gilad-Bachrach, A. Navot, N. Tishby, Margin based feature selection—theory and algorithms, in: ICML, 2004.

[9] M. Brown, W. Grundy, D. Lin, N. Cristianini, C. Sugnet, T. Furey, M. Ares, D. Haussle, Knowledge-based analysis of microarray gene expressions data by using support vector machines, Proceedings of the National Academy of Sciences 97 (1) (2000) 262–267.

[10] J. Goldberger, S. Roweis, G. Hinton, R. Salakhutdinov, Neighborhood component analysis, Neural Information Processing Systems (NIPS) 17 (2004) 513–520.

[11] J. Wang, P. Neskovic, L.N. Cooper, Improving nearest neighbor rule with a simple adaptive distance measure, Pattern Recognition Letters 28 (2007) 207–213.

[12] K. Weinberger, J. Blitzer, L. Saul, Distance metric learning for large margin nearest neighbor classification, Advances in Neural Information Processing Systems (NIPS) 18 (2006) 1473–1480.

[13] L. Breiman, Prediction games and arcing algorithms, Neural Computation 11 (1999) 1493–1517.

[14] L. Mason, J. Baxter, P. Bartlett, M.R. Frean, Boosting algorithms as gradient descent in function space, Neural Information Processing Systems 12 (2000) 512–518.

[15] C. Domeniconi, D. Gunopulos, J. Peng, Large margin nearest neighbor classifiers, IEEE Transactions on Neural Networks 16 (2005) 899–909.

[16] K.Q. Weinberger, L.K. Saul, Distance metric learning for large margin nearest neighbor classification, Journal of Machine Learning Research 10 (2009) 207–244.

[17] V. Vapnik, The Nature of Statistical Learning Theory, Springer-Verlag, 1995.

[18] Y. Freund, R.E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, Journal of Computer and System Sciences 55 (1997) 119–139.

[19] A. Asuncion, D. Newman, UCI machine learning repository, University of California, Irvine, School of Information and Computer Sciences, URL ⟨http://www.ics.uci.edu/mlearn/MLRepository.html⟩, 2007.

[20] R. Paredes, E. Vidal, Learning weighted metrics to minimize nearest-neighbor classification error, IEEE Transactions on Pattern Analysis and Machine Intelligence 28 (2006) 1100.

[21] X. Nguyen, M.J. Wainwright, M.I. Jordan, Divergences, surrogate loss functions and experimental design, in: Proceedings of NIPS 2005.

[22] Y.J. Sun, Iterative RELIEF for feature weighting: algorithms, theories, and applications, IEEE Transactions on Pattern Analysis and Machine Intelligence 29 (2007) 1035–1051.

[23] N. Cristianini, J. Shawe-Taylor, An Introduction to Support Vector Methods, Cambridge University Press, Cambridge, 2000.
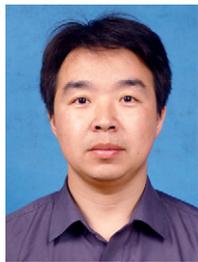
**Pengfei Zhu** is a master student with Harbin Institute of Technology. His research interests are focused on large-margin learning theory.



**Yongbin Yang** received his Bachelor and Master degrees from Department of Control Science and Engineering, Harbin Institute of Technology in 1982 and 1984, respectively. He has been with School of Energy Science and Engineering, Harbin Institute of Technology since 1982. Now he is an associate professor with this school. His main interests are focused on modeling and simulation of complex power systems.



**Daren Yu** was born in Datong, China, in 1966. He received the M.Sc. and D.Sc. degrees from Harbin Institute of Technology, Harbin, in 1988 and 1996, respectively. Since 1988, he has been working at the School of Energy Science and Engineering, Harbin Institute of Technology. His main research interests are in modeling, simulation, and control of power systems. He has published more than 100 conference and journal papers on power control and fault diagnosis.



**Qinghua Hu** received the master degree in power engineering from Harbin Institute of Technology, Harbin, China in 2002, and got his Ph.D. from Harbin Institute of Technology in 2008. Now he is an associate professor with Harbin Institute of Technology and a postdoctoral fellow with the Hong Kong Polytechnic University. His research interests are focused on data mining, knowledge discovery with fuzzy and rough techniques. He is a PC co-chair of RSCTC 2010 and serves as a referee for a great number of journals and conferences. He has authored or coauthored more than 70 journal and conference papers in the areas of machine learning, data mining and rough set theory.